



FlowVAT: Normalizing Flow Variational Inference with Affine-Invariant Tempering

Juehang Qin 
Rice University
qinjuehang@rice.edu

Shixiao Liang 
Rice University
liangsx@rice.edu

Christopher Tunnell 
Rice University
tunnell@rice.edu

Abstract

Multi-modal and high-dimensional posteriors present significant challenges for variational inference, causing mode-seeking behavior and collapse despite the theoretical expressiveness of normalizing flows. Traditional annealing methods require temperature schedules and hyperparameter tuning, falling short of the goal of truly black-box variational inference. We introduce FlowVAT, a conditional tempering approach for normalizing flow variational inference that addresses these limitations. Our method tempers both the base and target distributions simultaneously, maintaining affine-invariance under tempering. By conditioning the normalizing flow on temperature, we leverage overparameterized neural networks' generalization capabilities to train a single flow representing the posterior across a range of temperatures. This preserves modes identified at higher temperatures when sampling from the variational posterior at $T = 1$, mitigating standard variational methods' mode-seeking behavior. In experiments with 2, 10, and 20 dimensional multi-modal distributions, FlowVAT outperforms traditional and adaptive annealing methods, finding more modes and achieving better ELBO values, particularly in higher dimensions where existing approaches fail. Our method requires minimal hyperparameter tuning and does not require an annealing schedule, advancing toward fully-automatic black-box variational inference for complicated posteriors.

1 Introduction

Statistical inference forms the cornerstone of empirical research, allowing us to draw meaningful conclusions about unobserved variables from observed data [1]. This process is frequently framed in terms of posterior inference, where we have an observed dataset, \mathbf{X} , and aim to obtain a posterior distribution, $p(\boldsymbol{\theta}|\mathbf{X})$, over parameters $\boldsymbol{\theta}$, by combining a likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ with a prior distribution $p(\boldsymbol{\theta})$. Since exact inference is often intractable, researchers have developed various approximate methods. Variational Inference (VI) has emerged as a particularly powerful approach [1, 2, 3], wherein the posterior is approximated by a parameterized family of distributions, allowing us to find an approximate distribution $q_\phi(\boldsymbol{\theta})$ that closely resembles the true posterior $p(\boldsymbol{\theta}|\mathbf{X})$.

A significant advancement in this field has been the development of black-box variational inference [2, 3, 4], which provides a general-purpose framework applicable to diverse inference problems without requiring tailored variational families or methodologies for each specific scenario. When combined with highly expressive distribution families such as normalizing flows, these methods can theoretically represent even complex multimodal posteriors [3, 5]. However, despite the theoretical expressiveness of these models, variational methods in practice often exhibit mode-seeking behavior [6, 7], presenting a significant obstacle to the development of fully-automated variational inference systems capable of handling high-dimensional posteriors with multiple modes. Such high-dimensional and multimodal posteriors are often encountered in science, including in neutrino oscillation analyses [8, 9], searches for physics beyond the Standard Model [10], and gravitational wave inference [11].

To address these challenges, researchers have explored tempering approaches that can improve inference for complicated posteriors. Various tempering strategies and schedules have been investigated for both variational methods and Markov Chain Monte Carlo (MCMC) sampling [12, 13, 14]. The fundamental idea behind tempering is to modify either the likelihood or the posterior by raising it to the power of $1/T$, creating a shallower distribution that is easier to fit. The temperature parameter T can then be decreased until we recover the original target distribution at $T = 1$.

In this work, we introduce a novel approach that eliminates the need for annealing schedules and associated hyperparameters, bringing us one step closer to fully-automatic black-box variational inference capable of handling complicated multimodal posteriors. Our key insight is that by simultaneously tempering both the base distribution and the target posterior, the transformation represented by the normalizing flow can remain unchanged even as temperature changes. This means that the transformation represented by the normalizing flow does not need to vary as much across temperatures, thus promoting learning transformations that are generalizable across temperatures. Furthermore, by conditioning the normalizing flow on the temperature parameter, we leverage the generalization capabilities of overparameterized neural networks [15, 16] to train a single flow that can accurately represent the posterior distribution across a wide spectrum of temperatures. Crucially, this approach preserves modes identified at higher temperatures even at $T = 1$, thereby addressing one of the fundamental limitations of traditional variational methods. We call our method **FlowVAT**: normalizing **Flow** Variational inference with **A**ffine-invariant **T**empering, reflecting its integration of normalizing flows with our novel affine-invariant tempering approach to enhance variational inference. Code used to generate all data and figures in this paper will be made available on github.

2 Affine-Invariant Tempering for Normalizing Flow Variational Inference

There are two key motivations for our approach. First, we aim to leverage the demonstrated capacity of overparameterized neural networks to generalize effectively [15, 16]. Second, we seek to introduce a beneficial inductive bias that enhances inference quality. This is accomplished by applying temperature scaling to both the base distribution and the posterior distribution of a normalizing flow.

A normalizing flow, represented by the transformation \mathbf{f}_ϕ , defines a probability distribution through:

$$q_\phi(\boldsymbol{\theta}) = q(\mathbf{z}) \left| \frac{\partial \mathbf{f}_\phi}{\partial \mathbf{z}} \right|^{-1}, \quad (1)$$

where ϕ denotes the weights and parameters of the normalizing flow, \mathbf{z} represents the space of the base distribution, and $\boldsymbol{\theta}$ corresponds to the parameter space for inference. In variational inference, we optimize the normalizing flow using loss function that is minimized when the approximating distribution $q_\phi(\boldsymbol{\theta})$ converges to the true posterior $p(\boldsymbol{\theta}|\mathbf{X})$, such as the negative evidence lower bound (ELBO) [17, 2]. Maximizing the ELBO is equivalent to minimizing the Kullback-Leibler (KL) divergence [18]. The ELBO can be expressed as:

$$\mathcal{L}(\phi) = \mathbb{E}_{q(\mathbf{z})} \left[\log(p(\mathbf{X}, \mathbf{f}_\phi(\mathbf{z}))) - \log \left(q(\mathbf{z}) \left| \frac{\partial \mathbf{f}_\phi}{\partial \mathbf{z}} \right|^{-1} \right) \right]. \quad (2)$$

2.1 Theoretical Framework and Heuristic Arguments

To illustrate why tempering both the base and transformed distributions improves training, we examine a simple case where the base distribution is a standard Gaussian, $\mathcal{N}(\mu = 0, \sigma = \mathbb{I})$, and the posterior is a multivariate Gaussian. In such instances, the optimal transformation for a normalizing flow would be an affine transformation, as this can transform a standard Gaussian distribution into an arbitrary multivariate Gaussian:

$$\boldsymbol{\theta} = \mathbf{f}(\mathbf{z}) = \mathbf{A}\mathbf{z} + \mathbf{b}. \quad (3)$$

The Jacobian of this transformation is the matrix \mathbf{A} , resulting in the log determinant of the Jacobian matrix, $\log |\partial \mathbf{f} / \partial \mathbf{z}|$ being a constant. From eq. (2), we observe that a constant value of $\log |\partial \mathbf{f} / \partial \mathbf{z}|$ does not affect the gradient of the ELBO.

