

# FLOW-BASED VARIATIONAL MUTUAL INFORMATION: FAST AND FLEXIBLE APPROXIMATIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Mutual Information (MI) is a fundamental measure of dependence between random variables, but its practical application is limited because it is difficult to calculate in many circumstances. Variational methods offer one approach by introducing an approximate distribution to create various bounds on MI, which in turn is an easier optimization problem to solve. In practice, the variational distribution chosen is often a Gaussian, which is convenient but lacks flexibility in modeling complicated distributions. In this paper, we introduce new classes of variational estimators based on Normalizing Flows that extend the previous Gaussian-based variational estimators. Our new estimators maintain many of the same theoretical guarantees while simultaneously enhancing the expressivity of the variational distribution. We experimentally verify that our new methods are effective on large MI problems where discriminative-based estimators, such as MINE and InfoNCE, are fundamentally limited. Furthermore, we compare against a diverse set of benchmarking tests to show that the flow-based estimators often perform as well, if not better, than the discriminative-based counterparts. Finally, we demonstrate how these estimators can be effectively utilized in the Bayesian Optimal Experimental Design setting for online sequential decision making.

## 1 INTRODUCTION

Mutual Information (MI), a fundamental measure of dependence from information theory, has found utility in a variety of fields such as Bayesian optimal experimental design (Lindley, 1956; Foster et al., 2019), representation learning (Alemi et al., 2016; Chen et al., 2016), and structure learning (Vinh et al., 2011). Despite having a simple interpretation, MI typically lacks a closed form solution making it a difficult measure to use in many practical settings. Straightforward sample-based estimates can be used to approximate MI, but these are often inefficient both in terms of computation and sample complexity. Moreover, these so-called nested Monte Carlo (NMC) estimators exhibit large finite sample bias that decays slowly (Zheng et al., 2018; Rainforth et al., 2018).

To resolve the computational issues associated with MI, many so-called discriminative based estimators have been proposed. These distribution-free approximations typically utilize the Donsker-Varadhan lower bound of MI (Donsker & Varadhan, 1975), with some examples including MINE (Belghazi et al., 2018) and NWJ (Nguyen et al., 2010). While these estimators have seen wide adoption in practice, McAllester & Stratos (2020) show that they require an exponential number of samples to estimate large MI values. This poor sample complexity presents a fundamental limitation in regimes where the achievable information is likely to be high.

Generative variational methods, referred to simply as variational methods for the remainder of the paper, directly approximate the intractable target distribution with a tractable surrogate have seen promise in practical settings (Barber & Agakov, 2004; Foster et al., 2019; 2020). These variational methods make MI estimation more scalable than NMC based approaches and are not subject to the same fundamental limitation of approximating large MI as discriminative based estimators. Yet computation of these estimators can be prohibitively costly, particularly in sequential decision making regimes (Foster et al., 2021; Ivanova et al., 2021). Moreover, the accuracy of variational estimators is heavily dependent on the expressibility of the chosen surrogate distribution, which is often limited to a Gaussian approximation to control computation (Dahlke et al., 2023; Pacheco & Fisher III, 2019; Barber & Agakov, 2004).

We address the computational and accuracy aspects of variational MI estimators by introduction of two new estimators based on normalizing flows (Kobyzev et al., 2021). The proposed estimators flexibly adapt to complex target distributions and allow a straightforward tradeoff between accuracy and computational complexity. The first estimator, *Joint Variational Flow (JVF)*, builds on the *Joint Variational Gaussian (JVG)* estimator proposed in Dahlke et al. (2023). That previous work established that, for Gaussian approximating families, the estimator can be easily computed by moment-matching the target distribution. We relax the Gaussian assumption and show that the same efficient updates can be used in a flow-based estimator while simultaneously achieving more flexible density approximation.

Both JVG and JVF assume access to a joint distribution of the random quantities of interest; Gaussian in the first case and in the second a more flexible generative distribution given by reversing the normalizing flow. While this assumption of an underlying joint leads to efficient computation, it can restrict the MI approximation. The *Neural Variational Gaussian (NVG)* relaxes this assumption using a neural-network parameterized Gaussian distribution to approximate the posterior. We extend this estimator to *Neural Variational Flow (NVF)* to achieve a more expressible estimator that incorporates flow-based density estimates into the underlying approximation.

**Contributions** We summarize the main results of our paper as follows. We introduce 2 families of variational MI estimators (JVF and NVF) that incorporate normalizing flows to flexibly adapt to complex multimodal target distributions. We extend existing results to show that parameters of the underlying flow distribution can be efficiently computed via moment-matching operations. We provide a wide array of benchmarking experiments to highlight the strengths and weaknesses of each estimator as well as to compare to many of the state of the art MI estimators. Finally, we show that our variational estimators flexibly adapt to complex sequential decision making tasks with nonlinear non-Gaussian noise distributions.

## 2 PRELIMINARIES

Consider a joint distribution  $p(X, Y)$  over latent variable  $X$  and observation variable  $Y$ . Mutual information (MI) measures the dependence between these random quantities and is given by:

$$I(X, Y) = H_p(p(X)) - H_p(p(X|Y)) \quad (1)$$

where  $H_p(p(X)) = \mathbb{E}[-\log p(X)]$  is the *marginal entropy* and  $H_p(p(X|Y)) = \mathbb{E}[-\log p(X|Y)]$  is the *conditional entropy*. The entropy expectations are taken with respect to the joint  $p(X, Y)$ . Despite its simplicity, MI typically lacks a closed-form solution. Furthermore, nested Monte Carlo estimators of MI have large finite sample bias that is slow to decay (Zheng et al., 2018; Rainforth et al., 2018) making direct sample-based estimates of MI practically infeasible in many settings. To address these issues, variational methods introduce an approximate distribution,  $q$ , to estimate the underlying true distribution,  $p$  by minimizing Kullback-Leibler divergence:  $\min_q \text{KL}(q \| p)$ .

### 2.1 VARIATIONAL MUTUAL INFORMATION

The MI in Eqn. (1) requires knowledge of the marginal and posterior distributions. We consider an approach where both distributions are approximated,  $p(X) \approx q_{\text{marg}}(X)$  and  $p(X|Y) \approx q_{\text{post}}(X|Y)$ , in the *marginal and posterior approximation* (Foster et al., 2019; Dahlke et al., 2023):

$$I_{m+p}(X, Y) := H_p(q_{\text{marg}}(X)) - H_p(q_{\text{post}}(X|Y)) \approx I(X, Y) \quad (2)$$

where  $H_p(q(\cdot)) = \mathbb{E}_p[-\log q(\cdot)]$  is the cross-entropy. We can then minimize the following upper bound on absolute error (Foster et al., 2019):

**Lemma 2.1.** *For any model  $p(X, Y)$  and distributions  $q_{\text{marg}}(X)$ ,  $q_{\text{post}}(X | Y)$ , the following holds:*

$$|I_{m+p} - I| \leq \min_{q_{\text{marg}}} H_p(q_{\text{marg}}(X)) + \min_{q_{\text{post}}} H_p(q_{\text{post}}(X | Y)) + C$$

where  $C = -H_p(p(X)) - H_p(p(X | Y))$  does not depend on  $q_{\text{marg}}$  or  $q_{\text{post}}$ . Further, the RHS is 0 iff  $q_{\text{marg}}(X) = p(X)$  and  $q_{\text{post}}(X | Y) = p(X | Y)$  almost surely.

This bound can be made arbitrarily tight by finding good variational approximations of both  $p(X)$  and  $p(X | Y)$  and has the added benefit of only assuming that  $p(X, Y)$  can be sampled from. Thus,  $I_{m+p}$  can be applied in many settings, such as those with implicit (e.g. simulation-based) distributions.

## 2.2 JOINT VARIATIONAL GAUSSIAN (JVG) ESTIMATOR

Dahlke et al. (2023) show that under a Gaussian approximating distribution the upper bound in Lemma 2.1 can be efficiently minimized by moment-matching operations.

**Theorem 2.2** (Moment Matching = Optimization). *Let  $q(X, Y)$  be a joint Gaussian density:*

$$q(X, Y) = \mathcal{N} \left( \begin{bmatrix} X \\ Y \end{bmatrix} \middle| \mu := \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma := \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right) \quad (3)$$

then the marginal and posterior optimizing the bound in Lemma 2.1 are given by:

$$q_{\text{marg}}(X) = \mathcal{N}(X \mid \mu_x, \Sigma_{xx}), \quad q_{\text{post}}(X \mid Y) = \mathcal{N}(X \mid \mu_{x|y}, \Sigma_{x|y}), \quad \text{where} \quad (4)$$

$$\mu_{x|y} := \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (Y - \mu_y), \quad \Sigma_{x|y} := \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \quad (5)$$

where the mean and covariance are matched to the moments of the target density  $p(X, Y)$ :

$$\mathbb{E}_p[(X, Y)^T] = \mu, \quad \text{Cov}_p(X, Y) = \Sigma. \quad (6)$$

Furthermore, the value of  $I_{m+p}$  can be calculated in closed-form as the mutual information of the moment-matched Gaussian (Dahlke et al., 2023). We refer to this as the *Joint Variational Gaussian (JVG)* estimator:

$$I_{\text{JVG}} := H_p(q_{\text{marg}}(X)) - H_p(q_{\text{post}}(X|Y)) = \frac{1}{2} \log |\det(2\pi e \Sigma_{xx})| - \frac{1}{2} \log |\det(2\pi e \Sigma_{x|y})|. \quad (7)$$

These results ensure that, for Gaussian approximating  $q$ , the approximate MI  $I_{m+p}$  can be efficiently computed via simple moment-matching operations. The strength of this approach is that it does not require gradient descent and is very fast to compute in practice. The drawback of this approach is that it can only model linear dependence between random variables.

## 2.3 NEURAL VARIATIONAL GAUSSIAN (NVG) ESTIMATOR

To capture nonlinear dependence we relax the assumption that  $X$  and  $Y$  are jointly Gaussian (Foster et al., 2019). A base Gaussian variational distribution is assumed, but now the mean and variance of the posterior distribution are parameterized by a neural network which is a function of the observation variable  $y$ :

$$q(X) = \mathcal{N}(X \mid \mu_x, \Sigma_{xx}), \quad q(X \mid Y) = \mathcal{N}(X \mid \mu(Y), \Sigma(Y)) \quad (8)$$

where  $\mu(Y)$  and  $\Sigma(Y)$  are given by neural network function approximators trained to minimize the bound in Lemma 2.1. Marginal moments  $\mu_x$  and  $\Sigma_{xx}$  in Eqn. (8) are still learned by moment matching to the true distribution. This results in the *Neural Variational Gaussian (NVG)* estimator:

$$I_{\text{NVG}} := H_p(q(X)) - H_p(q(X|Y)) \approx \frac{1}{2} \log |\det(2\pi e \Sigma_{xx})| + \frac{1}{N} \sum_{i=1}^N \log \mathcal{N}(x_i \mid \mu(y_i), \Sigma(y_i)) \quad (9)$$

where  $\{(x_i, y_i)\} \sim p$ . The relaxed assumption of a Gaussian joint distribution allows for  $I_{\text{NVG}}$  to approximate nonlinear dependencies between  $X$  and  $Y$ , however this comes at the computational cost of needing to train a neural network. Furthermore,  $I_{\text{NVG}}$  is still constrained in its approximation power by assuming a Gaussian fit for each conditional value as well as the prior.

## 3 FLOW-BASED VARIATIONAL ESTIMATORS

Previous work considered  $I_{\text{JVG}}$  and  $I_{\text{NVG}}$  as jointly- and conditionally-Gaussian variational approximations for MI, respectively (Dahlke et al., 2023; Foster et al., 2019). Both of these approaches are limited by the expressibility of the Gaussian to an arbitrary distribution. To address this, we introduce a more flexible class of distributions based on normalizing flows. The aim of normalizing flows is to transform a simple base distribution, typically a Gaussian, via a diffeomorphism to achieve a more expressive distribution. Let  $Z \in \mathbb{R}^D$  be a random variable with density  $q_Z(Z)$  and let  $Z = f(X)$  be a diffeomorphism. Then the density of  $X$  is given by the change of basis formula (Bishop & Nasrabadi, 2006):

$$q_X(X) = q_Z(f(X)) |\det(\nabla_x f(X))| \quad (10)$$

where  $\nabla_x f(X)$  is the Jacobian of  $f(X)$  with respect to  $X$ . The choice of  $f(\cdot)$  can be any diffeomorphism, however, to train effectively we must be able to evaluate  $f$  and its log determinant efficiently. There have been many proposed flows (Kobyzev et al., 2021) including planar and radial (Rezende & Mohamed, 2016), coupling (Kingma et al., 2016), and auto-regressive flows (Kingma et al., 2016; Papamakarios et al., 2021). For this work we consider rational quadratic auto-regressive flows (Durkan et al., 2019) due to their success in related fields (Kobyzev et al., 2021).

