

# A Bayesian Bootstrap Framework for Mutual Information Neural Estimation: Bridging Classical Mutual Information Learning and Bayesian Nonparametric Learning

**Forough Fazeliasl**

*Department of Mathematical and Statistical Sciences  
University of Alberta, AB, Canada*

*fazeliasl@ualberta.ca*

**Michael Minyi Zhang**

*Department of Statistics and Actuarial Science  
University of Hong Kong Pok Fu Lam, Hong Kong*

*mzhang18@hku.hk*

**Linglong Kong**

*Department of Mathematical and Statistical Sciences  
University of Alberta, AB, Canada*

*lkong@ualberta.ca*

**Bei Jiang**

*Department of Mathematical and Statistical Sciences  
University of Alberta, AB, Canada*

*bei1@ualberta.ca*

**Reviewed on OpenReview:** <https://openreview.net/forum?id=mqGzGKXnFi>

## Abstract

In this work, we introduce a Bayesian bootstrap resampling framework for estimating mutual information (MI) via “mutual information neural estimation” (MINE), making MINE directly applicable in a Bayesian nonparametric learning (BNPL) framework. The resulting estimator shows low variability across batch sizes and high-dimensional settings, as demonstrated through extensive numerical studies. In particular, our proposed bootstrap version yields tighter and lower-variance estimates than the original MINE formulation, both theoretically and empirically. We further demonstrate its practical value in a downstream task by improving VAE-GAN training within BNPL, leading to higher-quality outputs. Beyond enabling MI-based BNPL, the proposed bootstrap estimator also performs competitively against leading frequentist state-of-the-art benchmarks. Overall, our findings establish the first principled framework for Bayesian bootstrap-based MI estimation and highlight its effectiveness as a reliable tool for future BNPL studies.

## 1 Introduction

Estimating and optimizing mutual information (MI) is important in many machine-learning problems, from representation learning to generative modeling. The MI is defined through a Kullback-Leibler (KL) divergence between the joint and marginal distributions to measure the dependency between two variables, making its direct optimization intractable in practice. For this reason, most practical methods use a variational lower bound (VLB) together with an auxiliary function, usually a discriminator that compares joint and marginal samples (Belghazi et al., 2018; Hjelm et al., 2018). Many studies have improved these VLBs to make estimations tighter and more stable, and using these bounds is now a standard approach in MI-based learning (Nguyen et al., 2010; Oord et al., 2018; Dorent et al., 2025).

However, these developments have been designed mainly for standard learning settings. Bayesian nonparametric learning (BNPL) is now growing as a flexible way to model complex data (Fong et al., 2019; Lyddon et al., 2018; 2019; Dellaporta et al., 2022; Fazeli-Asl et al., 2024), with the potential to provide posterior inference under model misspecification; yet none of the existing VLB-MI estimators can be used directly in MI-based BNPL procedures. BNPL relies on posterior sampling induced by randomized objective functions together with posterior bootstrap algorithms, whereas current MI estimators do not account for the underlying distribution being treated as a random probability measure. This gap leaves MI-based BNPL without a practical or compatible lower bound estimator.

In this paper, we fill this gap by introducing a Bayesian bootstrap framework to estimate the VLB of MI. The core idea is to apply Bayesian bootstrap resampling for Mutual Information Neural Estimation (MINE) (Belghazi et al., 2018), replacing uniform sample weights with Dirichlet-resampled weights. Beyond providing this structure that allows VLB optimization to be used in MI-based BNPL, which is the main contribution of this work, we also show that this Bayesian bootstrap produces even a tighter VLB compared with the original MINE. We further compare our estimator numerically with several state-of-the-art VLBs, including NWJ (Nguyen et al., 2010), CPC/InfoNCE (Oord et al., 2018), and JSD-LB (Dorent et al., 2025), and show that it is a strong competitor that offers improved stability and lower variance. To illustrate the practical value of our approach, we incorporate the proposed estimator into a learning problem already formulated within the BNPL framework, namely a VAE-GAN architecture from Fazeli-Asl & Zhang (2023). This provides a natural context to assess how our estimator performs in MI-based optimization and to demonstrate its effectiveness on real data. These results may encourage researchers to consider their learning problems within a Bayesian context.

The rest of the paper is organized as follows. Section 2 reviews standard VLB methods for MI estimation. Section 3 presents our main contribution, the Bayesian bootstrap framework DPMINE, and outlines its basic theoretical properties. Section 4 reports experiments showing that DPMINE provides stable, low-variance MI estimates and can also improve training in a BNP VAE-GAN model when applied within BNPL for data-synthesis purposes. Section 5 concludes with a brief discussion of future work. Additional proofs and experimental details appear in the Appendix.

## 2 Related Work

In this section, we briefly review standard VLBs for estimating MI. We then give a short note on their common use in standard data-generating tasks.

**Mutual Information Estimation.** Estimating MI in machine learning optimization problems is challenging, especially with high dimensional data or small mini-batches. Let  $(\mathbf{U}, \mathbf{V})$  be a pair of random variables with a joint cumulative distribution function (CDF)  $F_{\mathbf{U}\mathbf{V}}$  and marginal CDFs  $F_{\mathbf{U}}$  and  $F_{\mathbf{V}}$ , defined over the product space  $\mathcal{U} \otimes \mathcal{V}$ . The MI between  $\mathbf{U}$  and  $\mathbf{V}$  is given by:

$$\text{MI}(\mathbf{U}, \mathbf{V}) = D_{\text{KL}}(F_{\mathbf{U}\mathbf{V}}, F_{\mathbf{U}} \otimes F_{\mathbf{V}}) = H(\mathbf{U}) - H(\mathbf{U}|\mathbf{V}), \quad (1)$$

where  $D_{\text{KL}}(\cdot, \cdot)$  denotes the KL divergence, and  $H(\cdot)$  and  $H(\cdot|\cdot)$  represent marginal and conditional Shannon entropies, respectively.

Standard approaches to MI estimation either depend on parametric assumptions on the data distribution that rarely hold or rely on nonparametric methods such as  $k$ -nearest neighbors (Al-Labadi et al., 2022). Such nonparametric approaches become unstable as the dimension increases due to the curse of dimensionality (Sugiyama, 2012). These limitations have led to neural VLBs to estimate MI, which scale better but can still struggle when the true MI is large or when training relies on small mini-batch sizes (Dorent et al., 2025).

**Standard Neural VLBs.** The KL-MINE introduced by Belghazi et al. (2018) is one of the most widely used neural VLBs for MI estimation. For instance, Belghazi et al. (2018) introduced KL-MINE, which uses the Donsker-Varadhan (DV) representation (Donsker & Varadhan, 1983) of the KL divergence to form a lower bound of equation 1 as:

$$L_{\gamma}^{\text{KL}}(\mathbf{U}, \mathbf{V}) = \mathbb{E}_{F_{\mathbf{U}\mathbf{V}}}[T_{\gamma}(\mathbf{U}, \mathbf{V})] - \ln \mathbb{E}_{F_{\mathbf{U}} \otimes F_{\mathbf{V}}}[e^{T_{\gamma}(\mathbf{U}, \mathbf{V})}], \quad (2)$$

where  $\{T_\gamma\}_{\gamma \in \Gamma}$  be a set of continuous functions parameterized by a neural network on a compact domain  $\Gamma$  that maps  $\mathcal{U} \otimes \mathcal{V}$  to  $\mathbb{R}$  and  $\ln(\cdot)$  denotes the natural logarithm. In this framework, the auxiliary statistic  $T_\gamma$  serves as the key component that learns to discriminate between samples from  $F_{\mathbf{U}\mathbf{V}}$  and those from  $F_{\mathbf{U}} \otimes F_{\mathbf{V}}$ .