When training with a tempered posterior, we raise the original posterior to a power of $1/T$ [12], yielding:

$$p^{1/T}(\boldsymbol{\theta}|\mathbf{X}) = \frac{p^{1/T}(\mathbf{X}, \boldsymbol{\theta})}{\int p^{1/T}(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta}}. \quad (4)$$

This modifies our ELBO target to:

$$\mathcal{L}_T(\phi) = \mathbb{E}_{q_T(\mathbf{z})} \left[\log \left(p^{1/T}(\boldsymbol{\theta}|\mathbf{X}) \right) + \log \left(\int p^{1/T}(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} \right) - \log \left(q_T(\mathbf{z}) \left| \frac{\partial \mathbf{f}_{\phi, T}}{\partial \mathbf{z}} \right|^{-1} \right) \right]. \quad (5)$$

Since constant terms such as the marginalization integral do not affect the gradient, we can simplify to:

$$\mathcal{L}'_T(\phi) = \mathbb{E}_{q_T(\mathbf{z})} \left[\log \left(p^{1/T}(\boldsymbol{\theta}|\mathbf{X}) \right) - \log \left(q_T(\mathbf{z}) \left| \frac{\partial \mathbf{f}_{\phi, T}}{\partial \mathbf{z}} \right|^{-1} \right) \right]. \quad (6)$$

By additionally tempering the base distribution as $q_T(\mathbf{z}) = q^{1/T}(\mathbf{z})$, we obtain:

$$\mathcal{L}'_T(\phi) = \mathbb{E}_{q^{1/T}(\mathbf{z})} \left[\log \left(p^{1/T}(\mathbf{f}_{\phi, T}(\mathbf{z})|\mathbf{X}) \right) - \log \left(q^{1/T}(\mathbf{z}) \right) - \frac{1}{T} \log \left| \frac{\partial \mathbf{f}_{\phi, T}}{\partial \mathbf{z}} \right|^{-1} \right]. \quad (7)$$

For affine transformations as defined in eq. (3), the ELBO remains invariant except for the sampling distribution and constant or multiplicative factors. Importantly, if there exists a transformation for which the transformed posterior is exactly matched by the approximate distribution q (indicating a Gaussian posterior when using a standard Gaussian base distribution), then the transformation $\mathbf{f}_{\phi}(\mathbf{z})$ that achieves $\log(p^{1/T}(\mathbf{f}_{\phi}(\mathbf{z})|\mathbf{X})) = \log(q^{1/T}(\mathbf{z}))$ does so regardless of temperature.

Furthermore, we can represent a non-affine transformation as the composition of an affine transformation and an additional non-affine component. While our tempering approach still requires the non-affine component to adapt with temperature, the affine transformation can remain invariant. Heuristically, this affine-invariant tempering promotes consistent matching of the first two moments (mean and variance) across temperatures, particularly when using overparameterized neural networks that exhibit strong generalization properties.

Matching the first two moments across temperatures provides an additional benefit. In practice, the ELBO is estimated using Monte Carlo (MC) samples; as such, regions of parameter space with low sampling probability may become effectively inaccessible. This can result in local minima in the loss function when separated modes are not adequately sampled. When means and variances are matched across temperatures, our dual tempering approach promotes a sampling distribution that maintains approximately the same mean position across temperatures in the parameter space, but with a wider distribution at higher temperatures. This significantly reduces the probability of becoming trapped in local minima caused by inadequate coverage of the parameter space by the sampling distribution.

2.2 Evidence calculation using importance sampling

Using temperature as a conditional variable also enhances estimation of Bayesian evidence via importance sampling. In variational inference, we typically do not have access to the normalized posterior $p(\boldsymbol{\theta}|\mathbf{X})$, but rather to an unnormalized function proportional to the posterior density, given by the product of the likelihood function and the prior:

$$p'(\boldsymbol{\theta}) = p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta}). \quad (8)$$

The Bayesian evidence (marginal likelihood) is then given by:

$$\mathcal{Z} = \int p'(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (9)$$

The Bayesian evidence is particularly useful for model comparison using Bayes factors [1]. However, the evidence integral can be high-dimensional and thus intractable. There are various approaches to estimating the evidence integral, such as Nested Sampling [19, 20, 21]. Another approach is to formulate this integral as an importance sampling problem [22], such that

$$\mathcal{Z} = \int q_{\phi}(\boldsymbol{\theta}) \frac{p'(\boldsymbol{\theta})}{q_{\phi}(\boldsymbol{\theta})} d\boldsymbol{\theta} = \mathbb{E}_{q_{\phi}(\boldsymbol{\theta})} \left[\frac{p'(\boldsymbol{\theta})}{q_{\phi}(\boldsymbol{\theta})} \right], \quad (10)$$

where $q_\phi(\theta)$ is a normalized sampling distribution, in our case represented by the variational distribution. In practice, this approach is challenging, as it requires access to a proposal distribution that is normalized and covers the posterior sufficiently well while also remaining sufficiently concentrated. This is because if there are any regions of parameter space which are poorly covered by our approximate distribution, $q_\phi(\theta)$, then the importance ratio can have a heavy right tail; this is the motivation behind strategies such as Pareto Smoothed Importance Sampling [23]. On the other hand, using a very wide sampling distribution, such as the prior, would avoid a heavy right tail but would be extremely inefficient if the posterior is much more concentrated than the prior [22].

As FlowVAT results in a trained normalizing flow model with a conditional temperature parameter, we can avoid both of these issues by using the FlowVAT approximate posterior as a proposal distribution with. We can stabilize the importance ratios in eq. (10) using $T > 1$, thus decreasing the ratio $p'(\theta)/q_\phi(\theta)$ in the tails of the distribution. We demonstrate this approach in section 4. For particularly challenging posteriors, it is possible to combine this approach of a tempered sampling distribution with Pareto Smoothed Importance Sampling to improve the reliability of evidence estimates, and the Pareto shape diagnostic can be used to determine whether estimates are reliable.

3 Related Work

One key approach to mitigate mode-seeking behaviour in variational inference has been the use of loss functions which promote mode-covering behavior, such as Rényi Divergence [24], which generalizes the divergence used in Importance Weighted Autoencoders [25], tail-adaptive f -divergence [26], and the scale-invariant Alpha-Beta divergence [27]. However, in high-dimensional problems, due to the instability of mode-covering divergences, this approach often does not improve performance [7].

Another approach is to use tempered distributions. This approach with a linear annealing schedule was originally recommended when normalizing flows were introduced for variational inference [3]. AdaAnn [28], which we compare our work against, has been proposed more recently as an adaptive scheduler that adjusts the annealing schedule based on the change of the KL divergence. Relatedly, the approach of using temperature as a conditional variable has been explored in variational inference using a finite set of temperatures [12]. More recently, Annealing Flow [29] proposed a continuous normalizing flow model guided by annealed intermediate distributions. While these methods share our motivation of addressing multi-modal posteriors through tempering, our approach differs by simultaneously tempering both base and target distributions in an affine-invariant manner, allowing for more stable training across temperature scales, and by eschewing the use of an annealing schedule.

Our approach is also related to Temperature-Steerable Flows [30], which uses transformations with constant Jacobians to create normalizing flow models for sampling physical systems across thermodynamic states. However, our work focuses more generally on variational inference applications, and we rely on the flow to learn to generalize across temperatures while avoiding temperature dependence in the affine component, and expose temperature as a conditional variable to the neural network.