### 3.1 JOINT VARIATIONAL FLOW (JVF) ESTIMATOR

Theorem 2.2 states that moment matching yields optimal parameters of a variational estimator to minimize Lemma 2.1, but only applies for a Gaussian approximation. We introduce a new estimator, based on normalizing flows, that adds flexibility to the variational estimator but that satisfies the moment matching optimality conditions.

**Theorem 3.1** (Moment Matched Flow Distribution). *Let  $Z = f(X)$  and  $V = g(Y)$  be diffeomorphisms with Gaussian joint density given by:*

$$q_{Z,V}(Z, V) = \mathcal{N} \left( \begin{bmatrix} Z \\ V \end{bmatrix} \middle| \mu := \begin{bmatrix} \mu_z \\ \mu_v \end{bmatrix}, \Sigma := \begin{bmatrix} \Sigma_{zz} & \Sigma_{zv} \\ \Sigma_{vz} & \Sigma_{vv} \end{bmatrix} \right) \quad (11)$$

then the marginal and posterior optimizing the bound in Lemma 2.1 are given by:

$$\begin{aligned} q_{\text{marg}}(X) &= \mathcal{N}(f(X) \mid \mu_z, \Sigma_{zz}) |\det(\nabla_x f(X))|, \\ q_{\text{post}}(X \mid Y) &= \mathcal{N}(f(X) \mid \mu_{z|v}, \Sigma_{z|v}) |\det(\nabla_x f(X))|, \\ \text{where } \mu_{z|v} &:= \mu_z + \Sigma_{zv} \Sigma_{vv}^{-1} (g(Y) - \mu_v), \quad \Sigma_{z|v} := \Sigma_{zz} - \Sigma_{zv} \Sigma_{vv}^{-1} \Sigma_{vz} \end{aligned} \quad (12)$$

where the mean and covariance are matched to the moments of the target density  $p(X, Y)$ :

$$\mathbb{E}_p[(f(X), g(Y))^T] = \mu, \quad \text{Cov}_p(f(X), g(Y)) = \Sigma. \quad (13)$$

Theorem 3.1 establishes that the optimal flow distribution is moment-matched to the target distribution, providing a simple solution to parameters of the flow distribution. The joint density on  $X$  and  $Y$  also takes a convenient form given by the change of basis formula in Eqn. (10):

$$q_{X,Y}(X, Y) = q_{Z,V}(f(X), g(Y)) |\det(\nabla_x f(X))| |\det(\nabla_y g(Y))| \quad (14)$$

$$= \mathcal{N} \left( \begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma \right) |\det(\nabla_x f(X))| |\det(\nabla_y g(Y))| \quad (15)$$

where the Jacobian terms can be easily computed by construction of the flow. What remains is to train parameters of the flows to minimize Lemma 2.1. The following lemma establishes that this bound has a convenient form, which can be optimized via gradient descent.

**Lemma 3.2** (Flow Upper Bound). *Let  $p(X, Y)$  be an arbitrary target distribution and  $q_{X,Y}(X, Y)$  be the distribution of the form in Eqn. (14) with moment matched flow density  $q_{Z,V}(Z, V)$ . Then the bound in Lemma 2.1 is given by:*

$$|I_{m+p} - I| \leq \frac{1}{2} \log |\det(2\pi e \Sigma_{zz})| + \frac{1}{2} \log |\det(2\pi e \Sigma_{z|v})| - 2\mathbb{E}_{p_X} [\log |\det(\nabla_x f(X))|] + C \quad (16)$$

We see that the marginal and conditional entropy of  $q_{X,Y}$  can be expressed in terms of the underlying Gaussian entropy of  $q_{Z,V}(Z, V)$ , which has a closed form resulting in the first two terms of Eqn. (16). The last term is the log determinant of the Jacobian, which is easily computed by construction of the normalizing flow. The flows are trained to minimize Eqn. (16) and the parameters  $\mu$  and  $\Sigma$  are learned from moment matching at each step of the flow training. This results in the *Joint Variational Flow (JVF)* estimator:

$$I_{\text{JVF}} := H_p(q_{\text{marg}}(X)) - H_p(q_{\text{post}}(X|Y)) = \frac{1}{2} \log |\det(2\pi e \Sigma_{zz})| - \frac{1}{2} \log |\det(2\pi e \Sigma_{z|v})| \quad (17)$$

We see that evaluating the MI does not require the log determinant term introduced by the flow which is the property of MI invariance under invertible transformations (Czyż et al., 2023). While the resulting MI estimator captures only linear dependence (via Gaussian MI) it does so in the flow distribution after application of nonlinear flows. The resulting estimator captures nonlinear dependence in the original  $X$  and  $Y$ .

### 3.1.1 STABLE TRAINING ERROR UPPER BOUND

The bound in Lemma 3.2 only sees the log determinant of  $f(X)$ , and in practice we found that  $g(Y)$  tends to overfit during learning. To address this, we introduce a new bound which is a slight modification to Lemma 2.1.

**Lemma 3.3.** *Let  $p(X, Y)$  be any model and  $q_{X,Y}(X, Y)$  be of the form in Eqn. (14) with base distribution  $q_{Z,V}(Z, V)$ , then the following bound holds:*

$$|I_{JVF} - I| \leq 2H_p(q_{Z,V}(Z, V)) - 2\mathbb{E}_{p_X} [\log |\det(\nabla_x f(X))|] - 2\mathbb{E}_{p_Y} [\log |\det(\nabla_y g(Y))|] + C \quad (18)$$

For a Gaussian base  $q_{Z,V}(Z, V) = \mathcal{N}(\mu, \Sigma)$  the tightest bound of the form Eqn. (18) is given by the moment-matched flow distribution with  $\Sigma = \text{Cov}_p(f(X), g(Y))$  and takes the form:

$$|I_{JVF} - I| \leq \log |\det(2\pi e \Sigma)| - 2\mathbb{E}_{p_X} [\log |\det(\nabla_x f(X))|] - 2\mathbb{E}_{p_Y} [\log |\det(\nabla_y g(Y))|] + C \quad (19)$$

where  $C = 2H_p(p(X, Y))$ . This bound is tight when  $q_{X,Y}(X, Y) = p(X, Y)$  almost surely.

Eqn. (19) contains both the log determinant terms of each flow and in practice enables more stable training. For all experiments in Section 6, we minimize  $I_{JVF}$  with respect to this bound. Calculation of the MI estimator remains unchanged from Eqn. (17). The pseudocode for training JVF can be found in Appendix A.1 in Algorithm 1.

### 3.2 NEURAL VARIATIONAL FLOW (NVF) ESTIMATOR

The estimator  $I_{JVF}$  assumes that  $X$  and  $Y$  have a joint distribution of the form Eqn. (14). While the underlying joint is non-Gaussian it can still result in restrictions on expressibility. To address this we introduce neural network function approximators to the base (flow) distribution. We add a flow to the marginal,  $Z_{\text{marg}} = f_{\text{marg}}(X)$ , and to the posterior,  $Z_{\text{post}} = f_{\text{post}}(X)$  with the following distributions:

$$q_Z(Z_{\text{marg}}) = \mathcal{N}(\mu_z, \Sigma_{zz}), \quad q_{Z|Y}(Z_{\text{post}} | Y) = \mathcal{N}(\mu(Y), \Sigma(Y)) \quad (20)$$

where  $\mu(\cdot)$  and  $\Sigma(\cdot)$  are given by neural network function approximators. Then for two normalizing flows,  $f_{\text{marg}}(\cdot)$  and  $f_{\text{post}}(\cdot)$ , the corresponding distributions on  $X$  are:

$$q_X(X) = \mathcal{N}(f_{\text{marg}}(X) | \mu_z, \Sigma_{zz}) |\det(\nabla f_{\text{marg}}(X))| \quad (21)$$

$$q_{X|Y}(X | Y) = \mathcal{N}(f_{\text{post}}(X) | \mu(Y), \Sigma(Y)) |\det(\nabla f_{\text{post}}(X))| \quad (22)$$

Furthermore, the cross entropy of each distribution is given by:

$$H_p(q_X(X)) = H_{q_Z}(q_Z(Z_{\text{marg}})) - \mathbb{E}_{p_X} [\log |\det(\nabla_x f_{\text{marg}}(X))|] \quad (23)$$

$$H_p(q_{X|Y}(X | Y)) = \mathbb{E}_{p_X} [-\log \mathcal{N}(f_{\text{post}}(X) | \mu(Y), \Sigma(Y))] - \mathbb{E}_{p_X} [\log |\det(\nabla_x f_{\text{post}}(X))|] \quad (24)$$

where  $H_{q_Z}(q_Z(Z_{\text{marg}}))$  is the entropy of the underlying Gaussian. The optimal marginal moments  $\mu_z$  and  $\Sigma_{zz}$  are found via moment matching, but both flows and the posterior parameters are optimized to minimize Lemma 3.3. This results in the *Neural Variational Flow (NVF)* estimator:

$$I_{\text{NVF}} = \frac{1}{2} \log |\det(2\pi e \Sigma_{zz})| - \frac{1}{N} \sum_{i=1}^N \log |\det(\nabla_x f_{\text{marg}}(x_i))| \\ + \frac{1}{N} \sum_{i=1}^N \log \mathcal{N}(f_{\text{post}}(x_i) | \mu(y_i), \Sigma(y_i)) + \frac{1}{N} \sum_{i=1}^N \log |\det(\nabla_x f_{\text{post}}(x_i))| \quad (25)$$

where  $\{(x_i, y_i)\} \sim p$  and  $\Sigma_{zz} = \text{Cov}(f_{\text{marg}}(X))$ . The log determinant terms do not cancel in this case, since we apply separate transforms to the marginal and posterior, and so MI is not invariant under these transformations.  $I_{\text{NVF}}$  has both the flexibility of a flow-based distribution and nonlinear dependence modeling.  $I_{\text{NVF}}$  is capable of estimating much more complex distributions but comes at the cost of needing to train both flows and neural network parameters. The pseudocode for training NVF can be found in Appendix A.1 in Algorithm 2.

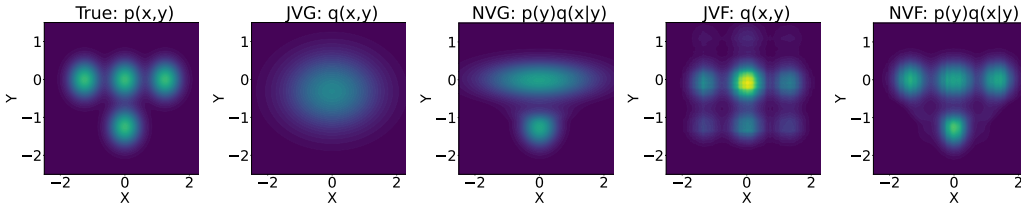


Figure 1: **Multi-modal Density** A four component Gaussian mixture model with nonlinear dependency is used to show the approximating capabilities of each estimator. JVG cannot model the nonlinearity or multi-modality. NVG is able to model the nonlinear dependence between rows, but cannot capture the multi-modality. JVF is able to model the multi-modality across the system but "hallucinates" extra modes to satisfy its linearity constraint. Finally, NVF models the multi-modality and the nonlinear dependency of the GMM.

#### 4 PRACTICAL APPLICATION OF VARIATIONAL MI ESTIMATORS

Estimators JVG, NVG, JVF, and NVF have different flexibility incorporated to improve accuracy in estimating MI. Moving from a joint to neural Gaussian increases modeling capability, but comes at the cost of adding a neural network to train and losing access to an explicit joint distribution. Although the joint distribution is not necessary for MI estimation, it can be useful to have a parameterized joint for applications such as generative models and classification. The second flexibility added to the estimators is flows which enhance the estimators capability of capturing multi-modality and nonlinear structures of the distribution but comes at the cost of needing to train two flows.

Figure 1 illustrates these properties on a four-component Gaussian mixture model. We see that JVG can only match the mean and spread of the data from the joint Gaussian assumption. NVG is able to capture the change of dependency from the top three modes to the bottom but is unable to capture the multimodality for any individual  $y$ . JVF on the other hand is able to model the multi-modality but still is constrained to linear dependence of the transformed variable from the joint assumption. This means that JVF "hallucinates" alternate modes to satisfy its linear dependence constraint. Finally, NVF has both capabilities and therefore is capable of capturing both the multi-modality of the distribution and the non-linearity of the dependence. It is important to note that for both the neural estimators no explicit joint  $q(X, Y)$  exists, instead the joint is estimated by  $p(Y)q(X | Y)$  for illustrative purposes.

The choice of estimator depends on the prior knowledge about the underlying model. JVG is suitable for simple cases, while NVG captures more complex dependencies. JVF is appropriate for multi-modal distributions with linear dependencies, and NVF handles both multi-modality and nonlinear dependencies. Table 1 summarizes the strengths and weaknesses of each estimator to guide the selection of the appropriate method.