The KL-MINE is then obtained by maximizing VLB given in equation 2 over  $\gamma \in \Gamma$ . In practice, equation 2 is estimated using the empirical cumulative distribution functions (ECDFs),  $F_{\mathbf{U}_{1:n}} := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{U}_i}$  and  $F_{\mathbf{V}_{1:n}} := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{V}_i}$ :

$$L_\gamma^{\text{KL}}(\mathbf{U}_{1:n}, \mathbf{V}_{1:n}) = \sum_{\ell=1}^n \frac{1}{n} T_\gamma(\mathbf{U}_\ell, \mathbf{V}_\ell) - \ln \sum_{\ell=1}^n \frac{1}{n} e^{T_\gamma(\mathbf{U}_\ell, \mathbf{V}_{\pi(\ell)})}, \quad (3)$$

where  $\{\pi(1:n)\}$  denotes a random permutation of  $\{1:n\}$ , a standard technique to empirically approximate  $\mathbb{E}_{F_{\mathbf{X}} \otimes F_{\mathbf{Y}}}(\cdot)$  in equation 2 (Belghazi et al., 2018; Hjelm et al., 2018). However, MINE is prone to high variances (Liao et al., 2020).

Another variant of MINE is based on a variational lower bound for the Jensen-Shannon (JS) divergence (Nowozin et al., 2016) between  $F_{\mathbf{U}\mathbf{V}}$  and  $F_{\mathbf{U}} \otimes F_{\mathbf{V}}$  (Hjelm et al., 2018; Jones et al., 2023). It is defined as

$$L_\gamma^{\text{JS}}(\mathbf{U}, \mathbf{V}) = \mathbb{E}_{F_{\mathbf{U}\mathbf{V}}}[-\zeta(-T_\gamma(\mathbf{U}, \mathbf{V}))] - \mathbb{E}_{F_{\mathbf{U}} \otimes F_{\mathbf{V}}}[\zeta(T_\gamma(\mathbf{U}, \mathbf{V}))], \quad (4)$$

where  $\zeta(\cdot) = \ln(1 + \exp(\cdot))$  denotes the softplus function. Although JS-MINE yields lower variance and more stable results, its connection to the KL-based definition of MI remains unclear (Dorent et al., 2025).

Dorent et al. (2025) addresses this concern by providing a new VLB on the KL divergence expressed as a function of the JS divergence, clarifying why maximizing the JS divergence is equivalent to maximizing a lower bound on MI. Their VLB also overcomes well-known drawbacks of NWJ (Nguyen et al., 2010), which often has high variance, and CPC/InfoNCE (Oord et al., 2018), whose estimates are restricted by the batch size.

**Alternative Neural Bounds.** There are additional neural variational bounds for MI estimation, including VNMC and VPCE (Ivanova et al., 2024), which provide upper and lower bounds, respectively. These methods rely on additional modeling assumptions, including the availability of at least one distribution in closed form, and estimate MI using nested Monte Carlo procedures based on Bayes’ rule. Consequently, they are commonly used in modeling settings that differ from distribution-free neural VLB estimators such as MINE, NWJ, CPC, or JS-based bounds, which operate directly on samples without imposing analytic distributional forms. We therefore view VNMC and VPCE as complementary approaches that target MI estimation under different assumptions, rather than methods intended for the same training setting considered in this paper.

**Two-Stage Estimators.** Several other neural MI estimators in the literature use a two-stage design, where a separate model is trained to estimate MI before it is used in an optimization task, see for example, Letizia et al. (2024); Shalev et al. (2022); Liao et al. (2020); Tsai et al. (2020); Song & Ermon (2019). These approaches can provide stable estimates, but they must be retrained whenever the joint distribution changes, which makes them impractical for learning settings where the data features change during training. These methods operate differently from single-stage estimators and cannot be used directly during training.

**Applying Standard VLBs to Improve GANs.** All the estimators discussed are often used within learning problems where MI guides the optimization. One example is to examine their behavior within the GAN family, which provides a natural setting for incorporating MI estimation into a downstream task. Models in this family, including GANs and VAE-GANs (Goodfellow et al., 2014; Larsen et al., 2016), use a generator to create samples, a discriminator to judge their quality, and in VAE-GANs an encoder that maps data to latent codes. MI has been introduced into these models to better align latent variables with generated outputs. InfoGAN (Chen et al., 2016) encouraged part of the latent code to remain informative about the output by adding an MI term to the objective. Building on this idea, Belghazi et al. (2018) used MINE to increase the MI between data and their encoded representations in a BiGAN-style model (Donahue et al., 2016; Dumoulin et al., 2016), improving data coverage and reducing mode collapse. These examples show that MI objectives can help GAN variants learn richer representations and produce more diverse samples.

### 3 Methodology

In this section, we present the main contribution of this work, a Bayesian bootstrap resampling for MINE. We begin with a brief background on BNPL and then present the resulting bootstrap algorithms. Finally, we illustrate how this framework can be applied in a downstream task.

#### 3.1 Dirichlet Process for the Modeling Data Distribution

The DP is an infinite generalization of the Dirichlet distribution that is considered on the sample space denoted as  $\mathfrak{X}$ , which possesses a  $\sigma$ -algebra  $\mathcal{A}$  comprising subsets of  $\mathfrak{X}$  (Ferguson, 1973).  $F$  follows a DP with parameters  $(a, H)$  with the notation  $F \sim DP(a, H)$ , if for any measurable partition  $A_1, \dots, A_k$  of  $\mathfrak{X}$  with  $k \geq 2$ , the joint distribution of the vector  $(F(A_1), \dots, F(A_k))$  follows a Dirichlet distribution characterized by parameters  $(aH(A_1), \dots, aH(A_k))$ . Moreover, it is assumed that  $H(A_j) = 0$  implies  $F(A_j) = 0$  with probability one. The base measure  $H$  captures the prior knowledge regarding the data distribution, while  $a$  signifies the strength or intensity of this knowledge.

As a conjugate prior, the Dirichlet process yields a posterior that is also a DP. Specifically, for  $n$  independent and identically distributed draws  $(\mathbf{X}_{1:n} \in \mathbb{R}^d)$  from the random probability measure  $F$ , the posterior distribution of  $F$  is

$$F^{\text{Pos}} := (F \mid \mathbf{X}_{1:n}) \sim DP(a + n, H_n^*). \quad (5)$$

The updated base measure is

$$H_n^* = \frac{a}{a+n} H + \frac{n}{a+n} F_{\mathbf{X}_{1:n}}, \quad (6)$$

where  $F_{\mathbf{X}_{1:n}}$  denotes the empirical cumulative distribution function of the sample  $\mathbf{X}_{1:n}$ .

It can be shown that a Dirichlet process has an infinite series representation and is therefore a discrete probability measure. Although the stick-breaking representation is a commonly employed series representation for DP inference (Sethuraman, 1994), it lacks the necessary normalization terms to convert it into a probability measure (Zarepour & Al-Labadi, 2012). Additionally, simulating from an infinite series is only feasible through using a random truncation approach to handle the terms within the series. To address these limitations, Ishwaran & Zarepour (2002) introduced an approximation of the DP in the form of a finite series, which allows for convenient simulation. In the context of posterior inference, this approximation is given by

$$F_N^{\text{Pos}} := \sum_{i=1}^N J_i^{\text{Pos}} \delta_{\mathbf{X}_i^{\text{Pos}}}, \quad (7)$$

where  $(J_{1:N}^{\text{Pos}}) \sim \text{Dirichlet}((a+n)/N, \dots, (a+n)/N)$ ,  $(\mathbf{X}_{1:N}^{\text{Pos}}) \stackrel{\text{iid}}{\sim} H_n^*$ , and  $\delta_{\mathbf{X}^{\text{Pos}}}$  is the Dirac delta measure. In this study, the variables  $J_i^{\text{Pos}}$  and  $\mathbf{X}_i^{\text{Pos}}$  represent the DP's weight and location, respectively. The sequence  $(F_N^{\text{Pos}})_{N \geq 1}$  converges in distribution to  $F^{\text{Pos}}$ , where  $F_N^{\text{Pos}}$  and  $F^{\text{Pos}}$  are random values in  $M_1(\mathbb{R}^d)$ , the space of probability measures on  $\mathbb{R}^d$  endowed with the topology of weak convergence (Ishwaran & Zarepour, 2002).