Normalizing flow models have also been used to accelerate other inference paradigms. This includes using flows as transport maps to accelerate MCMC [17] and nested sampling [31] by reparameterizing the problem to simplify the distribution being sampled, or as proposal distributions for importance nested sampling [32]. Our work focuses on variational inference, though these efforts could be viewed as complementary as improved variational posteriors that could potentially be used in conjunction with these methods. The use of normalizing flow models in conjunction with other approaches to posterior sampling also extends to estimations of Bayesian evidence, either using importance sampling of the harmonic mean estimator [33], or using additional terms in the loss function to train a specialized normalizing flow and only using samples within the bulk of the posterior [34].

Finally, our approach shares conceptual similarities with parallel tempering or replica-exchange MCMC sampling [13, 35, 36, 14], where multiple MCMC chains are instantiated at different temperatures, with communication between chains ensuring that those at low temperature can traverse regions with low probability density. While these methods operate in the MCMC framework, our work translates similar temperature-based intuitions to the variational setting, leveraging the generalization capabilities of neural networks to simultaneously represent the posterior across the entire temperature spectrum, and using the simultaneous tempering of the base and target distributions to promote positive inductive biases that allow variational posteriors at different temperatures to interact.

4 Experiments

In our experiments, we compare FlowVAT, where we temper both the posterior and the base distribution and use temperature as a conditional variable, with posterior-only tempering where we do not temper the base distribution, and standard variational inference with a normalizing flow model (NF VI). We also include linear annealing [3] and AdaAnn [28], where the temperature schedule of annealing is adaptively determined based on local estimates of KL divergence between intermediate target distributions. All methods are given the same number of training epochs; as such, the AdaAnn hyperparameters are tuned such that the annealing schedule reaches $T = 1$ before the training ends.

We find that our method is relatively insensitive to the tempering hyperparameters, and thus sample temperatures uniformly in the range $T \in [0.95, 10]$ when training. Models are trained for 10^4 epochs using either this temperature range, when using a conditional tempering method, or the respective annealing schedule. Following that, models are fine tuned for 5×10^3 epochs, sampling uniformly in the interval $T \in [0.95, 1.5]$ or using $T = 1$ for the conditional and non-conditional methods, separately. For both linear annealing and AdaAnn, we adopt a temperature range of $T \in [1, 100]$, consistent with prior work [3, 28]. All normalizing flow models are implemented using `flowjax` [37] as 10-layer Rational-Quadratic neural spline flows with coupling layers [38], with 16 knots and an interval width of 8. Each coupling layer is parameterized by a dense neural network with 5 hidden layers of width 1024 and `silu` activations [39, 40, 41]. The base distribution is a standard multivariate Gaussian over the latent space. Additional hyperparameters can be found in appendix B.2.

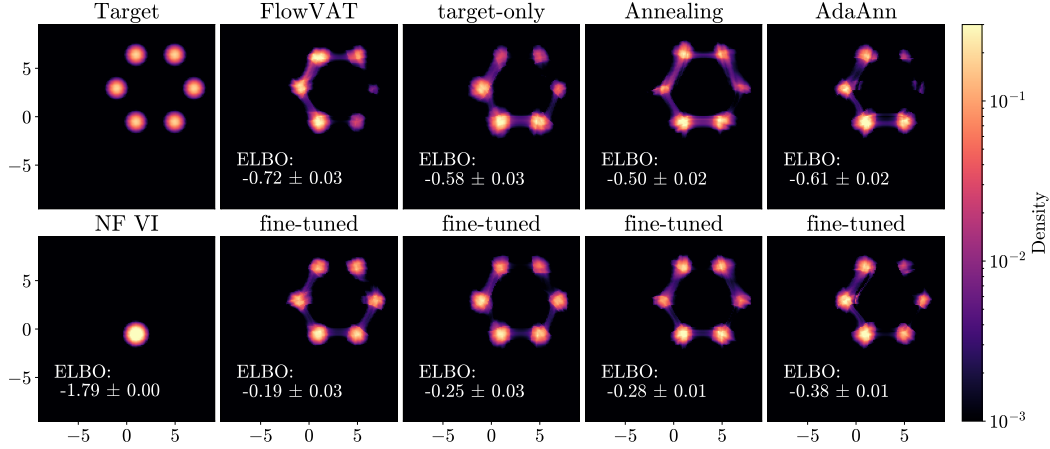


Figure 1: Comparison of different variational inference methods on 2D multimodal targets. The top row shows the target distribution and results from FlowVAT (this work), target-only tempering, traditional annealing, and AdaAnn. The bottom row includes normalizing flow VI and fine-tuned versions of each method, where additional training epochs are run at either $T = 1$ for the annealed methods or $T \in [0.95, 1.5]$ for the conditional tempering methods including FlowVAT. ELBO values with standard deviations are shown for each approach. We can see that our approach, FlowVAT, performs better than all other methods.

4.1 2D multimodal targets

In this experiment, we evaluate the ability of each algorithm to find separated modes in a 2D space. The target posterior has 6 Gaussian modes distributed in a ring; details of this distribution are given in appendix B.3. This target posterior, along with the variational posteriors after initial training and fine-tuning, are shown in fig. 1.

We find that all methods with tempering or annealing perform well after fine-tuning and finds all modes, while normalizing flow VI without tempering fails to find isolated modes. In addition, our conditional tempering method, FlowVAT, does not need the tempering hyperparameter to be tuned, presenting a key advantage for black-box VI, and also provides a slightly tighter ELBO.

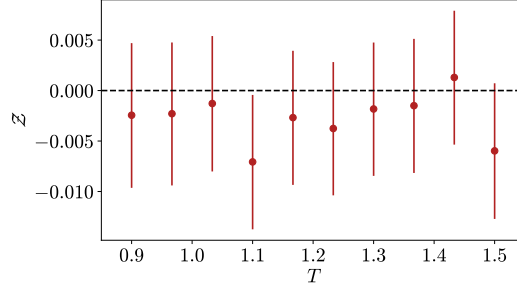


Figure 2: Estimated evidence via importance sampling (see eq. (10)) for the 2D multimodal target posterior along with 1σ errorbars. As the true posterior is well approximated by the variational posterior, we do not expect a strong temperature dependence.

We additionally estimate the true evidence using 5×10^4 importance samples, as described in eq. (10), with the FlowVAT variational posterior at various temperatures in the range $T \in [0.95, 1.5]$, as shown in fig. 2. We find that the estimated values are consistent with the ground truth, $Z = 0$. We do not observe a strong temperature dependence, as the variational posterior approximates the true posterior well, and thus even at temperatures slightly below one there is good coverage of the posterior.