Table 1: Flexibility of each estimator. Adding flows increases the variational density flexibility and adding neural parameters increases the dependency flexibility. Adding the neural parameters however removes the modeling of an explicit joint distribution.

Estimator	Density	Dependency	Training	Joint
$I_{JVG}$	Inflexible	Inflexible	None	Explicit
$I_{NVG}$	Inflexible	Flexible	$\mu(Y), \Sigma(Y)$	Implicit
$I_{JVF}$	Flexible	Inflexible	$f(X), g(Y)$	Explicit
$I_{NVF}$	Flexible	Flexible	$\mu(Y), \Sigma(Y), f_{\text{prior}}(X), f_{\text{post}}(X)$	Implicit

##### 4.1 BAYESIAN OPTIMAL EXPERIMENTAL DESIGN

A common application of MI is Bayesian Optimal Experimental Design (BOED) (Lindley, 1956) where the goal is to maximize information through a series of decisions,  $d$ , about a latent variable,  $X$ , given observations,  $Y$ . This model consists of an assumed prior on  $X$ ,  $p(X)$ , and a likelihood function,  $p(Y|X, d)$ . Each decision is quantified with the amount of information it will likely give, which in the context of BOED is called *expected information gain (EIG)*

$$EIG(d) = I_d(X, Y) = H_{p(X)} [p(X)] - H_{p(X, Y|d)} [p(X|Y, d)] \tag{26}$$

The optimal decision is the one that maximizes EIG:  $d^* = \operatorname{argmax}_d \text{EIG}(d)$ . In sequential experimental design, a sequence of decisions is made  $d_1, \dots, d_T$  and after each decision, an observation is made  $Y_1, \dots, Y_T$ . The observations are used to update the prior,  $p(X)$ , at the next time step using to posterior,  $p(X | Y_1, \dots, Y_t, d_1, \dots, d_t)$ , to include the observations from the previous decisions. When making  $T$  decisions, the goal is to maximize the *total expected information gain (TEIG)* (Ivanova et al., 2021)

$$\text{TEIG}_T(D) = \mathbb{E}_{p(X)p(h_T|X)} \left[ \sum_{t=1}^T \mathbb{E}_{p(X)p(Y_t|X, d_t, h_{t-1})} \left[ \log \frac{p(X | Y_t, d_t, h_{t-1})}{p(X | h_{t-1})} \right] \right] \quad (27)$$

where  $D = \{d_t\}_{t=1}^T$  is the set of decisions and  $h_t = \{(Y_i, d_i)\}_{i=1}^t$  is the *history* of previously taken decisions and their corresponding observations. We utilize  $I_{\text{JVG}}$ ,  $I_{\text{NVG}}$ ,  $I_{\text{JVF}}$ , or  $I_{\text{NVF}}$  as variational estimations of the EIG at each step in a greedy approach to maximizing TEIG. We learn the corresponding parameters for each model simultaneously with the decision which is an approach introduced by Foster et al. (2020) where they utilized the NVG estimator. We highlight this application for the use in Section 6 to show the utility of our estimators in practice.

## 5 RELATED WORK

The base MI approximation,  $I_{m+p}$ , considered in this work is one of many possible variational approximations. Choosing only one distribution to approximate with a variational distribution results in a bound of MI. These are considered in alternative work (Foster et al., 2019; Dahlke et al., 2023; Poole et al., 2019) but predominantly utilize only Gaussians as the variational distribution. Canonical correlation analysis (CCA) assumes a joint Gaussian variational distribution which is closely related to  $I_{\text{JVG}}$ . Recent work from (Butakov et al., 2024) has proposed a similar approach also utilizing normalizing flows which is closely related to  $I_{\text{JVF}}$ , but does not consider an estimator analogous to  $I_{\text{NVF}}$ . Furthermore, (Cheng et al., 2020) proposed an upper bound of MI that utilizes a variational distribution similar to that of  $I_{\text{NVG}}$ .

An alternative approach to approximating the distributions is to instead use a discriminator. The Donsker-Varadhan representation and the corresponding f-divergence representation lead to tight lower bounds on MI, resulting in a variety of approaches: DV (Donsker & Varadhan, 1975), MINE (Belghazi et al., 2018), and NWJ (Nguyen et al., 2010). Other methods use noise-contrastive estimation of the density ratio in MI which leads to a lower-bound of MI, called InfoNCE (van den Oord et al., 2019). Approaches that do not assume a discriminator or a density but instead compare  $k$ -nearest neighbors of samples have been considered called KSG (Kraskov et al., 2004). For a thorough review of MI approximations see (Poole et al., 2019).

In the space of BOED multiple approaches have been taken to select a series of actions that maximizes MI. Multiple greedy approaches have been considered where each design is chosen to maximize instantaneous EIG. The approach of simultaneously optimizing the variational distributions and decision via gradient decent has been considered (Foster et al., 2020). Other approaches, such as LFIRE (Kleinegesse et al., 2021), perform ratio estimation of MI with Bayesian optimization for the decisions. Batch optimization is the approach where all decisions are made before testing time and stay fixed regardless of the observed data, such as MINEBED (Kleinegesse & Gutmann, 2020). Finally, recent work takes a reinforcement learning approach where a policy-discriminator pair is learned (Foster et al., 2021; Ivanova et al., 2021; Blau et al., 2022; Huan & Marzouk, 2016). The discriminator plays the role of MI estimation for a select sequence of decision and observations while the policy is trained to make informative decisions based upon its history of decisions and observations.

## 6 EXPERIMENTS

We consider a wide range of experiments, from synthetic to application based. We start by comparing our MI estimation capabilities against many common estimators: DV (Donsker & Varadhan, 1975), MINE (Belghazi et al., 2018), NWJ (Nguyen et al., 2010), InfoNCE (van den Oord et al., 2019), CCA (Murphy, 2023), and KSG (Kraskov et al., 2004). The benchmarking tests are a large MI estimation (McAllester & Stratos, 2020) and a diverse set of 36 individual distributions with known ground truth

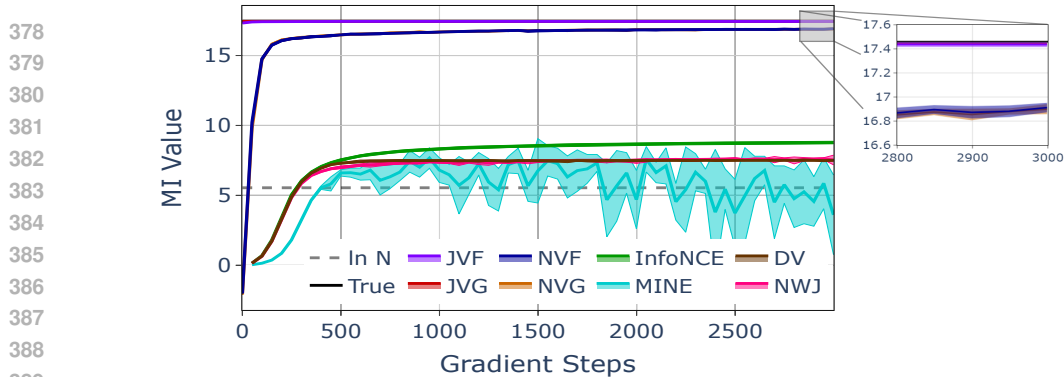


Figure 2: **Large MI** is difficult for many estimators to approximate. We see that the four distribution based estimators,  $I_{JVG}$ ,  $I_{NVG}$ ,  $I_{JVF}$ , and  $I_{NVF}$ , can accurately predict MI, but the four discriminative-based methods, MINE, InfoNCE, NWJ, and DV, all converge around  $\mathcal{O}(\log N)$ .

MI values (Czyż et al., 2023). We then look at the sequential decision making setting of location finding. All experiments were run on a high-performance computing cluster with nodes consisting of 2x AMD EPYC 7642 48-core (Rome) CPUs, 512GB of RAM, and NVIDIA V100S GPUs.

## 6.1 HIGH MUTUAL INFORMATION EXPERIMENT

We begin with a well-known difficult case of estimating large MI. McAllester & Stratos (2020) showed that any distribution-free high-confidence lower bound on MI, such as the commonly used discriminative-based methods, estimated from  $N$  samples cannot be larger than  $\mathcal{O}(\log N)$ . We create a Gaussian where  $X, Y \in \mathbb{R}^{15}$  with a correlation  $\rho = .95$ . The MI of this distribution can be exactly computed as:  $I(X, Y) = -\frac{Dim_X}{2} \ln(1 - \rho^2) = 17.459$ . We take a large set of 75,000 samples to train our distribution based estimators as well as a variety of discriminative based estimators for 3000 training steps using a batch size of  $N = 256$ . We utilize a 80-20 split of training and testing samples of the total samples. Figure 2 shows the best testing value for each estimator. We see that  $I_{JVG}$  and  $I_{JVF}$  converge nearly instantaneously to the true MI value and  $I_{NVG}$  and  $I_{NVF}$  both converge rapidly to high-quality estimates. For the discriminative-based methods, we see that they learn slower and converge to values drastically lower than the true MI.

## 6.2 MUTUAL INFORMATION BENCHMARK

Our next experiment is a collection of benchmarks from Czyż et al. (2023). They construct a diverse family of distributions with known ground-truth MI consisting of Gaussian, Uniform, and Student-T distributions that have MI invariant transformations applied to them (see Appendix B.2 for more details). We use  $N = 1,000$  samples per dimension, with a train-test split of 50-50. Figure 3 displays the average MI over 10 runs for each experiment. Due to the large number of experiments we focus on a selection here (see Appendix B.2 for all 36). For the majority of experiments we see that incorporating flow improves estimation accuracy, with JVF generally outperforming JVG. In the neural estimator case (NVF vs. NVG) the relative improvement is more limited, with a few important exceptions discussed next.

**Challenging cases** for all methods are the **spiral multinormal** and the **student-t distributions**. The spiral multinormal MI is substantially underestimated by all estimators, but we see that flow-based methods (JVF, NVF) typically outperform their non-flow counterparts (JVG, NVG) as well as the distribution-free baselines. For the student-t distributions, all methods produce unstable estimates due to high-variance (due to space we omit standard errors). The high variance for both JVG and NVG comes from the fact that for many cases moments of the student-t distribution are infinite or undefined, limiting effectiveness of moment-matching steps required by JVG and NVG. JVF and NVF can potentially address this issue with the flows, allowing the points to have finite moments, however this is particularly difficult to adequately transform the fat tails of the student-t with sparse samples. Czyż et al. (2023) introduce the **inverse hyperbolic sin transformation** as a means of normalizing the student-t tails, and thus creating finite moments. In this case we see  $I_{NVF}$  produces good results of the



True MI	0.41	0.41	0.33	0.41	0.41	0.41	0.41	1.02	1.02	1.02	1.02	1.02	1.02	1.02	1.02	1.02	0.43	0.29	0.43	0.29
MINE	-0.42	0.38	0.30	0.37	0.41	0.36	0.38	0.93	0.85	0.82	0.85	0.42	0.73	0.80	0.80	-0.02	0.15	0.25	0.10	
InfoNCE	-0.40	0.38	0.28	0.39	0.39	0.35	0.38	0.97	0.91	0.85	0.95	0.55	0.79	0.85	0.84	-0.15	0.21	0.36	0.21	
NWJ	-0.40	0.39	0.31	0.39	0.39	0.32	0.37	0.96	0.90	0.60	0.74	0.50	0.63	0.84	0.68	0.01	0.09	0.30	0.10	
DV	-0.40	0.38	0.30	0.38	0.40	0.35	0.38	0.96	0.90	0.87	0.92	0.50	0.78	0.85	0.85	-0.00	0.11	0.34	0.15	
CCA	-0.41	0.39	0.19	0.34	0.38	0.02	0.41	1.02	0.95	0.97	0.97	0.24	0.84	0.84	0.96	1.48	0.29	0.00	0.00	
KSG	-0.41	0.42	0.32	0.42	0.41	0.42	0.41	0.96	0.92	0.27	0.24	0.72	0.24	0.90	0.26	0.20	0.13	0.37	0.19	
JVG	-0.42	0.40	0.19	0.34	0.39	0.02	0.41	1.01	0.95	0.97	0.97	0.24	0.84	0.83	0.96	1.25	0.23	0.00	0.00	
JVF	-0.42	0.42	0.19	0.42	0.42	0.03	0.41	1.01	1.01	1.04	1.04	0.25	0.84	0.86	0.96	0.90	0.04	0.00	0.00	
NVG	-0.41	0.43	0.19	0.43	0.41	0.40	0.39	0.99	0.99	0.93	0.94	0.56	0.80	0.95	0.92	2.85	1.75	0.54	0.32	
NVF	-0.42	0.40	0.20	0.41	0.41	0.39	0.39	0.98	0.94	0.92	0.96	0.57	0.79	0.96	0.90	0.36	0.39	0.29		
Bivariate normal 1 x 1																				
Normal CDF @ Bivariate normal 1 x 1																				
Uniform 1 x 1 (additive noise=0.75)																				
Bimodal 1 x 1																				
Wiggly @ Bivariate normal 1 x 1																				
Swiss roll 2 x 1																				
Multinomial 3 x 3 (dense)																				
Normal CDF @ Multinomial 3 x 3 (2-pair)																				
Multinomial 3 x 3 (2-pair)																				
Half-cube @ Multinomial 25 x 25 (2-pair)																				
Spiral @ Multinomial 25 x 25 (2-pair)																				
Spiral @ Multinomial 3 x 3 (2-pair)																				
Spiral @ Normal CDF @ Multinomial 25 x 25 (2-pair)																				
Spiral @ Normal CDF @ Multinomial 3 x 3 (2-pair)																				
Student-t 2 x 2 (dof=1)																				
Student-t 3 x 3 (dof=2)																				
Asinh @ Student-t 2 x 2 (dof=1)																				
Asinh @ Student-t 3 x 3 (dof=2)																				

Figure 3: **MI Benchmark** runs against a selection of 19 different scenarios (for space; c.f. Appendix for all 36). Blue indicates an underestimate and red indicates an overestimate with the magnitude indicated by saturation. Blank cells indicate that the method produced divergent estimates.