Since the approximation in equation 7 is finite by construction, the remaining practical question concerns the choice of the truncation level  $N$ . Rather than fixing  $N$  in advance, an adaptive truncation rule is employed to determine the number of retained atoms so that additional terms contribute negligibly to the total mass. Following Zarepour & Al-Labadi (2012), let  $V_i \stackrel{\text{iid}}{\sim} \text{Gamma}\left(\frac{a+n}{j}, 1\right)$  for  $i = 1, \dots, j$ , and define the truncation level as

$$N = \inf \left\{ j : J_j^{\text{Pos}} = \frac{V_j}{\sum_{i=1}^j V_i} < \epsilon \right\}, \quad (8)$$

for a prescribed threshold  $\epsilon \in (0, 1)$ . This adaptive selection of  $N$  ensures that the finite approximation  $F_N^{\text{Pos}}$  captures the dominant posterior mass while remaining computationally tractable, and it is used throughout the remainder of the paper.

### 3.2 BNPL: Posterior Bootstrap Framework

Let  $(\mathbf{X}_{1:n})$  be IID random vectors in  $\mathbb{R}^d$  with common distribution  $F \in \mathcal{P}(\mathbb{R}^d)$ , where  $\mathcal{P}(\mathbb{R}^d)$  denotes the space of Borel probability measures on  $\mathbb{R}^d$ . Suppose we are interested in a parameter  $\psi \in \Psi \subseteq \mathbb{R}^p$  indexing a collection of parametric models

$$\mathcal{F}_\Psi = \{F_\psi(\mathbf{x}) : \psi \in \Psi\}, \quad (9)$$

for example, the distribution induced by a neural network family parameterized by  $\psi$ .

Letting  $F \notin \mathcal{F}_\Psi$ , which allows for the possibility of model misspecification, the nonparametric learning (NPL) framework defines the parameter of interest as a functional  $\dot{\psi} : \mathcal{P}(\mathbb{R}^d) \rightarrow \Psi$  of the data distribution, given by

$$\dot{\psi}(F) = \arg \min_{\Psi} \int L(\mathbf{x}, \psi) dF(\mathbf{x}), \quad (10)$$

where  $L(\mathbf{x}, \psi)$  is a user-chosen loss that determines the target statistic (Fong et al., 2019; Lyddon et al., 2018; 2019).

Since  $F$  is unknown, it is modeled as a random probability measure with a DP prior  $F \sim \text{DP}(\alpha, H)$ . The core idea of BNPL is that the posterior random measure  $F^{\text{Pos}}$  induces a posterior distribution on the parameter of interest through the functional  $\hat{\psi} : \mathcal{P}(\mathbb{R}^d) \rightarrow \Psi$ , given by

$$\hat{\psi}(F^{\text{Pos}}) = \arg \min_{\Psi} \int L(\mathbf{x}, \psi) dF^{\text{Pos}}(\mathbf{x}). \quad (11)$$

Because the DP is discrete and we use the finite approximation in equation 7, the expression in equation 11 admits the practical approximation

$$\hat{\psi}(F^{\text{Pos}}) \approx \arg \min_{\Psi} \sum_{i=1}^N J_i^{\text{Pos}} L(\mathbf{x}_i^{\text{Pos}}, \psi), \quad (12)$$

which corresponds to the Bayesian bootstrap procedure described in Fong et al. (2019, Algorithm 2), with the difference that Fong et al. (2019) employs a truncation of the stick-breaking series representation to obtain this finite approximation.

Dellaporta et al. (2022) and Fazeli-Asl et al. (2024) also independently introduced posterior bootstrap methods based on two different DP approximations for learning the data distribution. In both works, learning is carried out by minimizing the maximum mean discrepancy between the posterior distribution  $F^{\text{Pos}}$  and the distribution induced by a generator modeled as a parameterized neural network.

Overall, these posterior bootstrap methods operate within a robust learning framework and are designed to handle model misspecification.

### 3.3 DPMINE: Posterior Bootstrap Framework for MINE

We now present our main contribution by extending the Bayesian bootstrap idea of Fong et al. (2019) to the MINE. Although our setting differs from Fong et al. (2019), the connection is direct: the posterior on the MINE parameters is induced by placing a DP posterior on the unknown data distribution. This construction makes the original MINE applicable for MI-based BNPL. We first develop this approach for the KL-MINE objective and then show that the same steps apply to the JS-MINE case.

**KL-based representation:** Let  $(\mathbf{X}_{1:n})$  be  $n$  IID random variables with  $\mathbf{X} \in \mathbb{R}^d \sim F$ . For any continuous mappings  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ ,  $i = 1, 2$ , we are interested in estimating the MI between  $f_1(\mathbf{X})$  and  $f_2(\mathbf{X})$  using the KL-MINE objective from the BNPL perspective. In this setting, we define the parameter of interest as the functional  $\dot{\gamma} : \mathcal{P}(\mathbb{R}^d) \rightarrow \Gamma$ , given by

$$\dot{\gamma}(F) = \arg \min_{\Gamma} (-L_{\dot{\gamma}}^{\text{KL}}(f_1(\mathbf{X}), f_2(\mathbf{X}))). \quad (13)$$

**Algorithm 1:** MINE Posterior Bootstrap

**Input:** training data  $\mathbf{X}_{1:n}$ , mappings  $f_1, f_2$ , concentration parameter  $a$ , base measure  $H$ , truncation level  $N$ , number of bootstrap samples  $B$ , critic network  $T_\gamma$

1 **for**  $i = 1$  **to**  $B$  **do**

2     Select the  $i$ -th mini-batch  $\mathbf{X}_{\text{mb}}^{(i)}$  of size  $n_{\text{mb}}$ ;

3     Define  $H_{n_{\text{mb}}}^* = \frac{a}{a+n_{\text{mb}}} H + \frac{n_{\text{mb}}}{a+n_{\text{mb}}} F_{\mathbf{X}_{\text{mb}}^{(i)}}$ ;

4     Draw posterior samples  $\mathbf{X}_{1:N}^{\text{Pos}} \stackrel{\text{iid}}{\sim} H_{n_{\text{mb}}}^*$ ;

5     Draw Dirichlet weights  $(J_{1:N}^{\text{Pos}}) \sim \text{Dirichlet}((a+n)/N, \dots, (a+n)/N)$ ;

6     Compute  $L_\gamma^{\text{DP}}(\cdot, \cdot)$  using equation 15 (DPKL) or equation 17 (DPJS) with the critic  $T_\gamma$ , posterior samples  $\mathbf{X}_{1:N}^{\text{Pos}}$ , and Dirichlet weights  $J_{1:N}^{\text{Pos}}$ ;

7     Update the critic parameter

$$\hat{\gamma}^{(i)} = \arg \min_{\gamma \in \Gamma} \left\{ -L_\gamma^{\text{DP}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right\}$$

8     Compute the MI estimate

$$\text{MI}^{\text{DP}(i)} = L_{\hat{\gamma}^{(i)}}^{\text{DP}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}}))$$

9 **return** Posterior bootstrap sample  $\text{MI}^{\text{DP}(1:B)}$

Given a posterior draw  $F^{\text{Pos}} \sim \text{DP}(a+n, H_n^*)$ , the posterior for  $\dot{\gamma}(F)$  is obtained through the functional  $\hat{\gamma} : \mathcal{P}(\mathbb{R}^d) \rightarrow \Gamma$ , defined by

$$\hat{\gamma}(F^{\text{Pos}}) = \arg \min_{\gamma \in \Gamma} \left\{ -L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right\}, \quad (14)$$

where

$$L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) := \sum_{\ell=1}^N J_\ell^{\text{Pos}} T_\gamma(f_1(\mathbf{X}_\ell^{\text{Pos}}), f_2(\mathbf{X}_\ell^{\text{Pos}})) - \ln \sum_{\ell=1}^N J_\ell^{\text{Pos}} e^{T_\gamma(f_1(\mathbf{X}_\ell^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))}. \quad (15)$$

Finally, equation 16 represents the posterior belief about  $\text{MI}^{\text{KL}}(f_1(\mathbf{X}), f_2(\mathbf{X}))$ .

$$\text{MI}^{\text{DPKL}}(f_1(\mathbf{X}^{\text{Pos}}), f_2(\mathbf{X}^{\text{Pos}})) = L_{\hat{\gamma}}^{\text{DPKL}}(f_1(\mathbf{X}^{\text{Pos}}), f_2(\mathbf{X}^{\text{Pos}})). \quad (16)$$

Algorithm 1 provides a clear illustration of these steps.