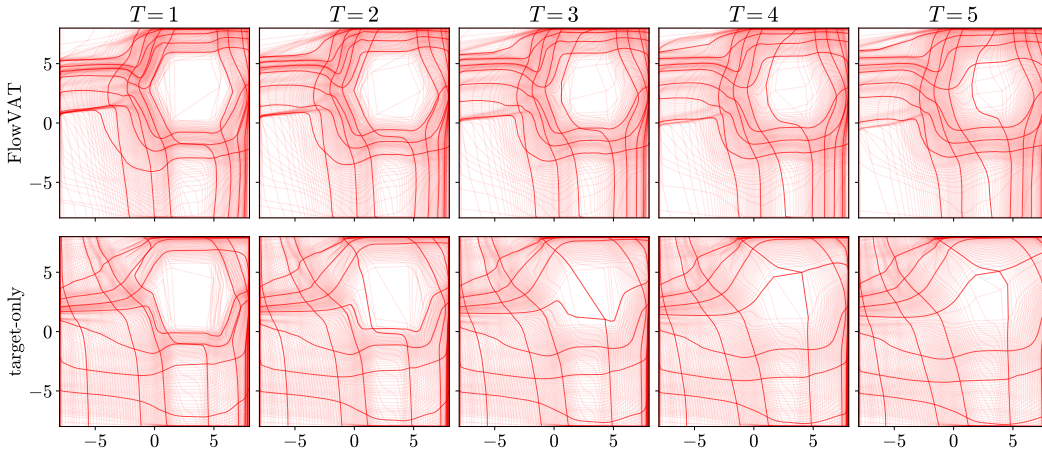


Figure 3: The transformation represented by the normalizing flow at various temperatures for the FlowVAT model (top), or target-only tempering (bottom). These plots are produced by transforming a regular grid using the normalizing flow models; as such, an identity transformation would be represented by a regular grid with a grid-spacing of 2. Models without fine-tuning are used to produce this diagram as we are including temperatures outside of the fine-tuning temperature range.

Finally, we can examine whether our theoretical arguments from section 2 regarding affine-invariant tempering are reflected in the trained flows. In particular, we would expect that the transformation changes much less as temperature varies when both the base and the target distributions are tempered, as is done with FlowVAT, as shown in fig. 3. We can see that the FlowVAT model indeed produces a normalizing flow where the transformations change much less with temperature when compared with target-only tempering, where the model has to fully adapt to the tempered posterior distribution.

4.2 Higher dimensional multimodal targets

To evaluate mode discovery performance in challenging settings, we construct synthetic Gaussian mixture posteriors (GMs) in 10 and 20 dimensions. Each mixture contains five equally weighted isotropic components with unit covariance. The mode locations are generated such that the pairwise distance between any two modes is at least the 0.99 quantile of the corresponding Gaussian, while

ensuring that each mode lies within the 0.999 quantile of at least one other mode. This setup is designed to mimic the clustered multimodal structures often observed in scientific posteriors, such as those in neutrino oscillation analyses [8, 9], searches for physics beyond the Standard Model [10], and gravitational wave inference [11].

The entire mixture is normalized so that the true log-evidence is zero. For each dimension, we generate 10 randomized instances of the GM and evaluate performance using 2000 samples from the variational posterior per trial. A mode is considered “captured” if more than 5% of the samples fall within the 0.9 quantile region of the corresponding Gaussian component. In addition, we report the estimated ELBO.

Table 1: Comparison of methods across datasets using number of modes captured and ELBO. Higher is better for both. The target posteriors contain 5 modes each. For the 10D posterior, standard variational inference performs best in terms of the ELBO, but does not find all modes reliably. FlowVAT outperforms all other methods for the more challenging 20D posterior, while performing the best out of all tempered/annealed methods for the 10D posterior.

Method	GM-10D		GM-20D	
	Modes	ELBO	Modes	ELBO
NF VI	4.5 ± 0.5	-0.25 ± 0.07	1.0 ± 0.0	-1.6 ± 0.0
FlowVAT (this work)	5.0 ± 0.0	-0.45 ± 0.05	3.4 ± 0.5	-1.2 ± 0.2
Target-only	5.0 ± 0.0	-0.55 ± 0.06	3.0 ± 0.7	-1.5 ± 0.2
Annealing	4.7 ± 0.4	-0.68 ± 0.08	1.8 ± 0.6	-1.6 ± 0.2
AdaAnn	5.0 ± 0.0	-0.49 ± 0.08	2.1 ± 0.7	-1.7 ± 0.2

Results are summarized in table 1. While no single method achieves uniformly superior performance across all metrics and dimensions, FlowVAT demonstrates strong and consistent performance. In 20 dimensions, it captures the most modes on average and achieves the highest ELBO, indicating a better fit to the target distribution. In 10 dimensions, where mode discovery is easier—as indicated by NF VI recovering most modes and achieving the highest ELBO—several methods including FlowVAT, target-only tempering, and AdaAnn successfully recover all modes.

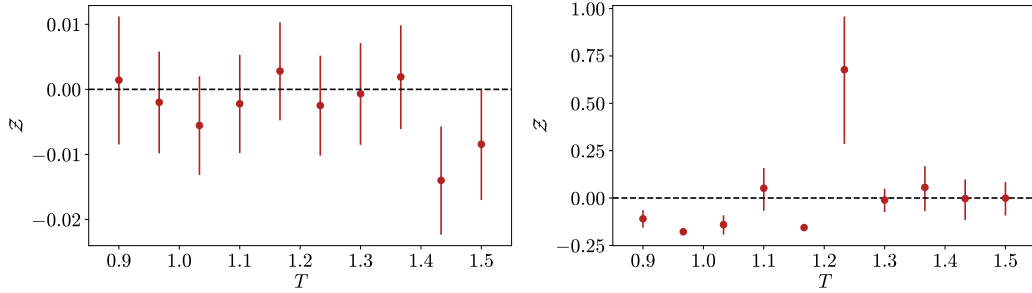


Figure 4: Estimated evidence via importance sampling (see eq. (10)) for instances of the 10D (left) and 20D (right) multimodal target posteriors along with 1σ errorbars. As the true posterior is well approximated by the variational posterior in the 10D case, there is no temperature dependence. In the more challenging 20D case, we can see that the estimated evidence converges for $T \geq 1.3$.

We hypothesize that FlowVAT’s advantage in higher dimensions stems from its affine-invariant tempering design, which reduces the burden on the flow model to learn complicated geometric transformations. In d dimensions, a full affine transformation involves $d(d+1)$ parameters. By tempering both the target and base distributions in an affine-invariant manner, FlowVAT allows the normalizing flow to focus on modeling the residual non-affine structure, rather than spending capacity on learning global affine transformations. This leads to improved performance, particularly in higher dimensions where such transformations become more complex. Additional implementation details for this experiment are provided in Appendix B.4.

We then estimate the Bayesian evidence using FlowVAT for one randomly chosen example for the 10 and 20 dimensional posteriors each using 5×10^4 importance samples. This can be seen in fig. 4. We

can see that in both cases, we can obtain values that are consistent with the ground truth of zero. In the 10 dimensional case, the computed evidence does not have a significant dependence on temperature, similar to the two dimensional case in section 4.1. In 20 dimensions, however, as the variational posterior is a poorer approximation of the target posterior, the computed evidence shows a clear bias at low temperature, with data points more than two standard deviations away from the truth values. Despite this, the evidence computed via importance sampling is stabilized for $T \geq 1.3$, once again giving values compatible with the ground truth of zero.

4.3 Benchmarks using eight-schools model

Table 2: Evidence Lower Bound (ELBO) comparison on the Eight Schools model. Higher is better.