MI, where all other methods seem to fail, suggesting the estimator is useful for fat tailed distributions with appropriate data regularization. Finally, we notice the **uniform additive noise** as the only case which the discriminative-based methods outperform distribution-based methods.

### 6.3 LOCATION FINDING

We now demonstrate our methods in application to BOED discussed in Section 4.1 on the location finding experiment (Ivanova et al., 2021). In this experiment there are multiple sources with unknown locations,  $X$ , producing a signal whose intensity decreases according to the inverse square law. The signals from all sources are summed together to create a total intensity,  $Y$  that can be nosily observed from any location,  $d$ . Our specific setting contains two sources in two dimensions, for a total of a four dimensional  $X \in \mathbb{R}^4$  and a one dimensional observation,  $Y \in \mathbb{R}^1$ . The goal is to make a sequence of  $T = 10$  decisions  $d_1, \dots, d_T \in \mathbb{R}^2$  that maximize the total expected information gain (Eqn. (27)). We consider iDAD (Ivanova et al., 2021) as a strong baseline in this experiment as iDAD performs non-myopic decision making and so should represent an upper bound on what is achievable by the greedy (myopic) methods JVG, JVG, NVG, and NVF. iDAD also requires more total computation time as it performs a significant offline training phase that other methods do not.

For our main results we see that incorporating flow uniformly improves MI estimation. When using the correct prior  $p(X) = \mathcal{N}(0, I)$  the JVG estimator has the lowest MI estimate due to assumption of a joint Gaussian distribution on  $X$  and  $Y$ . Incorporating flow (JVF) shows a significant improvement. The neural estimators show more flexibility but we see similar improvements when incorporating flow (NVF) to the NVG estimator. All greedy methods require less total computation time than the iDAD. Note that neural methods benefit from GPU acceleration that was not applied in our implementation of JVG or JVF. One drawback of amortized methods, such as iDAD, is that they do not adapt to model mismatch during the test phase; a result of the well-known *amortization gap* (Cremer et al., 2018). We test robustness of each method by using a mismatched prior in the test phase  $p(X) = \mathcal{N}(1, I)$ . The relative performance of all methods remains identical to the matched prior setting. But iDAD exhibits a 30% reduction in MI estimate whereas the greedy methods, which

Table 2: **Total Expected Information Gain** is reported for two settings of the location finding experiment, *matched prior* where the training prior matches the true prior and *mismatched prior* where the testing and training priors differ. We see that JVG, JVF, and NVG all fail to make informative decisions due to their limited flexibility and greedy decision making in both settings. We notice that iDAD performs well in the setting where its training data matches the testing data, but suffers a substantial amortization cost as the testing values are shifted. NVF yields informative decisions while being comparatively robust to model mismatch. Total computation time is lower for the variational methods when accounting for the offline training required by iDAD.

Method	Matched Prior	Mismatched Prior	Training Time	Deployment Time (s)
Random	$4.62 \pm 0.24$	$3.59 \pm 0.14$	-	-
iDAD (InfoNCE)	$7.54 \pm 0.19$	$5.29 \pm 0.21$	8826.683	$0.0256 \pm 0.0016$
JVG	$3.56 \pm 0.19$	$3.59 \pm 0.15$	-	$287 \pm 35$
JVF	$4.30 \pm 0.19$	$3.85 \pm 0.19$	-	$1730 \pm 207$
NVG	$4.61 \pm 0.19$	$3.99 \pm 0.20$	-	$186 \pm 21$
NVF	$5.10 \pm 0.20$	$4.748 \pm 0.18$	-	$1253 \pm 155$

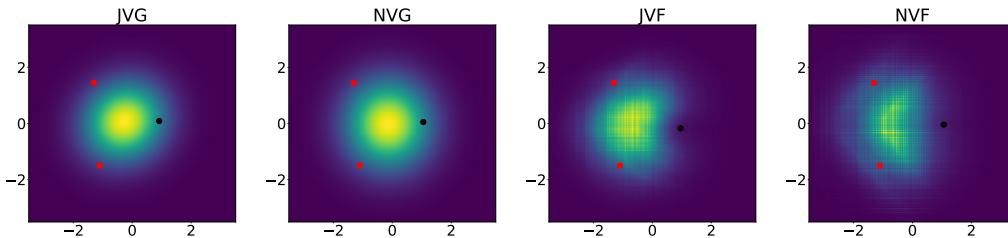


Figure 4: **Location Finding Posterior** estimations after a single observation. We notice that JVG, NVG both estimate the posterior to be a Gaussian. The introduction of flow to JVF and NVF increases flexibility to meaningfully approximate the posterior.

perform inference at test time either show no reduction in MI (JVG) or a smaller relative decrease (7%-13%).

We notice that in this experiment JVG, NVG, and JVF perform significantly worse than NVF. To shed some light on this, we can look at the posterior prediction after making a single observation (Fig.4). The result of this observation increases the belief that sources are located in a radius away from the origin. Neither JVG nor NVG are flexible enough to capture this belief and instead produce Gaussian posteriors. Including flow JVF and NVF have increased flexibility to capture the non-Gaussian belief. We see that the additional degrees of freedom in NVF allows for a higher posterior belief in the vicinity of the true target locations, whereas JVF concentrates more away from the targets.

## 7 DISCUSSION

We introduced two flow-based variational estimators,  $I_{JVF}$  and  $I_{NVF}$ , which are built on previous Gaussian estimators,  $I_{JVG}$  and  $I_{NVG}$ . The inclusion of flows increases the expressibility of the variational estimators. These estimators in general are just as accurate as widely used critic-based estimators, but are not limited at their capability of learning large MI values. Furthermore,  $I_{NVF}$  has the capability of estimating wide-tale distribution MI given appropriate data regularization, which are notoriously difficult distributions to estimate. Finally, we find that these methods are effective for sequential Bayesian experimental design. In the BOED setting these methods show robustness to model mismatch not exhibited by policy-based approaches while also yielding significantly lower computation time when training is accounted for.

**Limitations** Our proposed flow-based estimators have increased the flexibility of the distribution based variational MI estimators. However, there are still a few cases where we are unable to accurately model MI, such as additive noise in Section 6.2. Furthermore, each method adds increased number of parameters to train which can be costly, specifically at high dimensions. Future work plans to consider ways to increase training speed or remove the necessity of including neural networks.

## REFERENCES

- 540  
541  
542 Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational information  
543 bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- 544 David Barber and Felix Agakov. The IM algorithm: a variational approach to information maximiza-  
545 tion. *NIPS*, 16:201, 2004.
- 546 Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron  
547 Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference*  
548 *on machine learning*, pp. 531–540. PMLR, 2018.
- 550 Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer  
551 New York, 2006.
- 552 Tom Blau, Edwin V Bonilla, Iadine Chades, and Amir Dezfouli. Optimizing sequential experimental  
553 design with deep reinforcement learning. In *International conference on machine learning*, pp.  
554 2107–2128. PMLR, 2022.
- 556 Ivan Butakov, Alexander Tolmachev, Sofia Malanchuk, Anna Neopryatnaya, and Alexey Frolov.  
557 Mutual information estimation via normalizing flows, 2024. URL [https://arxiv.org/  
558 abs/2403.02187](https://arxiv.org/abs/2403.02187).
- 559 Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Info-  
560 gan: Interpretable representation learning by information maximizing generative adversarial nets.  
561 *Advances in neural information processing systems*, 29, 2016.
- 562 Pengyu Cheng, Weituo Hao, Shuyang Dai, Jiachang Liu, Zhe Gan, and Lawrence Carin. Club: A  
563 contrastive log-ratio upper bound of mutual information, 2020. URL [https://arxiv.org/  
564 abs/2006.12013](https://arxiv.org/abs/2006.12013).
- 566 Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders.  
567 In *International Conference on Machine Learning*, pp. 1078–1086. PMLR, 2018.
- 568 Paweł Czyż, Frederic Grabowski, Julia Vogt, Niko Beerenwinkel, and Alexander Marx. Be-  
569 yond normal: On the evaluation of mutual information estimators. In A. Oh, T. Nau-  
570 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural*  
571 *Information Processing Systems*, volume 36, pp. 16957–16990. Curran Associates, Inc.,  
572 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/  
573 file/36b80eae70ff629d667f210e13497edf-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/36b80eae70ff629d667f210e13497edf-Paper-Conference.pdf).
- 574 Caleb Dahlke, Sue Zheng, and Jason Pacheco. Fast variational estimation of mutual information for  
575 implicit and explicit likelihood models. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de  
576 Meent (eds.), *Proceedings of The 26th International Conference on Artificial Intelligence and*  
577 *Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pp. 10262–10278. PMLR,  
578 25–27 Apr 2023.
- 580 Monroe D Donsker and SR Srinivasa Varadhan. Asymptotic evaluation of certain markov process  
581 expectations for large time, i. *Communications on Pure and Applied Mathematics*, 28(1):1–47,  
582 1975.
- 583 Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows.  
584 In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garn-  
585 nett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Asso-  
586 ciates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
587 2019/file/7ac71d433f282034e088473244df8c02-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/7ac71d433f282034e088473244df8c02-Paper.pdf).
- 588 Adam Foster, Martin Jankowiak, Elias Bingham, Paul Horsfall, Yee Whye Teh, Thomas Rainforth,  
589 and Noah Goodman. Variational bayesian optimal experimental design. In *Advances in Neural*  
590 *Information Processing Systems 32*, pp. 14036–14047. 2019.
- 592 Adam Foster, Martin Jankowiak, Matthew O’Meara, Yee Whye Teh, and Tom Rainforth. A uni-  
593 fied stochastic gradient approach to designing bayesian-optimal experiments. In *International*  
*Conference on Artificial Intelligence and Statistics*, pp. 2959–2969. PMLR, 2020.

- 594 Adam Foster, Desi R. Ivanova, Ilyas Malik, and Tom Rainforth. Deep adaptive design: Amortizing  
595 sequential bayesian experimental design, 2021.
- 596
- 597 Xun Huan and Youssef M Marzouk. Sequential bayesian optimal experimental design via approximate  
598 dynamic programming. *arXiv preprint arXiv:1604.08320*, 2016.
- 599 Desi R. Ivanova, Adam Foster, Steven Kleinegesse, Michael U. Gutmann, and Tom Rainforth. Implicit  
600 deep adaptive design: Policy-based experimental design without likelihoods, 2021.
- 601
- 602 Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Im-  
603 proved variational inference with inverse autoregressive flow. In D. Lee, M. Sugiyama, U. Luxburg,  
604 I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29.  
605 Curran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper\\_](https://proceedings.neurips.cc/paper_files/paper/2016/file/ddeebdeefdb7e7e7a697e1c3e3d8ef54-Paper.pdf)  
606 [files/paper/2016/file/ddeebdeefdb7e7e7a697e1c3e3d8ef54-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/ddeebdeefdb7e7e7a697e1c3e3d8ef54-Paper.pdf).
- 607 Steven Kleinegesse and Michael U. Gutmann. Bayesian experimental design for implicit models by  
608 mutual information neural estimation, 2020.
- 609 Steven Kleinegesse, Christopher Drovandi, and Michael U Gutmann. Sequential bayesian exper-  
610 imental design for implicit models via mutual information. *Bayesian Analysis*, 16(3):773–802,  
611 2021.
- 612
- 613 Ivan Kobyzev, Simon J.D. Prince, and Marcus A. Brubaker. Normalizing flows: An introduction and  
614 review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43  
615 (11):3964–3979, 2021. doi: 10.1109/TPAMI.2020.2992934.
- 616 Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Phys.*  
617 *Rev. E*, 69:066138, Jun 2004. doi: 10.1103/PhysRevE.69.066138. URL [https://link.aps.](https://link.aps.org/doi/10.1103/PhysRevE.69.066138)  
618 [org/doi/10.1103/PhysRevE.69.066138](https://link.aps.org/doi/10.1103/PhysRevE.69.066138).
- 619 D. V. Lindley. On a measure of the information provided by an experiment. *The Annals of Mathemat-*  
620 *ical Statistics*, 27(4):986–1005, December 1956. ISSN 0003-4851.
- 621
- 622 David McAllester and Karl Stratos. Formal limitations on the measurement of mutual information.  
623 In *International Conference on Artificial Intelligence and Statistics*, pp. 875–884. PMLR, 2020.
- 624 Kevin P Murphy. *Probabilistic machine learning: Advanced topics*. MIT press, 2023.
- 625
- 626 XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals  
627 and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*,  
628 56(11):5847–5861, 2010.
- 629 Jason Pacheco and John Fisher III. Variational information planning for sequential decision making.  
630 In *AISTATS*, pp. 2028–2036, 2019.
- 631
- 632 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji  
633 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of*  
634 *Machine Learning Research*, 22(57):1–64, 2021. URL [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/19-1028.html)  
635 [19-1028.html](http://jmlr.org/papers/v22/19-1028.html).
- 636 Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. On variational  
637 bounds of mutual information. In *International Conference on Machine Learning*, pp. 5171–5180.  
638 PMLR, 2019.
- 639 Tom Rainforth, Robert Cornish, Hongseok Yang, and Andrew Warrington. On nesting monte carlo  
640 estimators. In *International Conference on Machine Learning*, pp. 4264–4273, 2018.
- 641
- 642 Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows, 2016.
- 643 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive  
644 coding, 2019.
- 645
- 646 Nguyen Xuan Vinh, Madhu Chetty, Ross Coppel, and Pramod P Wangikar. Globalmit: learning glob-  
647 ally optimal dynamic bayesian network with the mutual information test criterion. *Bioinformatics*,  
27(19):2765–2766, 2011.