**JS-based representation:** All steps used in the Bayesian bootstrap construction for KL-MINE extend simply to the JS-MINE objective. The DP-based analogue of the JS variational bound is given by

$$L_\gamma^{\text{DPJS}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) := \sum_{\ell=1}^N J_\ell^{\text{Pos}} \left[ -\zeta(-T_\gamma(f_1(\mathbf{X}_\ell^{\text{Pos}}), f_2(\mathbf{X}_\ell^{\text{Pos}}))) - \zeta(T_\gamma(f_1(\mathbf{X}_\ell^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))) \right], \quad (17)$$

The following theorem states that the proposed Bayesian VLB estimator is asymptotically no smaller, on average, than the standard VLB used in MINE. This results in a tighter VLB on the true MI and allows for a more accurate estimation of MI by maximizing the tighter bound over a compact domain  $\Gamma$ .

**Theorem 1** (Limiting expectation). *Considering DP posterior representations defined in equation 15 and equation 17. Let both  $T_\gamma$  and  $e^{T_\gamma}$  are  $M$ -bounded (i.e.,  $|T_\gamma|, |e^{T_\gamma}| \leq M$ ). Given the DP posterior approximation in equation 7, for any  $N > 0$ , we have,*

$$i. \quad \lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{\text{Pos}}} \left( L_{\hat{\gamma}}^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right) \geq L_{\hat{\gamma}}^{\text{KL}}(f_1(\mathbf{X}), f_2(\mathbf{X})),$$

$$ii. \lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{Pos}} \left( L_{\hat{\gamma}}^{DPJS}(f_1(\mathbf{X}_{1:N}^{Pos}), f_2(\mathbf{X}_{1:N}^{Pos})) \right) \geq L_{\hat{\gamma}}^{JS}(f_1(\mathbf{X}), f_2(\mathbf{X})).$$

In Theorem 1, the expectation is taken with respect to the nonparametric posterior on the data-generating process. As a result, the randomness in  $\hat{\gamma}$ , which arises from the posterior bootstrap, is automatically averaged out. Equivalently,  $\hat{\gamma}$  is treated as fixed conditional on a posterior draw, and the theorem characterises the resulting average behaviour. This treatment is standard in Bayesian nonparametric learning and posterior bootstrap analyses (Dellaporta et al., 2022).

### 3.4 Application: Integrating DPMINE into VAE-GAN Models within BNPL

The BNPWMMMD-GAN framework proposed by Fazeli-Asl & Zhang (2023) leverages representation learning within BNPL by placing a DP prior  $F \sim \text{DP}(a, H)$  to generate synthetic samples.

The model follows a standard VAE-GAN architecture, consisting of an encoder  $E_{\eta} : \mathbb{R}^d \rightarrow \mathbb{R}^p$ ,  $\eta \in \mathcal{H}$ , that maps data into a latent representation, a generator (decoder)  $G_{\omega} : \mathbb{R}^p \rightarrow \mathbb{R}^d$ ,  $\omega \in \Omega$ , that reconstructs and generates samples from latent codes, a code generator  $CG_{\omega'} : \mathbb{R}^q \rightarrow \mathbb{R}^p$ ,  $\omega' \in \Omega'$ , with  $q < p$  that injects additional latent noise to improve coverage of underrepresented regions of the latent space, and a discriminator  $D_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\theta \in \Theta$ , that distinguishes real from generated samples. These components are trained jointly using Wasserstein and MMD-based losses evaluated on posterior-weighted samples generated by DP approximation equation 7. Further details are provided in Appendix C.

Training is implemented by estimating the network parameters as functionals of  $F^{\text{Pos}}$ . Let

$$\Xi = \Omega \times \mathcal{H} \times \Theta \times \Omega'$$

denote the joint parameter space. The baseline BNPWMMMD-GAN objective is defined as

$$(\hat{\omega}, \hat{\eta}, \hat{\theta}, \hat{\omega}') (F^{\text{Pos}}) = \arg \min_{(\omega, \eta, \theta, \omega') \in \Xi} \left\{ L(G_{\omega}, E_{\eta}, F_N^{\text{Pos}}) + L(D_{\theta}, F_N^{\text{Pos}}) + L(CG_{\omega'}, F_N^{\text{Pos}}) \right\}. \quad (18)$$

To improve information extraction during encoding and preserve meaningful structure during decoding, this objective is refined by two regularization terms estimated via the proposed DPMINE framework, yielding

$$(\hat{\omega}, \hat{\eta}, \hat{\theta}, \hat{\omega}') (F^{\text{Pos}}) = \arg \min_{(\omega, \eta, \theta, \omega') \in \Xi} \left\{ L(G_{\omega}, E_{\eta}, F_N^{\text{Pos}}) + L(D_{\theta}, F_N^{\text{Pos}}) + L(CG_{\omega'}, F_N^{\text{Pos}}) \right. \\ \left. - L_{\hat{\gamma}_1}^{\text{DP}}(\mathbf{X}_{1:N}^{\text{Pos}}, E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}})) \right. \quad (19)$$

$$\left. - L_{\hat{\gamma}_2}^{\text{DP}}(E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}}), G_{\omega}(E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}}))) \right\}. \quad (20)$$

The regularization terms in equation 19 and equation 20, computed using the corresponding critics  $T_{\hat{\gamma}_1}$  and  $T_{\hat{\gamma}_2}$ , estimate  $\text{MI}(X, E_{\eta}(X))$  and  $\text{MI}(E_{\eta}(X), G_{\omega}(E_{\eta}(X)))$ , respectively. The associated critic parameters are estimated as posterior functionals,

$$\hat{\gamma}_1(F^{\text{Pos}}) = \arg \min_{\gamma_1 \in \Gamma_1} \left\{ -L_{\gamma_1}^{\text{DP}}(\mathbf{X}_{1:N}^{\text{Pos}}, E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}})) \right\}, \quad (21)$$

$$\hat{\gamma}_2(F^{\text{Pos}}) = \arg \min_{\gamma_2 \in \Gamma_2} \left\{ -L_{\gamma_2}^{\text{DP}}(E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}}), G_{\omega}(E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}}))) \right\}. \quad (22)$$

Overall, this formulation integrates DPMINE terms into a VAE-GAN training procedure, while maintaining the BNPL interpretation.

Using the notation of Section 3.3, the term in equation 21 corresponds to

$$f_1(\mathbf{X}_{1:N}^{\text{Pos}}) = \mathbf{X}_{1:N}^{\text{Pos}}, \quad f_2(\mathbf{X}_{1:N}^{\text{Pos}}) = E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}}), \quad (23)$$

Similarly, the term in equation 22 corresponds to

$$f_1(\mathbf{X}_{1:N}^{\text{Pos}}) = E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}}), \quad f_2(\mathbf{X}_{1:N}^{\text{Pos}}) = G_{\omega}(E_{\eta}(\mathbf{X}_{1:N}^{\text{Pos}})), \quad (24)$$

when  $G_{\omega}$  acts as the decoder.