Method	ELBO
NF VI	-31.337 ± 0.002
FlowVAT (this work)	-31.365 ± 0.048
Target-only	-31.483 ± 0.050
Annealing	-31.371 ± 0.005
AdaAnn	-31.344 ± 0.003

We then evaluate all methods on the non-centered Eight Schools model. The model is implemented in NumPyro [42, 43] following the formulation in posteriordb [44], serving as a benchmark for assessing posterior approximation quality.

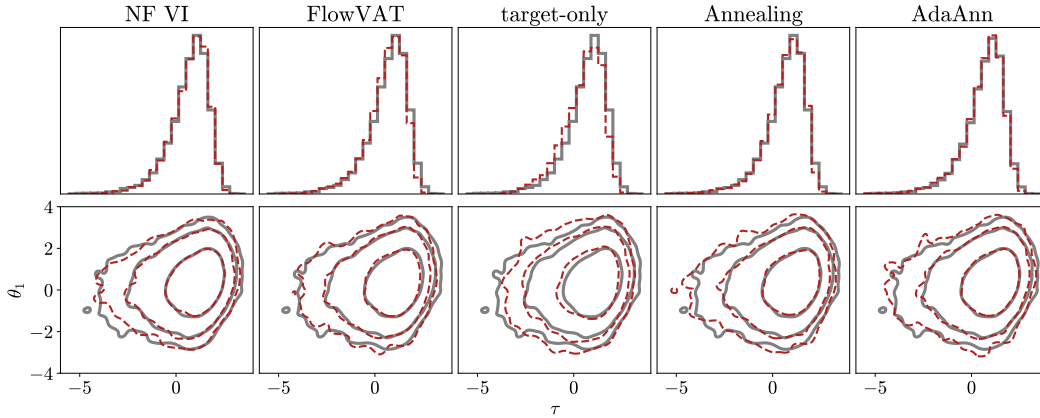


Figure 5: Comparison of variational inference methods on the eight schools model. Top row: Marginal histograms of θ_1 . Bottom row: Contour plots of the joint distribution between θ_1 and τ ; additional plots are shown in appendix B.5. Contours represent the 68%, 95%, and 99% credible regions, with thick gray lines for the reference posterior and dashed red lines for each variational approximation.

We run a single training run per method and estimate the ELBO using 5000 samples. The reported uncertainties in table 2 correspond to the standard error of the Monte Carlo estimate. To ensure adequate support coverage for this model, we set the spline transformation interval to 32 during training. All methods achieve similar ELBO values, and visual inspection of the posterior approximations (fig. 5) shows close agreement with the reference across both marginal and joint distributions.

Among the methods compared, tempering only the target distribution performs slightly worse in terms of both ELBO and posterior shape. This shows that ablating the tempering of the base distribution adversely affects performance, and is consistent with our hypothesis that jointly tempering the base and target in an affine-invariant way reduces the burden on the flow to learn affine transformations.

5 Discussion

Our experiments demonstrate that FlowVAT offers substantial advantages over existing variational inference methods, particularly with challenging posteriors. We discuss the key implications of these results and their broader impact.

The ability to handle complicated multi-modal posteriors in high dimensions represents a significant advancement in variational inference [7]. Traditional variational methods frequently struggle with this combination, often exhibiting mode-seeking behaviour even when using expressive variational families like normalizing flows [7]. By simultaneously tempering both the base distribution and the target posterior, FlowVAT effectively addresses this fundamental limitation without the need for an annealing schedule. As demonstrated in section 4, our approach captures more modes and achieves better ELBO values than competing methods, with this advantage becoming more pronounced as dimensionality increases. Such high-dimensional and multimodal posteriors are often encountered in science, including in neutrino oscillation analyses [8, 9], searches for physics beyond the Standard Model [10], and gravitational wave inference [11].

Beyond the performance improvements, FlowVAT represents a meaningful step toward truly hyperparameter-free black-box variational inference. A persistent challenge in practical applications of variational inference has been the need to carefully tune hyperparameters [5], particularly for difficult posteriors that require annealing. The requirement for such manual tuning undermines the "black-box" nature that makes variational inference attractive as a general-purpose inference method. Instead of requiring problem-specific schedule tuning, temperature is given to FlowVAT as a conditional variable, allowing the model to learn appropriate transformations across the entire temperature spectrum.

Our approach also introduces a beneficial inductive bias through affine-invariant tempering. As we argue for in section 2, this approach promotes consistent matching of the first two moments across temperatures. This property is particularly valuable in real-world inference problems where the target distribution may have complex structure and be located far from the origin in parameter space. The combination of this inductive bias with the generalization capabilities of overparameterized neural networks enables our model to efficiently transfer knowledge between different temperature regimes. In particular, we observe that ablation of the base distribution tempering adversely affects performance across all test cases, in support of our theoretical arguments.

In addition, the ability to sample from the variational posterior with a conditional temperature enables more reliable evidence estimation through importance sampling. Estimation of Bayesian evidence is useful for model comparison and hypothesis testing, but is traditionally very computationally intensive [34]. As demonstrated in our evidence computation experiments in section 4, FlowVAT enables fast and stable importance sampling of the evidence integral, and thus has the potential to accelerate model comparison and hypothesis testing as well.

In summary, FlowVAT addresses three key challenges in variational inference: capturing complicated multi-modal structure in high dimensions, enabling stable estimation of the Bayesian evidence, and reducing dependency on annealing hyperparameters. These advances bring us closer to the goal of developing truly automatic black-box variational inference methods capable of handling the diverse and challenging posteriors encountered in real-world applications.

5.1 Limitations and future work

One key limitation is that we have only explored inference in a continuous parameter space, as normalizing flows reparameterize such continuous spaces. Future work could explore architectural innovations needed to use such an approach when some parameters are discrete, for example by combining a conditional normalizing flow with a learned categorical distribution. In addition, we did not conduct experiments with large imbalances in the probability mass assigned to various modes, as capturing multiple modes is most crucial when modes have comparable probability mass. Finally, in this work, we focused on introducing the key ideas; while we are releasing the code used to produce experiments in this paper, we hope to release a user friendly module in the future.

6 Acknowledgment

We thank Ivy Li for his comments on the manuscript. This work is supported by the Department of Energy AI4HEP program and the National Science Foundation CAREER award PHY-204659.