Sue Zheng, Jason Pacheco, and John Fisher. A robust approach to sequential information theoretic planning. In *International Conference on Machine Learning*, pp. 5936–5944, 2018.

## A APPENDIX

### A.1 PSEUDOCODE

We briefly provide an outline of each methods pseudo code. So of the notation is lightly changed in an attempt to highlight parameters of models. For example, for NVF,  $f_{\text{prior}}(X)$  is referred to as  $f_{\theta_{\text{prior}}}(X)$  to highlight the parameters,  $\theta_{\text{prior}}$ , of the prior flow.

**for**  $i=1:K$  **do**

Sample  $\{x_j\}_{j=1:N} \sim p(X)$ ;  
 Sample  $\{y_j\}_{j=1:N} \sim p(y | x_j)$ ;  
 $\mu \leftarrow \frac{1}{N} \sum_{j=1}^N (f_{\theta}(x_j)^T, g_{\psi}(y_j)^T)^T$ ;  
 $\Sigma \leftarrow \frac{1}{N-1} \sum_{j=1}^N (f_{\theta}(x_j)^T, g_{\psi}(y_j)^T)(f_{\theta}(x_j)^T, g_{\psi}(y_j)^T)^T - \mu\mu^T$ ;  
 Loss  $\leftarrow \log |2\pi e \Sigma| - \frac{2}{N} \sum_{j=1}^N \log |\nabla_x f_{\theta}(x_j)| - \frac{2}{N} \sum_{j=1}^N \log |\nabla_y g_{\psi}(y_j)|$ ;  
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \text{Loss}$ ;  
 $\psi \leftarrow \psi - \beta \nabla_{\psi} \text{Loss}$ ;

**end**

$\mu \leftarrow \frac{1}{N} \sum_{j=1}^N (f_{\theta}(x_j)^T, g_{\psi}(y_j)^T)^T$ ;  
 $\Sigma \leftarrow \frac{1}{N-1} \sum_{j=1}^N (f_{\theta}(x_j)^T, g_{\psi}(y_j)^T)(f_{\theta}(x_j)^T, g_{\psi}(y_j)^T)^T - \mu\mu^T$ ;  
 $I_{JVF} \leftarrow \frac{1}{2} \log |2\pi e \Sigma_Z| - \frac{1}{2} \log |2\pi e \Sigma_{Z|V}|$

**Algorithm 1: JVF Pseudocode**

**for**  $i=1:K$  **do**

Sample  $\{x_j\}_{j=1:N} \sim p(X)$ ;  
 Sample  $\{y_j\}_{j=1:N} \sim p(y | x_j)$ ;  
 $\mu_Z \leftarrow \frac{1}{N} \sum_{j=1}^N f_{\theta_{\text{prior}}}(x_j)$ ;  
 $\Sigma_{ZZ} \leftarrow \frac{1}{N-1} \sum_{j=1}^N f_{\theta_{\text{prior}}}(x_j) f_{\theta_{\text{prior}}}(x_j)^T - \mu_Z \mu_Z^T$ ;  
 Loss  $\leftarrow \log |2\pi e \Sigma_{ZZ}| - \frac{1}{N} \sum_{j=1}^N \log |\nabla_x f_{\theta_{\text{prior}}}(x_j)| -$   
 $\frac{1}{N} \sum_{j=1}^N (\log \mathcal{N}(f_{\theta_{\text{post}}}(x_j) | \mu_{\phi}(y_i), \Sigma_{\phi}(y_i)) + \log |\nabla_x f_{\theta_{\text{post}}}(x_j)|)$ ;  
 $\theta_{\text{prior}} \leftarrow \theta_{\text{prior}} - \alpha \nabla_{\theta_{\text{prior}}} \text{Loss}$ ;  
 $\theta_{\text{post}} \leftarrow \theta_{\text{post}} - \alpha \nabla_{\theta_{\text{post}}} \text{Loss}$ ;  
 $\phi \leftarrow \phi - \beta \nabla_{\phi} \text{Loss}$ ;

**end**

$\mu_Z \leftarrow \frac{1}{N} \sum_{j=1}^N f_{\theta_{\text{prior}}}(x_j)$ ;  
 $\Sigma_{ZZ} \leftarrow \frac{1}{N-1} \sum_{j=1}^N f_{\theta_{\text{prior}}}(x_j) f_{\theta_{\text{prior}}}(x_j)^T - \mu_Z \mu_Z^T$ ;  
 $I_{JVF} \leftarrow \log |2\pi e \Sigma_{ZZ}| - \frac{1}{N} \sum_{j=1}^N \log |\nabla_x f_{\theta_{\text{prior}}}(x_j)| +$   
 $\frac{1}{N} \sum_{j=1}^N (\log \mathcal{N}(f_{\theta_{\text{post}}}(x_j) | \mu_{\phi}(y_i), \Sigma_{\phi}(y_i)) + \log |\nabla_x f_{\theta_{\text{post}}}(x_j)|)$

**Algorithm 2: NVF Pseudocode**

### A.2 SECTION 2 PROOFS

**Lemma 2.1.** For any model  $p(X, Y)$  and distributions  $q_{\text{marg}}(X)$ ,  $q_{\text{post}}(X | Y)$ , the following holds:

$$|I_{m+p} - I| \leq \min_{q_{\text{marg}}} H_p(q_{\text{marg}}(X)) + \min_{q_{\text{post}}} H_p(q_{\text{post}}(X | Y)) + C$$

where  $C = -H_p(p(X)) - H_p(p(X | Y))$  does not depend on  $q_{\text{marg}}$  or  $q_{\text{post}}$ . Further, the RHS is 0 iff  $q_{\text{marg}}(X) = p(X)$  and  $q_{\text{post}}(X | Y) = p(X | Y)$  almost surely.

702 *Proof.* We recreate the proof from Foster et al. (2019)

$$703 |I_{m+p} - I| = |H_p(q_{\text{marg}}(X)) - H_p(q_{\text{post}}(X | Y)) - H_p(p(X)) + H_p(p(X | Y))| \quad (28)$$

$$704 = |-H_p(p(X)) + H_p(q_{\text{marg}}(X)) + H_p(p(X | Y)) - H_p(q_{\text{post}}(X | Y))| \quad (29)$$

$$705 = |\text{KL}(p(X) \parallel q_{\text{marg}}(X)) - \text{KL}(p(X | Y) \parallel q_{\text{post}}(X | Y))| \quad (30)$$

$$706 \leq |\text{KL}(p(X) \parallel q_{\text{marg}}(X))| + |\text{KL}(p(X | Y) \parallel q_{\text{post}}(X | Y))| \quad (31)$$

$$707 = -H_p(p(X)) + H_p(q_{\text{marg}}(X)) - H_p(p(X | Y)) + H_p(q_{\text{post}}(X | Y)) \quad (32)$$

$$708 = H_p(q_{\text{marg}}(X)) + H_p(q_{\text{post}}(X | Y)) + C \quad (33)$$

709 where  $C = -H_p(p(X)) - H_p(p(X | Y))$   $\square$

710 To prove Theorem 2.2 and Theorem 3.1, we rely on results from Dahlke et al. (2023). In Section A.4,  
711 we include the related theorems and proofs necessary for the results of this paper as a convenience to  
712 the reader.

713 **Theorem 2.2** (Moment Matching = Optimization). *Let  $q(X, Y)$  be a joint Gaussian density:*

$$714 q(X, Y) = \mathcal{N} \left( \begin{bmatrix} X \\ Y \end{bmatrix} \middle| \mu := \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma := \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right) \quad (34)$$

715 then the marginal and posterior optimizing the bound in Lemma 2.1 are given by:

$$716 q_{\text{marg}}(X) = \mathcal{N}(X | \mu_x, \Sigma_{xx}), \quad q_{\text{post}}(X | Y) = \mathcal{N}(X | \mu_{x|y}, \Sigma_{x|y}), \quad \text{where} \quad (35)$$

$$717 \mu_{x|y} := \mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \quad \Sigma_{x|y} := \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx} \quad (36)$$

718 where the mean and covariance are matched to the moments of the target density  $p(X, Y)$ :

$$719 \mathbb{E}_p[(X, Y)^T] = \mu, \quad \text{Cov}_p(X, Y) = \Sigma. \quad (37)$$

720 *Proof.* It suffices to verify the assumption of the posterior expected statistics being a linear com-  
721 bination of joint statistics (Eqn. (97)) is satisfied. Recall the sufficient statistics of a multivariate  
722 Gaussian

$$723 T(X, Y) = \begin{bmatrix} X \\ Y \\ \text{vec}(XX^T) \\ \text{vec}(XY^T) \\ \text{vec}(YY^T) \end{bmatrix}$$

724 In this case  $\tau_1(Y) = y$  and  $\tau_2(Y) = \text{vec}(yy^T)$ . We now verify that the expected value under  
725  $q_{\text{post}}(X | Y)$  of each term in the sufficient statistic is a linear function of  $\tau_1(Y)$  and  $\tau_2(Y)$

726 1.  $x$

$$727 \mathbb{E}_{q_{\text{post}}(X|Y)}[x] = \mu_{x|y} = \mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y)$$

728 2.  $y$

$$729 \mathbb{E}_{q_{\text{post}}(X|Y)}[y] = y$$

730 3.  $xx^T$

$$731 \begin{aligned} 732 \mathbb{E}_{q_{\text{post}}(X|Y)}[xx^T] &= \Sigma_{x|y} + \mu_{x|y}\mu_{x|y}^T \\ 733 &= \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx} + (\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y))(\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y))^T \\ 734 &= \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx} + \mu_x\mu_x^T + \dots \\ 735 &\quad \mu_x(y - \mu_y)^T\Sigma_{yy}^{-1}\Sigma_{yx} + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y)\mu_x + \dots \\ 736 &\quad \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y)(y - \mu_y)^T\Sigma_{yy}^{-1}\Sigma_{yx} \end{aligned}$$

737 4.  $xy^T$

$$738 \begin{aligned} 739 \mathbb{E}_{q_{\text{post}}(X|Y)}[xy^T] &= (\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y))y^T \\ 740 &= (\mu_x - \Sigma_{xy}\Sigma_{yy}^{-1}\mu_y)y^T + \Sigma_{xy}\Sigma_{yy}^{-1}yy^T \end{aligned}$$

5.  $yy^T$

$$\mathbb{E}_{q_{\text{post}}(X|Y)} [yy^T] = yy^T$$

So the statistics are linear functions of  $\tau_1(Y) = y$  and  $\tau_2(Y) = yy^T$  so  $q_{\text{post}}(X | Y)$  satisfies the conditions of Theorem A.5 and moment matching the joint  $q(X, Y) = \mathcal{N}(m, \Sigma)$  yields the optimal  $q_{\text{post}}(X | Y)$ . Furthermore, the joint Gaussian sufficient statistics are

$$T(X, Y) = [X, Y, \text{vec}(XX^T), \text{vec}(XY^T), \text{vec}(YY^T)]^T$$

and the sufficient statistics of a marginal distribution are

$$T(X) = [X, \text{vec}(XX^T)]^T$$

which are simply the first and third sufficient statistic from the joint. Therefore moment matching the joint Gaussian trivially moment matches the marginal, giving the optimal  $q_{\text{marg}}(X)$ .  $\square$