## 4 Experimental Results

In this section, the procedure is evaluated on both simulated and real datasets. All experiments are conducted using mini-batches, with the batch size chosen appropriately for each setting. To place a DP prior on the data distribution  $F$ , its hyperparameters must first be specified. Unless an experiment is designed to study their specific influence, the base measure is taken to be a multivariate normal distribution estimated from each mini-batch. For a mini-batch  $\mathbf{X}_{\text{mb}}^{(i)}$ ,  $i \in \{1, \dots, B\}$ , the sample mean and covariance are

$$\bar{\mathbf{X}}_{\text{mb}} = \frac{1}{B} \sum_{i=1}^B \mathbf{X}_{\text{mb}}^{(i)}, \quad S_{\mathbf{X}_{\text{mb}}} = \sum_{i=1}^B (\mathbf{X}_{\text{mb}}^{(i)} - \bar{\mathbf{X}}_{\text{mb}})(\mathbf{X}_{\text{mb}}^{(i)} - \bar{\mathbf{X}}_{\text{mb}})^T, \quad (25)$$

and the base measure is defined as

$$H := N(\bar{\mathbf{X}}_{\text{mb}}, S_{\mathbf{X}_{\text{mb}}}).$$

The concentration parameter is then selected as the *maximum a posteriori* (MAP) estimate of  $a$ , obtained by maximizing the log-likelihood of  $F^{\text{Pos}}$  with respect to  $a$ . Full details of this computation are provided in Fazeli-Asl & Zhang (2023, Section 4.2). Additionally, the value of  $N$  in the DP approximation of equation 7 is chosen as prescribed in Section 3.1, using  $\epsilon = 10^{-4}$ .

We start by presenting simulation studies that evaluate the KL-based and JS-based versions of our Bayesian bootstrap framework for the MI estimation. We then show how the better-performing estimator can be used to refine the training of VAE-GANs within the BNPL framework.

### 4.1 MI Estimation using Bayesian Bootstrap Resampling

We follow the staircase protocol of Letizia et al. (2024), adapted to our experimental setting, to construct the toy examples. In our design, we let

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}) \in \mathbb{R}^{2k}, \quad \mathbf{U}, \mathbf{V} \in \mathbb{R}^k,$$

and consider estimating  $\text{MI}(\mathbf{U}, \mathbf{V})$  for  $k \in \{1, 2, 3, 4, 5\}$ . Increasing  $k$  produces the successive staircase levels: the first stage uses  $k = 1$ , the second uses  $k = 2$ , and this continues until the fifth stage uses  $k = 5$ . This yields five distinct increases in the true MI, giving a clear benchmark for comparing estimator performance as more dimensions are added.

**Binary-Gaussian Example.** For each value of  $k$ , we generate a binary-Gaussian example by defining

$$\mathbf{U} = \text{sign}(\mathbf{Z}), \quad \mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_k), \quad \mathbf{V} = \mathbf{U} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, 0.2\mathbf{I}_k).$$

In the notation of Section 3.3, where the MI was written as  $\text{MI}(f_1(\mathbf{X}), f_2(\mathbf{X}))$ , this setting corresponds to the choice  $f_1(\mathbf{X}) = \mathbf{U}$  and  $f_2(\mathbf{X}) = \mathbf{V}$ .

Finally, each stairstep provides 4000 training steps, where a step means one gradient update of the neural estimation tasks. At every step, we draw a mini-batch from the distribution associated with the current value of  $k$ . Additional toy examples following the same protocol are provided in the appendix.

#### 4.1.1 DPMINE versus MINE

As the main contribution of this paper is to apply Bayesian bootstrap resampling to the mini-batches used in MINE, we begin by examining whether our Bayesian version can outperform the original estimator. Figure 1 shows our results, where we use a mini-batch size of 16. The solid curves display the running averages of each estimator; for BNP-MINE, this running average corresponds to the posterior expectation induced by the Bayesian bootstrap and represents the expected behavior of the MI posterior, which converges to the true MI in the KL-based setting. The shaded regions show the raw mini-batch variability. In the KL setting, our BNP-MINE estimator gives a tighter and more stable lower bound compared to FNP-MINE, which numerically supports the theoretical results in Theorem 1(i). The FNP curves vary widely and can rise

above the true MI, reflecting the instability often seen in MINE training. In contrast, the BNP-KL curve has less variance and remains close to the true MI, leading to a more consistent estimate.

The BNP-JS also shows reduced variance and avoids the large jumps observed in the FNP-JS, even though it is not as tight as the KL version. Overall, the BNP framework leads to more stable and reliable MI estimates, with the KL formulation giving the strongest improvement.

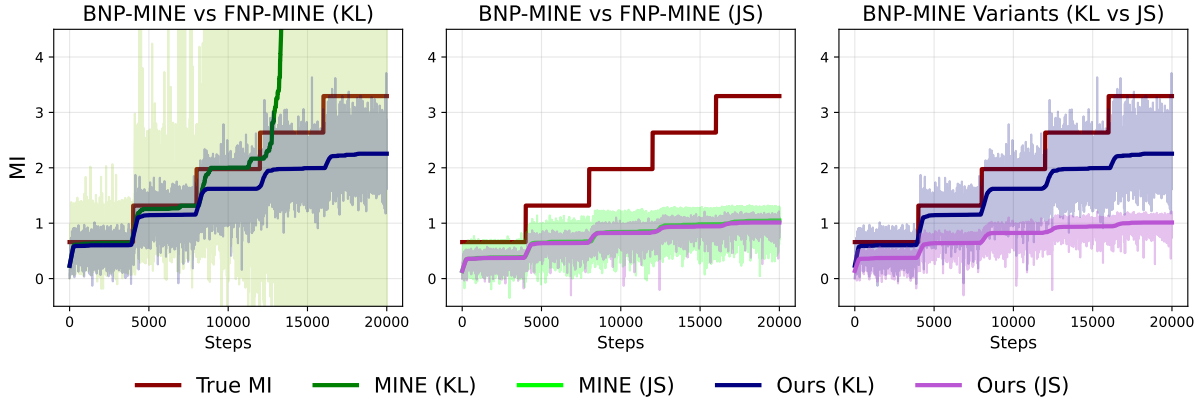


Figure 1: Staircase scheme for the binary-Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE and FNP-MINE (KL and JS) are compared using mini-batch size 16, with  $F \sim \text{DP}(a_{\text{MAP}}, H)$  and  $H := N(\bar{\mathbf{X}}_{\text{mb}}, S_{\mathbf{X}_{\text{mb}}})$ , across dimensions  $k = 1, \dots, 5$ .

#### 4.1.2 Sensitivity to different mini-batch sizes

To examine how different mini-batch sizes affect the estimator, Figure 2 shows BNP-MINE under several choices of batch size. Across all settings, the estimates remain close to the true MI, and larger batches tend to produce smoother and slightly more accurate curves. Overall, the method remains stable, and the effect of the batch size is modest.

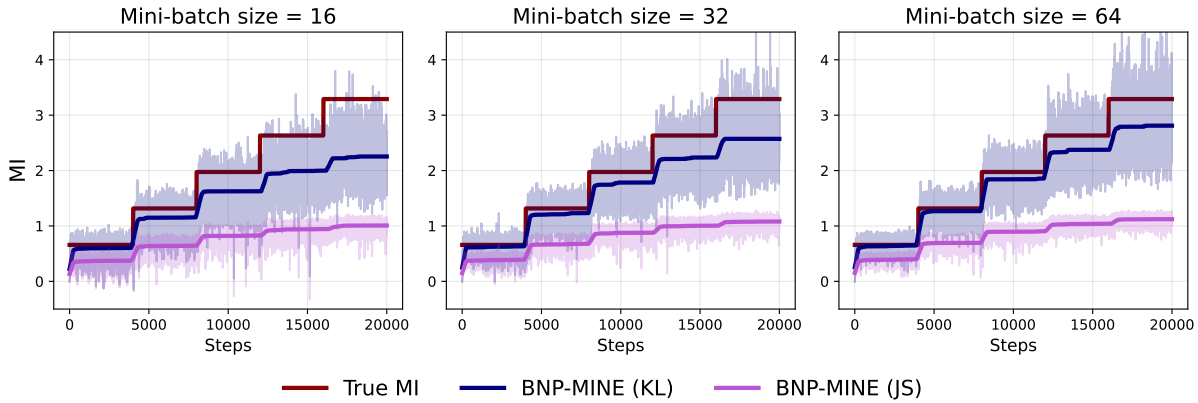


Figure 2: Staircase scheme for the binary-Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE is evaluated across different mini-batch sizes, with  $F \sim \text{DP}(a_{\text{MAP}}, H)$  and  $H := N(\bar{\mathbf{X}}_{\text{mb}}, S_{\mathbf{X}_{\text{mb}}})$ , across dimensions  $k = 1, \dots, 5$ .