References

- [1] Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 1995.
- [2] Rajesh Ranganath, Sean Gerrish, and David M. Blei. Black box variational inference. In *International Conference on Artificial Intelligence and Statistics*, 2013.
- [3] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- [4] Ryan Giordano, Martin Ingram, and Tamara Broderick. Black box variational inference with a deterministic objective: Faster, more accurate, and even more black box. *J. Mach. Learn. Res.*, 25:18:1–18:39, 2023.
- [5] Abhinav Agrawal, Daniel R Sheldon, and Justin Domke. Advances in black-box vi: Normalizing flows, importance weighting, and optimization. *Advances in Neural Information Processing Systems*, 33:17358–17369, 2020.
- [6] Roman Soletskyi, Marylou Gabri  , and Bruno Loureiro. A theoretical perspective on mode collapse in variational inference, 2024.
- [7] Akash Kumar Dhaka, Alejandro Catalina, Manushi Welandawe, Michael R Andersen, Jonathan Huggins, and Aki Vehtari. Challenges and opportunities in high dimensional variational inference. *Advances in Neural Information Processing Systems*, 34:7787–7798, 2021.
- [8] Ivan Esteban, MC Gonzalez-Garcia, Michele Maltoni, Ivan Martinez-Soler, and Jordi Salvado. Updated constraints on non-standard interactions from global analysis of oscillation data. *Journal of High Energy Physics*, 2018(8):1–33, 2018.
- [9] Pilar Coloma, MC Gonzalez-Garcia, Michele Maltoni, Jo  o Paulo Pinheiro, and Salvador Urrea. Global constraints on non-standard neutrino interactions with quarks and electrons. *Journal of High Energy Physics*, 2023(8):1–42, 2023.
- [10] Andrew Fowlie, Artur Kalinowski, Malgorzata Kazana, Leszek Roszkowski, and Y. L. Sming Tsai. Bayesian Implications of Current LHC and XENON100 Search Limits for the Constrained MSSM. *Phys. Rev. D*, 85:075012, 2012.
- [11] Maximilian Dax, Stephen R. Green, Jonathan Gair, Nihar Gupte, Michael P  rrer, Vivien Raymond, Jonas Wildberger, Jakob H. Macke, Alessandra Buonanno, and Bernhard Sch  lkopf. Real-time inference for binary neutron star mergers using machine learning. *Nature*, 639(8053):49–53, 2025.
- [12] Stephan Mandt, James McInerney, Farhan Abrol, Rajesh Ranganath, and David M. Blei. Variational tempering. In *International Conference on Artificial Intelligence and Statistics*, 2014.
- [13] Robert H. Swendsen and Jian-Sheng Wang. Replica monte carlo simulation of spin-glasses. *Phys. Rev. Lett.*, 57:2607–2609, 1986.
- [14] Nikola Surjanovic, Saifuddin Syed, Alexandre Bouchard-Cot  e, and Trevor Campbell. Parallel tempering with a variational reference. In *Neural Information Processing Systems*, 2022.
- [15] Andrew Gordon Wilson. Deep learning is not so mysterious or different, 2025.
- [16] Rylan Schaeffer, Mikail Khona, Zachary Robertson, Akhilan Boopathy, Kateryna Pistunova, Jason W. Rocks, Ila Rani Fiete, and Oluwasanmi Koyejo. Double descent demystified: Identifying, interpreting & ablating the sources of a deep learning puzzle, 2023.
- [17] Matthew Hoffman, Pavel Sountsov, Joshua V. Dillon, Ian Langmore, Dustin Tran, and Srinivas Vasudevan. NeuTra-lizing Bad Geometry in Hamiltonian Monte Carlo Using Neural Transport. *arXiv e-prints*, page arXiv:1903.03704, 2019.
- [18] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- [19] Greg Ashton et al. Nested sampling for physical scientists. *Nature*, 2, 2022.
- [20] John Skilling. Nested Sampling. *AIP Conference Proceedings*, 735(1):395–405, 2004.
- [21] John Skilling. Nested sampling for general Bayesian computation. *Bayesian Analysis*, 1(4):833 – 859, 2006.
- [22] Fernando Llorente, Luca Martino, David Delgado, and Javier Lopez-Santiago. Marginal likelihood computation for model selection and hypothesis testing: an extensive review. *SIAM Rev.*, 65:3–58, 2020.
- [23] Aki Vehtari, Daniel Simpson, Andrew Gelman, Yuling Yao, and Jonah Gabry. Pareto smoothed importance sampling. *Journal of Machine Learning Research*, 25(72):1–58, 2024.
- [24] Yingzhen Li and Richard E. Turner. Rényi divergence variational inference. In *Neural Information Processing Systems*, 2016.
- [25] Yuri Burda, Roger Baker Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *CoRR*, abs/1509.00519, 2015.
- [26] Dilin Wang, Hao Liu, and Qiang Liu. Variational inference with tail-adaptive f-divergence. In *Neural Information Processing Systems*, 2018.
- [27] Jean-Baptiste Regli and Ricardo Silva. Alpha-beta divergence for variational inference. *arXiv preprint arXiv:1805.01045*, 2018.
- [28] Emma R Cobian, Jonathan D Hauenstein, Fang Liu, and Daniele E Schiavazzi. Adaann: Adaptive annealing scheduler for probability density approximation. *International Journal for Uncertainty Quantification*, 13(3):39–68, 2023.
- [29] Dongze Wu and Yao Xie. Annealing flow generative model towards sampling high-dimensional and multi-modal distributions. *arXiv preprint arXiv:2409.20547*, 2024.
- [30] Manuel Dibak, Leon Klein, Andreas Krämer, and Frank Noé. Temperature steerable flows and boltzmann generators. *Physical Review Research*, 4(4):L042005, 2022.
- [31] Dorian Amaral, Shixiao Liang, Juehang Qin, and Christopher Tunnell. Fast bayesian inference for neutrino non-standard interactions at dark matter direct detection experiments. *Machine Learning: Science and Technology*, 6(1):015049, 2025.
- [32] Johannes U Lange. nautilus: boosting Bayesian importance nested sampling with deep learning. *Monthly Notices of the Royal Astronomical Society*, 525(2):3181–3194, 08 2023.
- [33] Alicja Polanska, Matthew A. Price, Davide Piras, Alessio Spurio Mancini, and Jason D. McEwen. Learned harmonic mean estimation of the Bayesian evidence with normalizing flows, 5 2024.
- [34] Rahul Srinivasan, Marco Crisostomi, Roberto Trotta, Enrico Barausse, and Matteo Breschi. Bayesian evidence estimation from posterior samples with normalizing flows. *Phys. Rev. D*, 110(12):123007, 2024.
- [35] Enzo Marinari and Giorgio Parisi. Simulated tempering: a new monte carlo scheme. *EPL*, 19:451–458, 1992.
- [36] David J. Earl and Michael W. Deem. Parallel tempering: theory, applications, and new perspectives. *Physical chemistry chemical physics : PCCP*, 7 23:3910–6, 2005.
- [37] Daniel Ward, Tennessee Hickling, and Matthew Mould. danielward27/flowjax: v12.2.0, May 2024.
- [38] Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows. *ArXiv*, abs/1906.04032, 2019.
- [39] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [40] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks : the official journal of the International Neural Network Society*, 107:3–11, 2017.
- [41] Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions, 2017.
- [42] Du Phan, Neeraj Pradhan, and Martin Jankowiak. Composable effects for flexible and accelerated probabilistic programming in numpyro. *arXiv preprint arXiv:1912.11554*, 2019.

- [43] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal probabilistic programming. *J. Mach. Learn. Res.*, 20:28:1–28:6, 2019.
- [44] Måns Magnusson, Jakob Torgander, Paul-Christian Bürkner, Lu Zhang, Bob Carpenter, and Aki Vehtari. posteriordb: Testing, benchmarking and developing bayesian inference algorithms. *arXiv preprint arXiv:2407.04967*, 2024.

A Broader Impacts

FlowVAT represents an advance in variational inference, which has seen significant use in scientific research. As such, we expect that improving the performance of variational inference will create positive societal impacts by accelerating research and scientific discovery. We do not expect this work to produce negative societal impacts beyond knock-on effects from research that uses our methods.