Moment matching the joint Gaussian yields the optimal  $q_{\text{marg}}$  and  $q_{\text{post}}$ . With this we can derive the formula for the MI estimate  $I_{\text{JVG}}$ . We utilize another results from Dahlke et al. (2023)

**Lemma A.1.** *Analytic Entropy*

Let  $p(X)$  be any distribution and  $q(X)$  be in the exponential family with constant base measure,  $h(X) = C$ , which is analytically moment matched to  $p(X)$  and  $\hat{q}(X)$  is empirically moment matched, then

$$H_p(q(X)) = H_q(q(X)) \quad \hat{H}_p(\hat{q}(X)) = H_{\hat{q}}(\hat{q}(X)) \quad (38)$$

Lemma A.1 allow us to replace the cross entropy terms with simply the entropy of the corresponding Gaussian, allowing for us to have a closed form equation for the JVG estimation

$$I_{\text{JVG}}(X, Y) := H_p(q_{\text{marg}}(X)) - H_p(q_{\text{post}}(X | Y)) \quad (39)$$

$$= H_{q_{\text{marg}}}(q_{\text{marg}}(X)) - H_{q_{\text{post}}}(q_{\text{post}}(X | Y)) = \frac{1}{2} \log |2\pi e \Sigma_x| - \frac{1}{2} \log |2\pi e \Sigma_{x|y}| \quad (40)$$

### A.3 SECTION 3 PROOFS

**Theorem 3.1** (Moment Matched Flow Distribution). *Let  $Z = f(X)$  and  $V = g(Y)$  be diffeomorphisms with Gaussian joint density given by:*

$$q_{Z,V}(Z, V) = \mathcal{N} \left( \begin{bmatrix} z \\ v \end{bmatrix} \middle| \mu := \begin{bmatrix} \mu_z \\ \mu_v \end{bmatrix}, \Sigma := \begin{bmatrix} \Sigma_{zz} & \Sigma_{zv} \\ \Sigma_{vz} & \Sigma_{vv} \end{bmatrix} \right) \quad (41)$$

then the marginal and posterior optimizing the bound in Lemma 2.1 are given by:

$$q_{\text{marg}}(X) = \mathcal{N}(f(X) | \mu_z, \Sigma_{zz}) |\nabla_x f(X)|, \quad q_{\text{post}}(X | Y) = \mathcal{N}(f(X) | \mu_{z|v}, \Sigma_{z|v}) |\nabla_x f(X)|, \\ \text{where } \mu_{z|v} := \mu_z + \Sigma_{zv} \Sigma_{vv}^{-1} (g(Y) - \mu_v), \quad \Sigma_{z|v} := \Sigma_{zz} - \Sigma_{zv} \Sigma_{vv}^{-1} \Sigma_{vz} \quad (42)$$

where the mean and covariance are matched to the moments of the target density  $p(X, Y)$ :

$$\mathbb{E}_p[(f(X), g(Y))^T] = \mu, \quad \text{Cov}_p(f(X), g(Y)) = \Sigma. \quad (43)$$

*Proof.* Let  $Z = f(X)$  and  $V = g(Y)$  be diffeomorphisms with Gaussian joint density, then the joint in  $X$  and  $Y$  can be found by a simple change of variable (Bishop & Nasrabadi, 2006).

$$q_{X,Y}(X, Y) = \mathcal{N} \left( \begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma \right) |\nabla_x f(X)| |\nabla_y g(Y)| \quad (44)$$

We now must show that  $q_{X,Y}(X, Y)$  satisfies the moment matching condition

$$\mathbb{E}_{q_{X|Y}(X|Y)} [T(X, Y)] = \sum_i^k g_i(\eta) \tau_i(Y)$$

$q_{X,Y}(X, Y)$  is in the exponential family with base measure  $h(X, Y) = |\nabla_x f(X)| |\nabla_y g(Y)|$  and  $T(X, Y) = [f(X), g(Y), \text{vec}(f(X)f(X)^T), \text{vec}(f(X)g(Y)^T), \text{vec}(g(Y)g(Y)^T)]$  which are the sufficient statistics we must consider for Theorem A.5. Furthermore, we can derive the posterior

$$\begin{aligned}
q_{X|Y}(X | Y) &= \frac{q_{X,Y}(X, Y)}{q_Y(Y)} = \frac{\mathcal{N}\left(\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) |\nabla_x f(X)| |\nabla_y g(Y)|}{\int \mathcal{N}\left(\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) |\nabla_x f(X)| |\nabla_y g(Y)| dx} \\
&= \frac{\mathcal{N}\left(\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) |\nabla_x f(X)|}{\int \mathcal{N}\left(\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) |\nabla_x f(X)| dx} \\
&= \frac{\mathcal{N}\left(\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) |\nabla_x f(X)|}{\int \mathcal{N}\left(\begin{bmatrix} z \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) dz} \\
&= \frac{\mathcal{N}\left(\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} \middle| \mu, \Sigma\right) |\nabla_x f(X)|}{\mathcal{N}(g(Y) | \mu_V, \Sigma_{VV})} \\
&= \mathcal{N}(f(X) | \mu_{Z|V}, \Sigma_{Z|V}) |\nabla_x f(X)|
\end{aligned}$$

where  $\mu_{Z|V} = \mu_Z + \Sigma_{ZV}\Sigma_{VV}^{-1}(v - \mu_V) = \mu_Z + \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V)$  and  $\Sigma_{Z|V} = \Sigma_{ZZ} - \Sigma_{ZV}\Sigma_{VV}^{-1}\Sigma_{VZ}$ . Using these, we can easily check the linearity condition of each  $T(X, Y)$ .

1.  $f(X)$

$$\begin{aligned}
\mathbb{E}_{q_{X|Y}(X|Y)} [f(X)] &= \int \mathcal{N}(f(X) | \mu_{Z|V}, \Sigma_{Z|V}) |\nabla_x f(X)| f(X) dx \\
&= \int \mathcal{N}(Z | \mu_{Z|V}, \Sigma_{Z|V}) z dz \\
&= \mu_{Z|V} = \mu_Z + \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V)
\end{aligned}$$

2.  $g(Y)$

$$\mathbb{E}_{q_{X|Y}(X|Y)} [g(Y)] = g(Y)$$

3.  $f(X)f(X)^T$

$$\begin{aligned}
\mathbb{E}_{q_{X|Y}(X|Y)} [f(X)f(X)^T] &= \int \mathcal{N}(f(X) | \mu_{Z|V}, \Sigma_{Z|V}) |\nabla_x f(X)| f(X)f(X)^T dx \\
&= \int \mathcal{N}(Z | \mu_{Z|V}, \Sigma_{Z|V}) z z^T dz \\
&= \Sigma_{Z|V} + \mu_{Z|V}\mu_{Z|V}^T \\
&= \Sigma_{ZZ} - \Sigma_{ZV}\Sigma_{VV}^{-1}\Sigma_{ZV}^T + \dots \\
&\quad (\mu_Z + \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V))(\mu_Z + \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V))^T \\
&= \Sigma_{ZZ} - \Sigma_{ZV}\Sigma_{VV}^{-1}\Sigma_{ZV}^T + \mu_Z\mu_Z^T + \dots \\
&\quad \mu_Z(g(Y) - \mu_V^T)\Sigma_{VV}^{-1}\Sigma_{ZV}^T + \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V)\mu_Z + \dots \\
&\quad \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V)(g(Y) - \mu_V)^T \Sigma_{VV}^{-1}\Sigma_{ZV}^T
\end{aligned}$$

4.  $f(X)g(Y)^T$

$$\begin{aligned}
\mathbb{E}_{q_{X|Y}(X|Y)} [f(X)g(Y)^T] &= \mathbb{E}_{q_{X|Y}(X|Y)} [f(X)] g(Y)^T \\
&= (\mu_Z + \Sigma_{ZV}\Sigma_{VV}^{-1}(g(Y) - \mu_V)) g(Y)^T \\
&= (\mu_Z - \Sigma_{ZV}\Sigma_{VV}^{-1}\mu_V) g(Y)^T + \Sigma_{ZV}\Sigma_{VV}^{-1}g(Y)g(Y)^T
\end{aligned}$$



864 5.  $g(Y)g(Y)^T$

$$865 \mathbb{E}_{q_{X|Y}(X|Y)} [g(Y)g(Y)^T] = g(Y)g(Y)^T$$

866  
867 Since these are all linear functions of  $\tau_1(Y) = g(Y)$  and  $\tau_2(Y) = g(Y)g(Y)^T$ , then  $q_{X|Y}(X|Y)$   
868 satisfies the conditions of Theorem A.5 and moment matching the joint  $q_{Z,V}(Z, V) = \mathcal{N}(m, \Sigma)$   
869 yields the optimal  $q_{X|Y}(X|Y)$ . Furthermore, the joint sufficient statistics are

$$870 T(X, Y) = [f(X), g(Y), \text{vec}(f(X)f(X)^T), \text{vec}(f(X)g(Y)^T), \text{vec}(g(Y)g(Y)^T)]$$

871 and the sufficient statistics of a marginal distribution are

$$872 T(X) = [f(X), \text{vec}(f(X)f(X)^T)]$$

873 which are simply the first and third sufficient statistic from the joint. Therefore moment matching the  
874 joint trivially moment matches the marginal, giving the optimal  $q_X(X)$ .  $\square$

875  
876 **Lemma 3.2** (Flow Upper Bound). *Let  $p(X, Y)$  be an arbitrary target distribution and  $q_{X,Y}(X, Y)$   
877 be the distribution of the form in Eqn. (14) with moment matched flow density  $q_{Z,V}(Z, V)$ . Then the  
878 bound in Lemma 2.1 is given by:*

$$879 |I_{m+p} - I| \leq \frac{1}{2} \log |2\pi e \Sigma_{zz}| + \frac{1}{2} \log |2\pi e \Sigma_{z|v}| - 2\mathbb{E}_{p_X} [\log |\nabla_x f(X)|] + C \quad (45)$$

880  
881 *Proof.* The upper bound in Lemma 2.1 is

$$882 |I_{m+p} - I| \leq H_p(q_X(X)) + H_p(q_{X|Y}(q(X|Y))) + C$$

883 This means we need access to the cross-entropy terms of the flow distribution

$$884 H_p(q_X(X)) = - \int p(X) \log q_X(X) dx = - \int p(X) \log q_Z(f(X)) |\nabla_x f(X)| dx \quad (46)$$

$$885 = H_p(q_Z(Z)) - \mathbb{E}_p [\log |\nabla_x f(X)|] = H_{q_Z}(q_Z(Z)) - \mathbb{E}_p [\log |\nabla_x f(X)|] \quad (47)$$

$$886 = \frac{1}{2} \log |2\pi e \Sigma_{zz}| - \mathbb{E}_p [\log |\nabla_x f(X)|] \quad (48)$$

$$887 H_p(q_{X|Y}(X|Y)) = - \int p(X, Y) \log q_{X|Y}(X|Y) dx \quad (49)$$

$$888 = - \int p(X, Y) \log q_Z(f(X) | g(Y)) |\nabla_x f(X)| dx \quad (50)$$

$$889 = H_p(q_{Z|V}(Z|V)) - \mathbb{E}_p [\log |\nabla_x f(X)|] \quad (51)$$

$$890 = H_{q_{Z|V}}(q_{Z|V}(Z|V)) - \mathbb{E}_p [\log |\nabla_x f(X)|] \quad (52)$$

$$891 = \frac{1}{2} \log |2\pi e \Sigma_{z|v}| - \mathbb{E}_p [\log |\nabla_x f(X)|] \quad (53)$$

892 Here, we recognize that  $H_p(q_{Z|V}(Z|V))$  is a moment matched Gaussian, so Lemma A.1 applies to  
893 the cross-entropy. With the cross entropy terms, we can now plug into the upper bound

$$894 |I_{m+p} - I| \leq H_p(q_X(X)) + H_p(q_{X|Y}(q(X|Y))) + C \quad (54)$$

$$895 = \frac{1}{2} \log |2\pi e \Sigma_{zz}| + \frac{1}{2} \log |2\pi e \Sigma_{z|v}| - 2\mathbb{E}_p [\log |\nabla_x f(X)|] + C \quad (55)$$

896  $\square$

897 The equations for the cross entropies, Eqn. (48) and Eqn. (53), allow for us to explicitly write the  
898 form of  $I_{\text{JVF}}$

$$899 I_{\text{JVF}}(X, Y) = H_p(q_X(X)) - H_p(q(X|Y)) \quad (56)$$

$$900 = \frac{1}{2} \log |2\pi e \Sigma_{zz}| - \mathbb{E}_p [\log |\nabla_x f(X)|] - \frac{1}{2} \log |2\pi e \Sigma_{z|v}| + \mathbb{E}_p [\log |\nabla_x f(X)|] \quad (57)$$

$$901 = \frac{1}{2} \log |2\pi e \Sigma_{zz}| - \frac{1}{2} \log |2\pi e \Sigma_{z|v}| \quad (58)$$

902 We see that this is simply the MI of the Gaussian  $q_{Z,V}(Z, V)$  which is the property of MI invariance  
903 under invertible transformations (Czyż et al., 2023).