#### 4.1.3 Sensitivity to different choices of DP hyperparameters

To examine how the DP hyperparameters affect our procedure, we study a range of concentration values  $a \in \{0.005, 25, 250, 2500\}$  as well as  $a_{\text{MAP}}$ , under two choices of the centering measure  $H$ . In both sets of experiments,  $a_{\text{MAP}}$  produces stable estimates in the sense that the MI remains close to the true value.

Figure 3(a) shows that when  $H$  is taken to be a standard normal distribution, increasing  $a$  forces the prior to dominate the data, and the estimator gradually converges to the MI of a standard normal model, which is zero. In contrast, Figure 3(b) shows that when  $H$  is defined using the sample mean and sample covariance of each mini-batch, the prior remains aligned with the data, and the estimator stays accurate even for large  $a$ . These results, which follow directly from the posterior base measure in equation 6, show that reasonable choices of  $H$  keep the method reliable even for large  $a$ . Although strongly misspecified choices of  $H$  can push the estimates toward the wrong target, the method remains robust to the choice of  $H$  when  $a$  is sufficiently small.

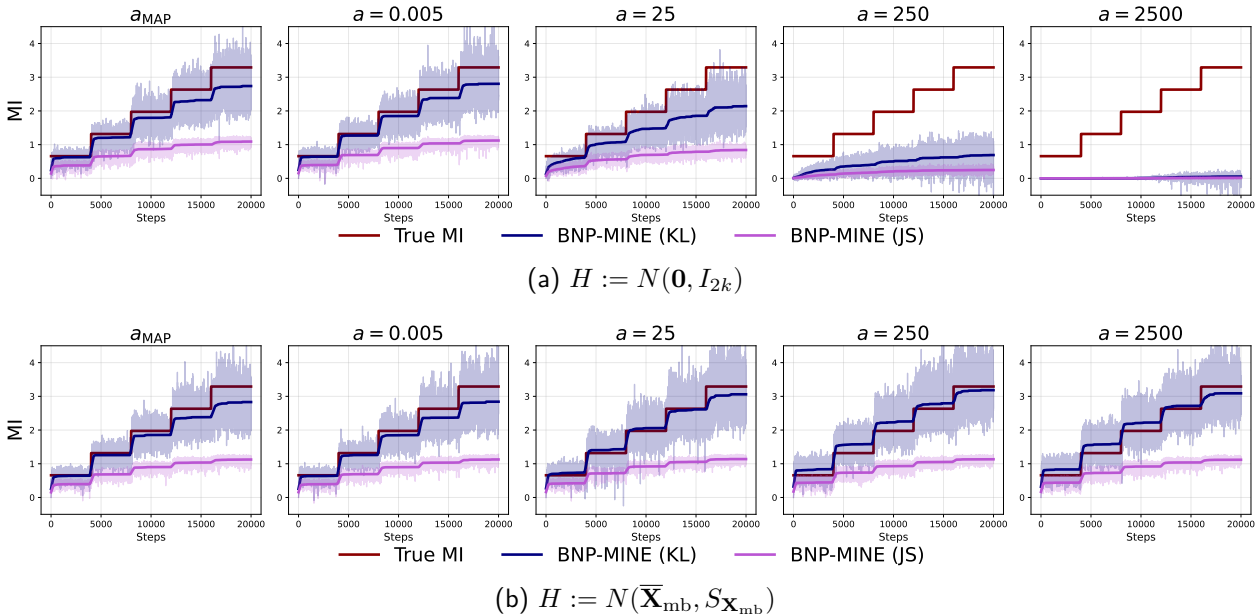


Figure 3: Staircase scheme for the binary-Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE is evaluated using a mini-batch size of 64, with  $F \sim \text{DP}(a, H)$  under different choices of the concentration parameter  $a$  and the centering measure  $H$ , across dimensions  $k = 1, \dots, 5$ .

#### 4.1.4 Comparison with further differentiable MI-VLB estimators

In addition to the traditional MINE estimator (Belghazi et al., 2018), we compare our method against several influential differentiable MI-VLB approaches. NWJ (Nguyen et al., 2010) is unbiased but typically yields a loose and high-variance lower bound. CPC/InfoNCE (Oord et al., 2018) offers stable optimization but is fundamentally constrained by  $\ln(n_{\text{mb}})$ , restricting its ability to capture high MI. We also include JSD-LB (Dorent et al., 2025), the most recent state-of-the-art estimator providing a tight and tractable lower bound connecting JS divergence to MI. As our experiments show, BNP-MINE delivers a more stable and reliable lower-bound estimate than NWJ and JSD-based methods, and although CPC appears smoother, its structural upper bound prevents consistent performance across settings. Additional results in the appendix demonstrate that CPC becomes non-tight in several benchmarks, while our estimator consistently preserves the staircase structure with substantially reduced variance.

**Sensitivity to high dimensionality.** To examine how well our method handles high-dimensional data, we test it with  $k = 10, 100, 1000$  as shown in Figure 5. In all cases, DPMINE produces smooth and stable curves with much lower noise compared to the baseline estimators. As the dimension increases, MINE and NWJ become noticeably more unstable, while DPMINE continues to give a steady estimate of the MI lower bound. These results show that our estimator remains reliable even in very high dimensional settings.

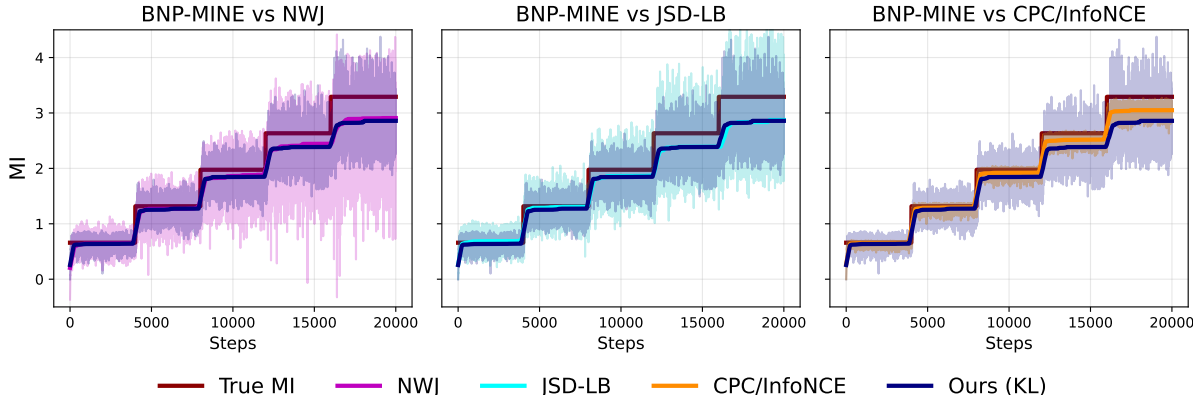


Figure 4: Staircase scheme for the binary-Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE is compared against different baselines using a fixed mini-batch size of 64, with  $F \sim \text{DP}(a_{\text{MAP}}, H)$  and  $H := N(\bar{\mathbf{X}}_{\text{mb}}, S_{\mathbf{X}_{\text{mb}}})$ , across dimensions  $k = 1, \dots, 5$ .