B Experiments and Details

B.1 Computational resource

All experiments were run in a containerized environment with 8 CPU cores and 64 GB of RAM. A single NVIDIA A40 GPU (48 GB memory) was used for all experiments. The training time remains similar (~ 18 min) for all the methods, which is expected as the model architecture stays unchanged.

B.2 Experiment hyperparameters

Table 3: Additional hyperparameters that were not included in the main text.

Hyperparameter	2D	10D	20D	Eight schools
Optimizer	adamW	adamW	adamW	adamW
Pre-training learning rate	5×10^{-6}	5×10^{-6}	5×10^{-6}	5×10^{-6}
fine-tuning learning rate	10^{-6}	10^{-6}	10^{-6}	10^{-6}
AdaAnn tolerance	0.02	0.1	0.15	0.1

For the linear annealing and AdaAnn methods, the temperature T is adjusted every 100 epochs.

B.3 2D multimodal distribution parameters

The 2D multimodal distribution is centered around $[3, 3]$, with 6 evenly spaced modes placed on a ring. The radius of the ring is 4, and the standard deviation of each mode is 0.38, chosen such that NF VI consistently fails to find all but one mode.

B.4 Target posterior generation for high dimension multimodal experiment

The generation of the high dimensional multimodal targets rely on the quantiles of multivariate normal distribution. In d -dimensional space, the squared Euclidean norm of a sample from $\mathcal{N}(0, \sigma^2 I_d)$ follows a scaled chi-squared distribution:

$$\|x\|^2 \sim \sigma^2 \cdot \chi_d^2.$$

To define a threshold distance that contains a specified quantile (e.g., 95%) of the Gaussian mass, we compute the square root of the corresponding chi-squared quantile:

$$\delta_{\text{quantile}} = \sigma \cdot \sqrt{F_{\chi_d^2}^{-1}(\alpha)},$$

where $F_{\chi_d^2}^{-1}(\alpha)$ is the inverse CDF (quantile function) of the chi-squared distribution with d degrees of freedom, and $\alpha \in (0, 1)$ is the desired quantile level (e.g., $\alpha = 0.95$).

This threshold is used to set reasonable values for `min_dist` or `max_dist` when sampling Gaussian centers in \mathbb{R}^d , ensuring that the modes are neither too tightly clustered nor overly dispersed. As mentioned in section 4.2, a 0.99 quantile is used as `min_dist` and a 0.999 quantile is used as `max_dist`. The Gaussian centers are generated by sampling points in a hypercube using the algorithm described in Algorithm 1.

A preliminary visual check is done to inspect if the model fits the multimodal target. An example is shown in fig. 6. Each sample from the normalizing flow model is assigned to the closest center, and the distribution of the absolute distance to the center is compared with the one from a Gaussian distribution.

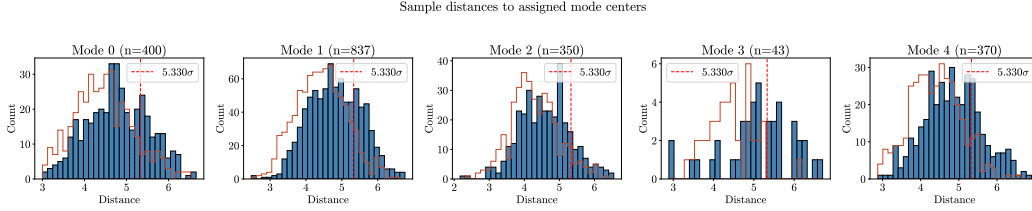


Figure 6: An example of visual check of the training result for FlowVAT on 20D Gaussian mixture distribution. The distribution of absolute distance between the center of the modes and the assigned samples is shown for each mode. The blue histogram represents the samples from the normalizing flow model and the red steps histogram represents the expected distribution from a 20D Gaussian distribution. The dashed red vertical line marks the 0.9 quantile. Here, Mode 3 is considered not captured as the ratio of the samples within the threshold is smaller than 0.05.

Algorithm 1 Generation of Gaussian Centers

Require: Number of centers K , dimension d , minimum distance d_{\min} , maximum distance d_{\max} , batch size M , random seed

Ensure: Set of K well-separated Gaussian centers $\mathcal{C} = \{c_1, \dots, c_K\} \subset \mathbb{R}^d$

```

1: Initialize random number generator with given seed
2: Sample initial center  $c_1 \sim \mathcal{U}([-5, 5]^d)$ 
3:  $\mathcal{C} \leftarrow \{c_1\}$ 
4: while  $|\mathcal{C}| < K$  and max_tries not exceeded do
5:   Sample  $M$  candidates  $\{x_1, \dots, x_M\} \sim \mathcal{U}([-5, 5]^d)$ 
6:   for each candidate  $x_i$  do
7:     Compute distances  $D_i = \{\|x_i - c\| : c \in \mathcal{C}\}$ 
8:     if  $\min D_i > d_{\min}$  and  $\min D_i < d_{\max}$  then
9:        $\mathcal{C} \leftarrow \mathcal{C} \cup \{x_i\}$ 
10:    if  $|\mathcal{C}| = K$  then
11:      break
12:    end if
13:  end if
14: end for
15: end while
16: return  $\mathcal{C}$ 

```

B.5 Posterior plots for the eight schools experiment

In the full corner plots of the inferred posterior versus reference posterior (fig. 7, fig. 8, fig. 9, fig. 10 and fig. 11), contours represent the 68% and 95% credible regions. Thick gray lines represent the reference posterior and dashed red lines represent variational approximation.