**Lemma 3.3.** Let  $p(X, Y)$  be any model and  $q_{X,Y}(X, Y)$  of the form in Eqn. (14) with base distribution  $q_{Z,V}(Z, V)$ , then the following bound holds:

$$|I_{JVF} - I| \leq 2H_p(q_{Z,V}(Z, V)) - 2\mathbb{E}_{p_X} [\log |\nabla_x f(X)|] - 2\mathbb{E}_{p_Y} [\log |\nabla_y g(Y)|] + C \quad (59)$$

For a Gaussian base  $q_{Z,V}(Z, V) = \mathcal{N}(\mu, \Sigma)$  the tightest bound of the form Eqn. (18) is given by the moment-matched flow distribution with  $\Sigma = \text{Cov}_p(f(X), g(Y))$  and takes the form:

$$|I_{JVF} - I| \leq \log |2\pi e \Sigma| - 2\mathbb{E}_{p_X} [\log |\nabla_x f(X)|] - 2\mathbb{E}_{p_Y} [\log |\nabla_y g(Y)|] + C \quad (60)$$

where  $C = 2H_p(p(X, Y))$ . This bound is tight when  $q_{X,Y}(X, Y) = p(X, Y)$  almost surely.

*Proof.* Since JVF assumes and parameterizes a joint distribution  $q_{X,Y}(X, Y)$ , we have access to the symmetric definition of MI

$$I_{JVF}(X, Y) = H_p(q_X(X)) - H_p(q_{X|Y}(X | Y)) \quad (61)$$

$$= H_p(q_X(X)) - H_p(q_{X,Y}(X, Y)) + H_p(q_Y(Y)) \quad (62)$$

$$= H_p(q_Y(Y)) - H_p(q_{Y|X}(y | x)) = I_{JVF}(Y, X) \quad (63)$$

For this proof, we simply use  $I_{JVF}(X, Y)$  and  $I_{JVF}(Y, X)$  to differentiate between the two symmetric forms of JVF but do not necessarily hold this notation constant elsewhere. So we will apply Lemma 2.1 twice

$$|I_{JVF}(X, Y) - I(X, Y)| \leq 2|I_{JVF}(X, Y) - I(X, Y)| \quad (64)$$

$$= |I_{JVF}(X, Y) - I(X, Y)| + |I_{JVF}(Y, X) - I(Y, X)| \quad (65)$$

$$\leq H_p(q_X(X)) + H_p(q_{X|Y}(X | Y)) + C_1 \quad (66)$$

$$+ H_p(q_Y(Y)) + H_p(q_{Y|X}(y | x)) + C_2 \quad (67)$$

$$= 2H_p(q_{X,Y}(X, Y)) + C \quad (68)$$

where  $C = C_1 + C_2$ ,  $C_1 = -H_p(p(X)) - H_p(p(X | Y))$ , and  $C_2 = -H_p(p(Y)) - H_p(p(y | x))$  which come directly from bound applied to each term. We also were able to combine  $H_p(q_X(X))$  and  $H_p(q_{Y|X}(y | x))$  into  $H_p(p_{X,Y}(X, Y))$  and likewise for  $H_p(q_Y(Y))$  and  $H_p(q_{X|Y}(X | Y))$ . Since this was two applications of Lemma 2.1, we know that moment matching is optimal for bound and tight when  $q_X(X) = p(X)$ ,  $q_Y(Y) = p(Y)$ ,  $q_{X|Y}(X | Y) = p(X | Y)$ , and  $q_{Y|X}(y | x) = p(y | x)$ . Finally, deriving the explicit form of the bound, we get

$$|I_{JVF}(X, Y) - I(X, Y)| \leq 2H_p(q_{X,Y}(X, Y)) + C \quad (69)$$

$$= 2\mathbb{E}_p [-\log q_{X,Y}(X, Y)] + C \quad (70)$$

$$= 2\mathbb{E}_p [-\log q_{Z,V}(f(X), g(Y)) |\nabla_x f(X)| |\nabla_y g(Y)|] + C \quad (71)$$

$$= 2H_p(q_{Z,V}(Z, V)) - 2\mathbb{E}_p [\log |\nabla_x f(X)| |\nabla_y g(Y)|] + C \quad (72)$$

$$= \log |2\pi e \Sigma| - 2\mathbb{E}_p [\log |\nabla_x f(X)|] - 2\mathbb{E}_p [\log |\nabla_y g(Y)|] + C \quad (73)$$

Where again, we used Lemma A.1 since  $q_{Z,V}(Z, V)$  is a moment matched Gaussian.  $\square$

#### A.4 EXTERNAL SUPPORTING PROOFS

To prove the optimality of moment matching for our estimators in Theorem 2.2 and Theorem 3.1, we relied on the generalized results from Dahlke et al. (2023). We include related results from that paper in this section for convenience of understanding the proofs of Theorem 2.2 and Theorem 3.1.

**Theorem A.2.** Let  $q_m(X)$  be in the exponential family with statistics  $T(X)$ , then for any  $p(X)$ , the optimal  $I_{marg}^*$  is given by moment matching:

$$\mathbb{E}_{q_m(X)} [T(X)] = \mathbb{E}_{p(X)} [T(X)]$$

*Proof.* Since  $H_p(X)$  is constant in  $q_m$  we have,

$$\underset{q_m}{\operatorname{argmin}} H_p(q_m(X)) = \underset{q_m}{\operatorname{argmin}} H_p(q_m(X)) - H_p(X) = \underset{q_m}{\operatorname{argmin}} \text{KL}(p(X) \| q(X)) \quad (74)$$

It is a known result, as proven in (Bishop & Nasrabadi, 2006), that for exponential families,  $\mathbb{E}_{q_m(X)} [T(X)] = \mathbb{E}_{p(X)} [T(X)]$  minimizes  $\text{KL}(p(X) \| q(X))$ .  $\square$

**Lemma A.3.** If  $q_p(X | Y)$  takes the form of

$$q_{X|Y}(X | Y) = q_{X|Y}(X | Y; \eta) = \frac{q_{X,Y}(X, Y; \eta)}{q_Y(Y; \eta)} \quad (75)$$

where  $q_{X,Y}(X, Y; \eta)$  is in the exponential family, then the minimization of  $H_p(q_p(X|Y))$  occurs when

$$\mathbb{E}_{p(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]] = \mathbb{E}_{p(X,Y)} [T(X, Y)] \quad (76)$$

*Proof.* The goal is to minimize  $H_p(q_p(X|Y))$  where  $q_p(X|Y)$  is generated from  $q(X, Y; \eta)$  in the exponential family. We will find the minimizing parameters of this distributions. We appeal to the property of exponential families that  $\frac{\partial}{\partial \eta} A(\eta) = \mathbb{E}_{q(X,Y)} [T(X, Y)]$

$$\frac{\partial}{\partial \eta} (H_p(q_p(X|Y))) = -\frac{\partial}{\partial \eta} \int p(X, Y) \log(q_p(X|Y)) = -\int p(X, Y) \frac{\partial}{\partial \eta} \log\left(\frac{q(X, Y; \eta)}{q(Y; \eta)}\right) \quad (77)$$

$$= -\int p(X, Y) \frac{\partial}{\partial \eta} (\log(h(X, Y)) + \eta^T T(X, Y) - A(\eta) - \log(q(Y; \eta))) \, dx dy \quad (78)$$

$$= -\int p(X, Y) \left( T(X, Y) - \frac{\partial}{\partial \eta} A(\eta) - \frac{\partial}{\partial \eta} \log(q(Y; \eta)) \right) \, dx dy \quad (79)$$

$$= -\mathbb{E}_{p(X,Y)} [T(X, Y)] + \mathbb{E}_{q(X,Y)} [T(X, Y)] + \int p(X, Y) \frac{1}{q(Y; \eta)} \frac{\partial}{\partial \eta} \left( \int h(X', Y) \exp(\eta^T T(X', Y) - A(\eta)) \, dx' \right) \, dx dy \quad (80)$$

$$= -\mathbb{E}_{p(X,Y)} [T(X, Y)] + \mathbb{E}_{q(X,Y)} [T(X, Y)] + \int p(X, Y) \frac{1}{q(Y; \eta)} \left( \int q(X', Y; \eta) \left( T(X', Y) - \frac{\partial}{\partial \eta} A(\eta) \right) \, dx' \right) \, dx dy \quad (81)$$

$$= -\mathbb{E}_{p(X,Y)} [T(X, Y)] + \mathbb{E}_{q(X,Y)} [T(X, Y)] + \int p(X, Y) \left( \int q(X'|Y) (T(X', Y) - \mathbb{E}_{q(X,Y)} [T(X, Y)]) \, dx' \right) \, dx dy \quad (82)$$

$$= -\mathbb{E}_{p(X,Y)} [T(X, Y)] + \mathbb{E}_{p(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]] \quad (83)$$

The zero derivative yields  $\mathbb{E}_{p(X,Y)} [T(X, Y)] = \mathbb{E}_{p(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]]$  which is the stationary condition. It now remains to show that the objective is convex in  $\eta$ . Expanding the form of  $H_p(q(X | Y))$  we have the objective,

$$\min_{\eta} -\mathbb{E}_p [\log(h(X, Y)) + \eta^T T(X, Y) - A(\eta) - \log(q(Y; \eta))] \quad (84)$$

The term  $\eta^T T(X, Y)$  is linear in  $\eta$ . Convexity of  $A(\eta)$  in  $\eta$  is a standard property of the exponential family, however we will show a constructive proof that  $A(\eta) + \log(q(Y; \eta))$  is convex using Hölder's inequality. Let  $\eta = \lambda \eta_1 + (1 - \lambda) \eta_2$  where  $\lambda \in [0, 1]$  and  $\eta_1, \eta_2$  in the convex set of valid exponential family parameters of  $q$  then:

$$A(\eta) + \log(q(Y; \eta)) = A(\eta) + \log\left(\int h(X, Y) \exp(\eta^T T(X, Y) - A(\eta)) \, dx\right) \quad (85)$$

$$= A(\eta) + \log\left(\exp(-A(\eta)) \int h(X, Y) \exp(\eta^T T(X, Y)) \, dx\right) \quad (86)$$

$$= \log\left(\int h(X, Y) \exp(\eta^T T(X, Y)) \, dx\right) \quad (87)$$

$$= \log\left(\int (h(X, Y) \exp(\eta_1^T T(X, Y)))^\lambda (h(X, Y) \exp(\eta_2^T T(X, Y)))^{(1-\lambda)} \, dx\right) \quad (88)$$

$$\leq \lambda \log\left(\int h(X, Y) \exp(\eta_1^T T(X, Y)) \, dx\right) + (1 - \lambda) \log\left(\int h(X, Y) \exp(\eta_2^T T(X, Y)) \, dx\right) \quad (89)$$

$$= \lambda(A(\eta_1) + \log q(Y; \eta_1)) + (1 - \lambda)(A(\eta_2) + \log q(Y; \eta_2)) \quad (90)$$

Thus convexity holds in  $\eta$  and the stationary conditions are globally optimal.

$$\mathbb{E}_{p(X,Y)} [T(X, Y)] = \mathbb{E}_{p(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]] \quad (91)$$

□

**Theorem A.4.** Let  $q(X, Y)$  be in the exponential family with sufficient statistics,  $T(X, Y) = [\tau(X), \tau(Y), \tau(X, Y)]^T$  where  $\tau(X)$  are the sufficient statistics dependent only on  $x$ ,  $\tau(Y)$  only on  $y$ , and  $\tau(X, Y)$  on both. Further, let the posterior expected statistics be a linear combination of marginal statistics as in,

$$\mathbb{E}_{q_p(X|Y)} [T(X, Y)] = \sum_i^k g_i(\eta) \tau_i(Y) \quad (92)$$

where  $\tau_i(Y)$  is the  $i^{\text{th}}$  component of  $\tau(Y)$  and  $g_i(\eta)$  are functions of only the parameter  $\eta$ . Then, the optimal variational distribution,  $q_p$ , for  $I_{\text{post}}$  is defined by joint moment matching:  $\mathbb{E}_{p(X, Y)} [T(X, Y)] = \mathbb{E}_{q(X, Y)} [T(X, Y)]$ .