Table 1: Runtime of neural variational lower-bound estimators across data dimensions and training epochs.

Estimator	Epoch	Dimension ( $d = 2 \times k$ )			
		2	20	200	2000
MINE, NWJ, CPC/InfoNCE, JSD-LB	500	5 sec	5 sec	5 sec	5 sec
	1500	12 sec	12 sec	12 sec	12 sec
DPMINE	500	42 sec	42 sec	42 sec	42 sec
	1500	1 min 25 sec	1 min 25 sec	1 min 25 sec	1 min 25 sec

## 4.2 Improving VAE-GAN training using DPMINE within BNPL

To implement the BNPWMMMD+DPMINE, we consider the Gaussian kernel function, defined as  $k_{G_\sigma}(\mathbf{X}, \mathbf{Y}) = \exp(-\frac{\|\mathbf{X}-\mathbf{Y}\|^2}{2\sigma^2})$  with bandwidth parameter  $\sigma$ , in the MMD distance given in Appendix B. We search for the appropriate bandwidth parameter  $\sigma$  over a fixed grid of values,  $\sigma \in \{2, 5, 10, 20, 40, 80\}$ . We then compute the mixture of Gaussian kernels, denoted as  $k(\cdot, \cdot) = \sum_{\sigma} k_{G_\sigma}(\cdot, \cdot)$ . This selection of kernel function and bandwidth has been shown to yield satisfactory performance in training MMD-based GANs, as mentioned in Li et al. (2015); Fazeli-Asl et al. (2024); Fazeli-Asl & Zhang (2023). Then, we consider next examples to investigate the proposed approach.

### 4.2.1 Toy Example

**Coil Dataset:** We first look at a toy example that showcases the effectiveness of the BNPWMMMD+DPMINE in mitigating mode collapse. In this example, we simulate 5000 samples in 3D space, denoted as  $(X(t), Y(t), Z(t))$ , where  $X(t) = 6 \cos t$ ,  $Y(t) = 6 \sin t$ , and  $Z(t) = t$ , with  $t$  ranging from  $-2\pi$  to  $4\pi$ . We then normalize all datasets to a range between  $-1$  and  $1$ . This normalization ensures compatibility with the hyperbolic tangent activation function, which is used in the generator’s last layer in all compared models. We also used a latent dimension of 100 with a sub-latent dimension of 10 in this example. Additionally, we provide the results of the BNPWMMMD to display the basic model’s performance in covering data space in the absence of DPMINE.

Figure 6 illustrates the effect of incorporating DPMINE into the BNPWMMMD-GAN framework. Compared to BNPWMMMD alone, BNPWMMMD+DPMINE produces generated samples that more closely follow the geometric structure of the real data, with reduced dispersion and improved alignment along the underlying manifold. This improvement is particularly evident in the random generation results, where DPMINE helps preserve global structure and continuity.

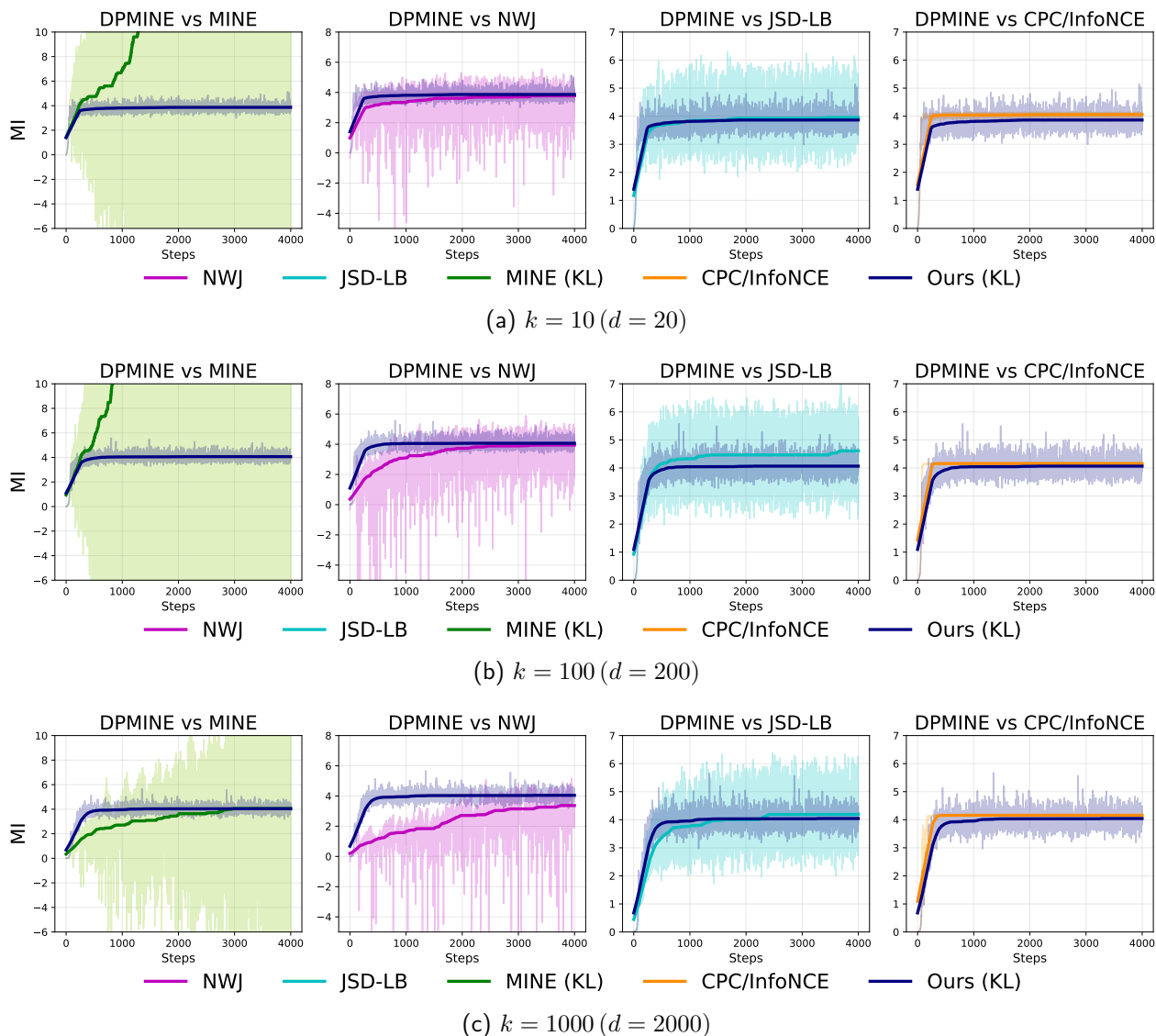


Figure 5: Binary-Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE uses mini-batch size 64, with  $F \sim \text{DP}(a_{\text{MAP}}, H)$  and  $H = \mathcal{N}(\bar{\mathbf{X}}_{\text{mb}}, \mathbf{S}_{\mathbf{X}_{\text{mb}}})$ . Experiments are run for  $k = 10, 100, 1000$ , using 4000 steps per dimension.

In the reconstruction task, the inclusion of DPMINE leads to more faithful reconstructions with smoother trajectories and fewer distortions, indicating improved information retention between the data and latent spaces. Overall, these results show that DP-MINE enhances representation quality, stabilizes training, and reduces mode collapse.

#### 4.2.2 Real Example

**COVID-19 Dataset:** This dataset comprises of 3D chest CT images of 1000 subjects diagnosed with lung infections after testing positive for COVID-19<sup>1</sup>. All images in this dataset are stored in the DICOM format and have a resolution of 16 bits per pixel, with dimensions of  $512 \times 512$  pixels in grayscale. We randomly selected 200 patients, resulting in a total of 91,960 images. As part of the preprocessing step, we first stored each patient’s data in a Neuroimaging Informatics Technology Initiative (NIFTI) format file. Then, we

<sup>1</sup>The dataset is freely available online at <https://doi.org/10.7910/DVN/6ACUZJ> (license: CC0 1.0).