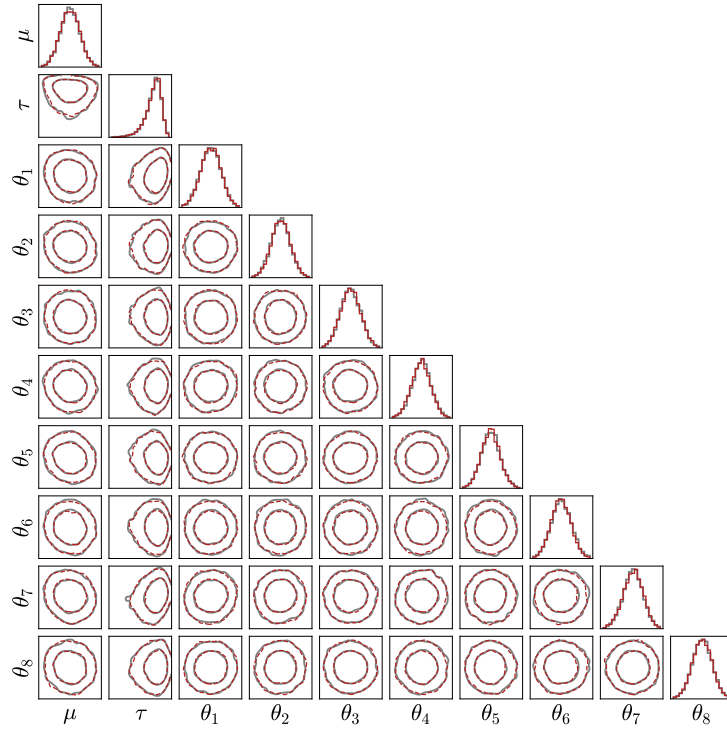


Figure 7: Full corner plot of NF VI for the Eight Schools model

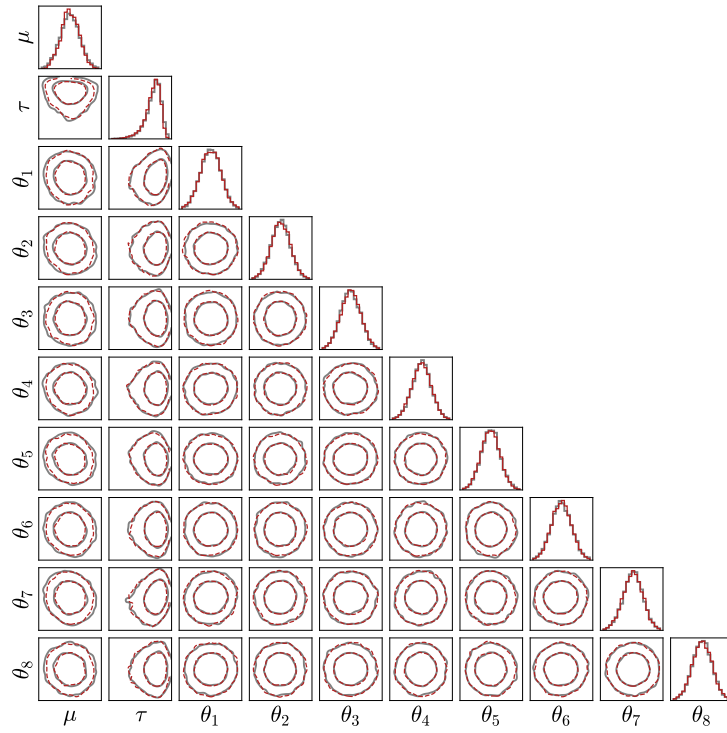


Figure 8: Full corner plot FlowVAT for Eight Schools model

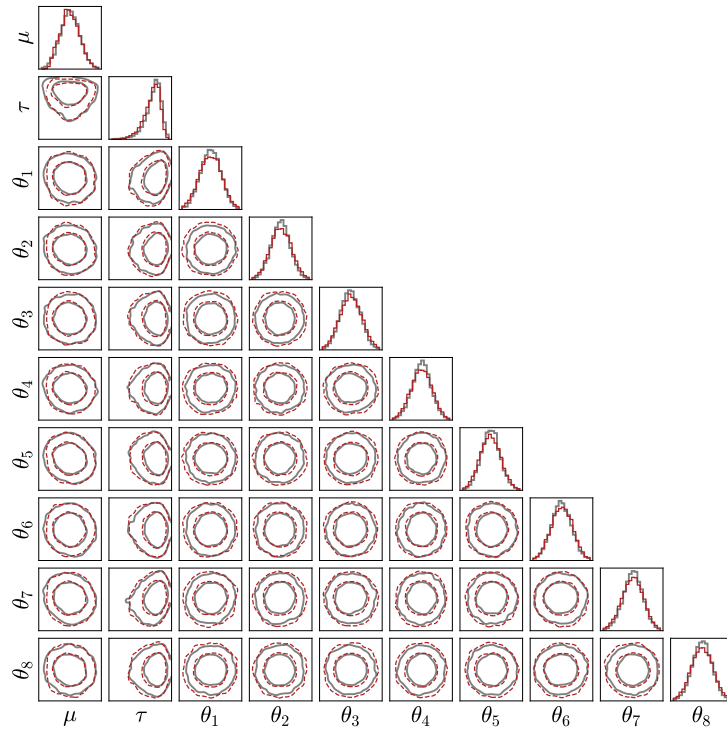


Figure 9: Full corner plot target-only tempering for Eight Schools model

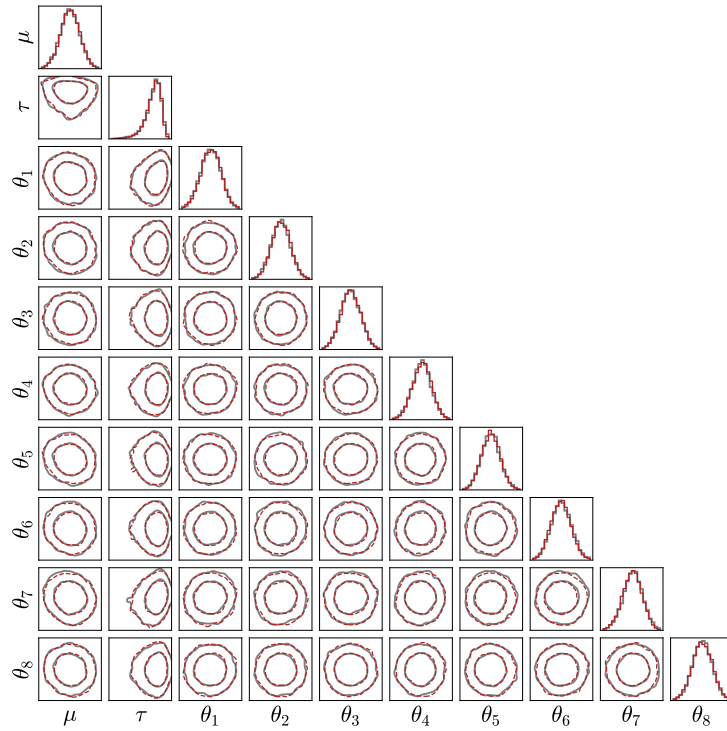


Figure 10: Full corner plot Annealing for Eight Schools model

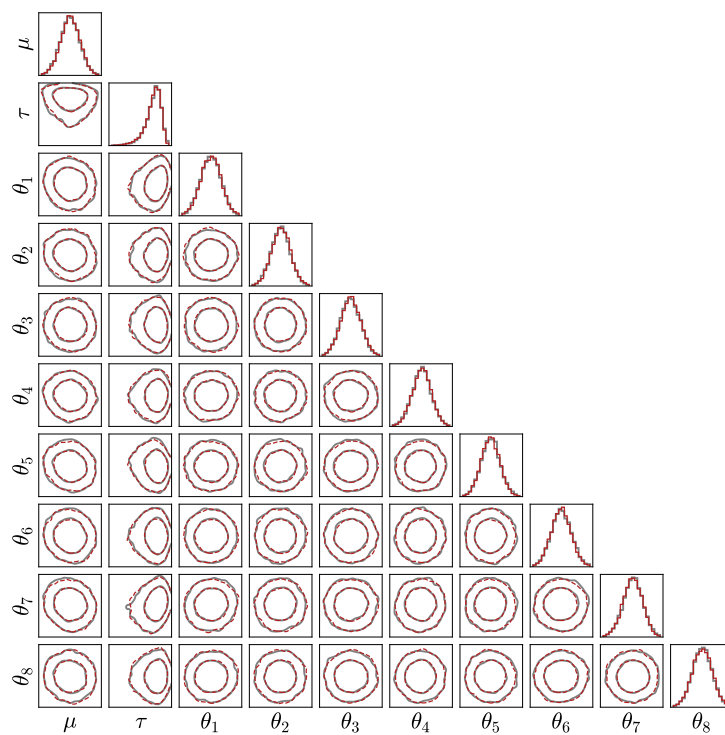


Figure 11: Full corner plot AdaAnn for Eight Schools model