*Proof.* From Lemma A.3, we know that  $\mathbb{E}_{p(X, Y)} [T(X, Y)] = \mathbb{E}_{p(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]]$  is the optimality condition. Let us now show that the condition in Eqn. (92) implies that joint moment matching satisfies the optimality condition of Eqn. (76)

$$\mathbb{E}_{p(X, Y)} [T(X, Y)] = \mathbb{E}_{q(X, Y)} [T(X, Y)] \quad (93)$$

$$= \mathbb{E}_{q(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]] = \mathbb{E}_{q(Y)} \left[ \sum_i^k g_i(\eta) \tau_i(Y) \right] \quad (94)$$

$$= \sum_i^k g_i(\eta) \mathbb{E}_{q(Y)} [\tau_i(Y)] = \sum_i^k g_i(\eta) \mathbb{E}_{p(Y)} [\tau_i(Y)] \quad (95)$$

$$= \mathbb{E}_{p(Y)} \left[ \sum_i^k g_i(\eta) \tau_i(Y) \right] = \mathbb{E}_{p(Y)} [\mathbb{E}_{q_p(X|Y)} [T(X, Y)]] \quad (96)$$

So with Lemma A.3 and the assumption of the posterior expected statistics being a linear combination of joint statistics (Eqn. (92)) results in  $\mathbb{E}_{p(X, Y)} [T(X, Y)] = \mathbb{E}_{q(X, Y)} [T(X, Y)]$  being the optimal conditions.  $\square$

**Theorem A.5.** Let  $q_{\text{marg}}(X)$  and  $q(X, Y)$  be exponential family distributions where the joint sufficient statistics are  $T(X, Y) = [\tau(X), \tau(Y), \tau(X, Y)]^T$  and  $\tau(X)$  are the sufficient statistics dependent only on  $x$ ,  $\tau(Y)$  only on  $y$ , and  $\tau(X, Y)$  on both. Further, let  $q(X, Y)$  satisfy the linear conditional expectations property

$$\mathbb{E}_{q_{\text{post}}(X|Y)} [T(X, Y)] = \sum_i^k g_i(\eta) \tau_i(Y) \quad (97)$$

where  $\eta$  are the natural parameters of  $q(X, Y)$  and  $g_i(\eta)$  are arbitrary functions dependent only on the natural parameters. Then, moment matching the joint  $q(X, Y)$  and marginal  $q_m(X)$

$$\mathbb{E}_{p(X, Y)} [T(X, Y)] = \mathbb{E}_{q(X, Y)} [T(X, Y)]$$

$$\mathbb{E}_{p(X)} [T(X)] = \mathbb{E}_{q_{\text{marg}}(X)} [T(X)]$$

yield optimal  $q_{\text{post}}(X | Y) \propto q(X, Y)$  and  $q_{\text{marg}}(X)$  that minimize the bound on  $I_{m+p}$  in Lemma 2.1.

*Proof.* We break this down into the two cases of Theorem A.2 and Theorem A.4. Notice that the variational distributions  $q_m(X)$  and  $q_p(X | Y)$  need not share a common joint  $q(X, Y)$ . So, let us use different natural parameters,  $\eta_1$  and  $\eta_2$ , for each (i.e.  $q_m(X) = q(X; \eta_1)$  and  $q_p(X|Y) = q(X|Y; \eta_2)$ ). We optimize the bound in Lemma 2.1 with respect to both natural parameters, beginning with  $\eta_1$ :

$$\frac{\partial}{\partial \eta_1} (-\mathbb{E}_{p(X, Y)} [\log q(X; \eta_1) + \log q(X | Y; \eta_2)] + C) = -\frac{\partial}{\partial \eta_1} \mathbb{E}_{p(X, Y)} [\log q(X; \eta_1)] \quad (98)$$

This is exactly the condition in Theorem A.2 which we know is solved by moment matching the marginal. Likewise, for  $\eta_2$ :

$$\frac{\partial}{\partial \eta_2} (-\mathbb{E}_{p(X, Y)} [\log q(X; \eta_1) + \log q(X | Y; \eta_2)] + C) = -\frac{\partial}{\partial \eta_2} \mathbb{E}_{p(X, Y)} [\log q(X | Y; \eta_2)] \quad (99)$$

The above is the start of the proof for Lemma A.3 in Eqn. (77) and along with Eqn. (92) in Theorem A.4, we get that moment matching the joint finds the optimal  $q_p$ . Therefore, the optimization of Lemma 2.1 simply reduces to moment matching the marginal and the joint.  $\square$

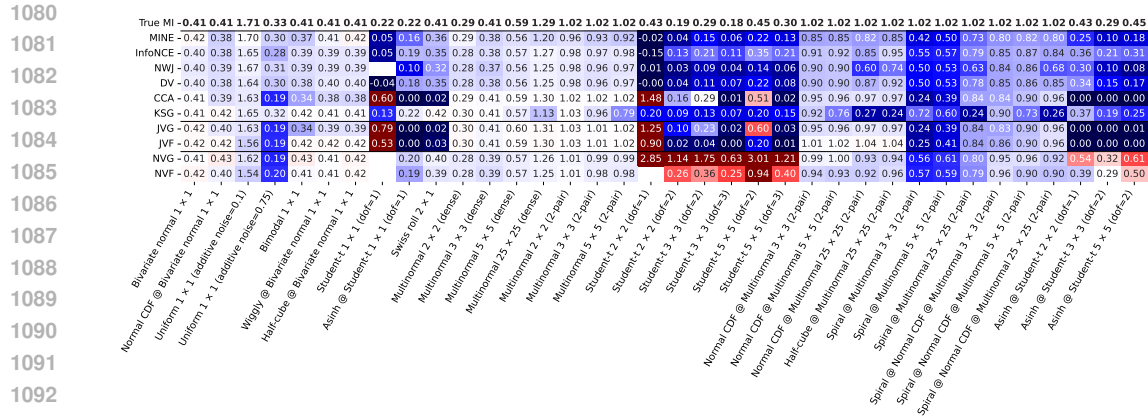


Figure 5: **MI Benchmark** runs against a selection of 36 different scenarios. Blue indicates an underestimate and red indicates an overestimate with the magnitude indicated by saturation. Blank cells indicate that the method produced divergent estimates.

## B EXPERIMENT DETAILS

### B.1 HIGH MUTUAL INFORMATION EXPERIMENT

We consider a synthetic test case where we can create a distribution with high mutual information. Consider a joint Gaussian where  $X, Y \in \mathbb{R}^{15}$

$$p(X, Y) = \mathcal{N} \left( \begin{bmatrix} x \\ y \end{bmatrix} \middle| \mu := \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma := \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right) \quad (100)$$

where  $\mu = 0$  is the zero vector,  $\Sigma_{xx} = \Sigma_{yy} = I$ , and  $\Sigma_{xy} = \Sigma_{yx} = \rho I$ . Then  $I(X, Y) = -\frac{Dim_X}{2} \ln(1 - \rho^2)$  so we are free to choose the value of MI by changing the dimension or by changing the correlation constant. We select  $\rho = .95$  to achieve a MI of  $I(X, Y) = 17.459$ . We train all estimators using a batch size of  $N = 256$  for a total of 3000 steps to ensure all estimators converged.

The critic based methods (DV, MINE, InfoNCE, and NWJ) all use a neural network with two hidden layers of the structure [16, 8] with ReLU activation functions. Czyż et al. (2023) utilized a  $lr = .1$  and batch size of  $N = 256$  which we kept the same.

The normalizing flows utilized in JVG and NVF are Rational Quadratic Splines (Durkan et al., 2019) which learn 128 knots to parameterize the spline. The spline is bounded between  $-8$  and  $8$  where it is linear outside. The knots are learned from a neural network of size [8, 8] with ReLU activation functions. To prevent overfitting of the flows, we utilize dropout with a rate of .2 and L2 Regularization with a factor of  $10^{-5}$ .

The neural parameters used in NVG and NVF are the mean and variance parameterized by two neural networks of size [16, 8] with relu activation functions. The output of the neural network parameterizing the covariance is a lower triangular matrix to approximate the cholesky decomposition of  $\Sigma$  where the diagonal elements have a Softplus activation applied to them to ensure the resulting matrix is semi-positive definite. These parameters did not have dropout performed but did have the same L2 Regularization with a factor of  $10^{-5}$ . NVG, JVF, and NVF all utilized a learning rate of .005 and reduced learning rate by a factor of .1 on test loss if no improvement was seen over 250 testing steps.

### B.2 MUTUAL INFORMATION BENCHMARK

Czyż et al. (2023) created a large selection of MI benchmark tests with ground truth values via MI invariance properties under certain transformations. We will give a brief description of the experiments but in-depth information can be found in their paper.

**Bivariate Normal:** A simple  $1 \times 1$  dimensional Gaussian with  $\rho = .75$ .

**Uniform Margins:** The Gaussian CDF,  $\Psi(\cdot)$ , is applied to Gaussian variables  $X$  and  $Y$ , the resulting

distributions are marginally uniform but not jointly uniform. Denoted as 'Normal CDF @ P'.

**Half-Cube Map:** The transformation  $F(X) = |x|^{3/2}$  is applied to Gaussian variables with the goal of lengthening the tails. Denoted as 'Half-cube @ P'.

**Asinh Mapping:** To shorten tails, the inverse hyperbolic sine function  $\text{asinh}(X) = \log(X + \sqrt{1 + X^2})$  is applied. Denoted as 'Asinh @ P'.

**Wiggly Mapping:** To model non-uniform lengthscales, the mapping  $F(X) = x + \sum a_i \sin(w_i x + \phi_i)$  is applied. Denoted as 'Wiggly @ P'.

**Bimodal Variables:** In the inverse CDF of a two component GMM is applied to  $X$  and  $Y$  to create multimodality.

**Additive Noise:** The random variable  $X \sim \text{Uniform}(0, 1)$  is considered along with noise  $N \sim \text{Uniform}(-\epsilon, \epsilon)$ , then  $Y = X + N$ . Denoted 'Uniform (additive noise= $\epsilon$ )'.

**Swiss Roll Embedding:**  $X, Y \sim \text{Uniform}(0, 1)$  and the Swiss roll embedding is performed on  $X$ .

**Multivariate Normal (2-pair):**  $X$  and  $Y$  are jointly Gaussian distributed with  $\text{Cor}(X_1, Y_1) = \text{Cor}(X_1, Y_2) = .8$  and 0 correlation everywhere else.

**Multivariate Normal (Dense):**  $X$  and  $Y$  are jointly Gaussian distributed with all off diagonal correlations set to .5.

**Multivariate Student:** A multivariate Student-T distribution with degrees of freedom  $\nu$ .

**Spiral:** The transformation  $F(X) = \exp(vA\|x\|^2)x$  where  $A$  is a skew-symmetric matrix which 'mixes' the dimensions and  $v$  is the rate at which the mixing occurs. Denoted 'Spiral @ P'.

The dimensions  $X$  and  $Y$  range from 1 to 25 to give a wide selection of applications. Each estimator is given 1000 samples per dimension with a train test split of 50-50. Each estimator is evaluated on the testing data every 250 gradients steps and early stopping is performed if the estimator test loss has not decreased within the previous 500 steps. All data is pre-processed by centering the data and scaling each dimension by its variance as performed in (Czyż et al., 2023). The same network structures are used as in the Large MI experiment.

### B.3 LOCATION FINDING

In this experiment, we have 2 hidden objects in  $\mathbb{R}^2$  which we wish to learn their locations  $X = \{X_1, X_2\}$ . Each source emits a signal that decays with respect to the inverse square law. From and position,  $d$ , we can observe the superposition of all signals together

$$\mu(X, d) = b + \frac{1}{m + \|X_1 - d\|^2} + \frac{1}{m + \|X_2 - d\|^2} \quad (101)$$

where  $b = .1$  is the background noise and  $m = 10^{-4}$  controls the maximum signal. We can then normally observe the log intensity of the signal from the likelihood

$$\log y | X, d \sim \mathcal{N}(\log(\mu(X, d)), \sigma^2) \quad (102)$$

where  $\sigma = .5$  controls the noise in the observation. We consider two settings for our experiments, *matched prior* and *mismatched prior*. In the case of the matched prior, every method assumes the prior  $p(X) = \mathcal{N}(0, I)$  which corresponds to the true testing prior where testing  $x$  are sampled from. In the mismatched prior setting, every method still assumes that  $p(X) = \mathcal{N}(0, I)$  but the true testing  $x$  are sampled from  $\mathcal{N}(1, I)$ . In practice, we often assume a convenient prior but that prior likely does not correspond well with the real world so an ideal BOED experiment should be able to adapt to this mismatch with minimal loss in information to still provide informative decision across a sequence of decisions. We notice that in this experiment JVG, NVG, and JVF perform significantly worse than NVF. To shed some light on this, we can look at the posterior prediction after making a single observation (Fig.4). The result of this observation increases the belief that sources are located in a radius away from the origin. Neither JVG nor NVG are flexible enough to capture this belief and instead produce Gaussian posteriors. Including flow JVF and NVF have increased flexibility to capture the non-Gaussian belief. We see that the additional degrees of freedom in NVF allows for a higher posterior belief in the vicinity of the true target locations, whereas JVF concentrates more away from the targets.