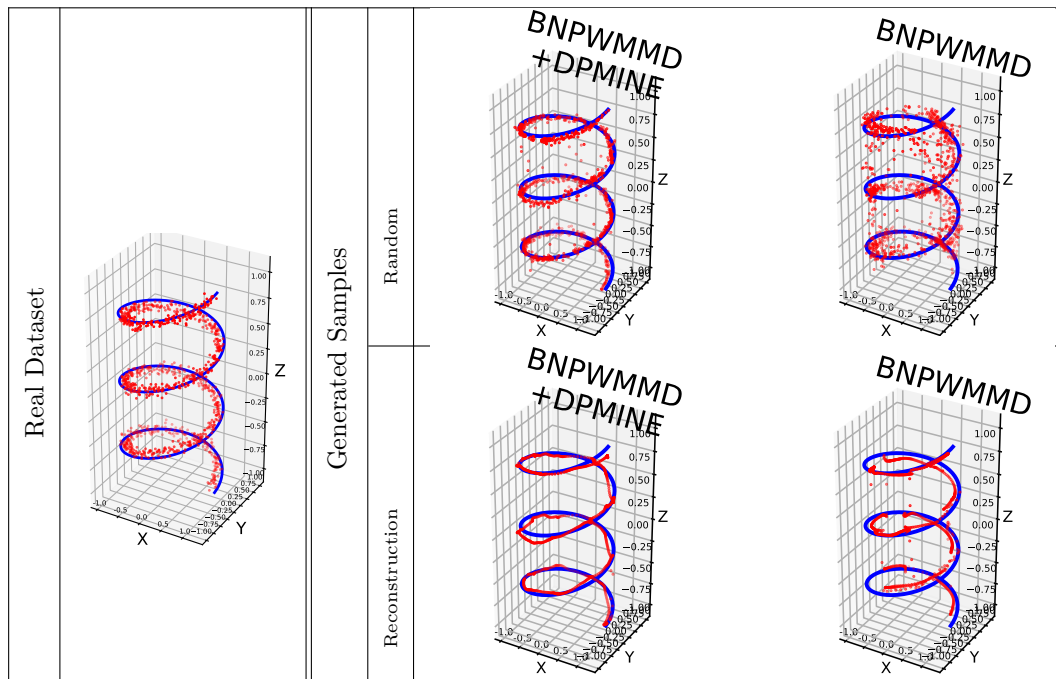


Figure 6: 1000 randomly generated and reconstructed samples using BNPWMMMD+DPMINE, compared with BNPWMMMD, after 5000 epochs for the coil example.

converted each NIFTI file into a 3D image with a dimension of  $64 \times 64 \times 64$ . Each dimension represents the axial, sagittal, and coronal views of the lungs. After normalizing the dataset, we used a mini-batch size of 16 and trained all compared models for 7500 epochs. We also used a latent dimension of  $p = 1000$  as suggested in Kwon et al. (2019) for high-dimensional cases, with a sub-latent dimension of  $q = 100$  as used in Fazeli-Asl & Zhang (2023).

To illustrate the image generation performance of BNPWMMMD+DPMINE in comparison with BNPWMMMD alone, Figure 7 presents a 3D visualization of real and randomly generated chest CT samples. The results indicate that BNPWMMMD+DPMINE produces sharper and more coherent images. For a more comprehensive view of the generated samples, we included additional slices from each dimension representation in Appendix E.2.2. We also compared the reconstruction capability of different models by displaying slices of a reconstructed sample from each model in the Appendix. The BNPWMMMD+DPMINE yields reconstructions that are visually closer to the training data under BNPWMMMD+DPMINE.

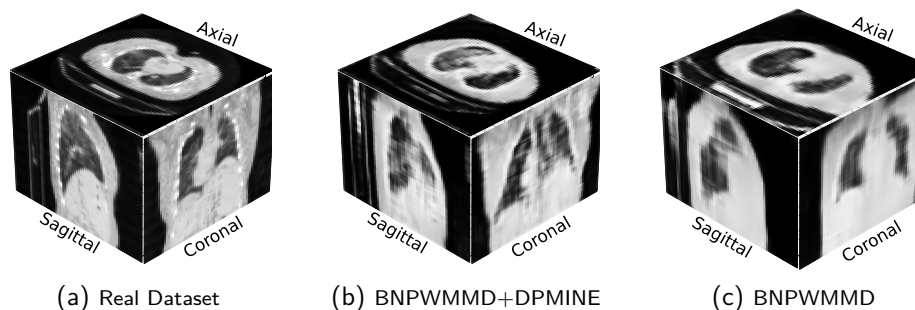


Figure 7: A 3D visualization of a real sample and a randomly generated sample using BNPWMMMD+DPMINE after 7500 epochs, compared with BNPWMMMD, in the COVID-19 example.

However, additional quantitative tools are essential for assessing the similarity between real and generated samples. Dimensionality reduction provides a useful way to visualize how well generated samples align with the real data distribution.

In Figure 8(a), the t-SNE embedding (Van der Maaten & Hinton, 2008) of the real data and 200 randomly generated samples shows that BNPWMMD+DPMINE achieves broader and more consistent coverage of the real data manifold compared to BNPWMMD alone. In contrast, samples generated by BNP-WMMD appear more concentrated and exhibit gaps relative to the real data distribution.

Figure 8(b) provides a complementary view using features obtained from a custom two-layer linear encoder. This representation supports the same conclusion: samples from BNPWMMD+DPMINE overlap more closely with real samples and better capture the underlying structure, indicating improved representation quality and reduced mode collapse.

To quantify dissimilarity between real and generated samples in the two-dimensional feature space, we consider the real features  $\mathbf{f}_r := (\text{Feature1}_{1:200,r}, \text{Feature2}_{1:200,r})$  and the generated features  $\mathbf{f}_g := (\text{Feature1}_{1:200,g}, \text{Feature2}_{1:200,g})$ . We report the mean Fréchet Inception Distance (FID) and Kernel Inception Distance (KID) over 100 replications, together with variability measured as  $\pm 2\sigma$ , using both the t-SNE embedding and a custom encoder representation.

In addition, Table 2 reports empirical MMD values computed in the same two-dimensional feature space, using the same kernel and bandwidth parameters as those employed in the BNPWMMD training objective. We intentionally avoid reporting MMD in the original data space, since BNPWMMD-based models explicitly minimize data-space MMD during training, which would render such comparisons unfair.

Across all feature representations and metrics, BNPWMMD+DP-MINE consistently attains lower mean FID, KID, and MMD values than BNPWMMD alone. While the absolute differences are moderate and the  $\pm 2\sigma$  intervals indicate non-negligible variability, the improvements are systematic across both t-SNE and custom-encoder embeddings. These trends are consistent with the visualizations in Figure 8 and suggest improved alignment between generated and real feature distributions, together with reduced mode collapse.

We further assess perceptual image quality using the Multi-Scale Structural Similarity Index (MS-SSIM) (Wang et al., 2003). MS-SSIM is reported as a mean  $\pm 2\sigma$  over 100 replications, with higher values indicating better structural fidelity. Although the compared models are unconditional generators, the use of MS-SSIM is reasonable due to the strong anatomical alignment of axial, sagittal, and coronal lung CT slices across patients. The higher MS-SSIM score achieved by BNPWMMD+DP-MINE indicates improved visual consistency relative to BNPWMMD.

Additional implementation details are provided in Appendix E.3. Appendix E.2.2 further extends the evaluation to a brain MRI dataset, demonstrating that the observed improvements generalize across datasets and imaging modalities.

Table 2: Comparison of statistical scores (mean  $\pm 2\sigma$ ) for the COVID-19 example over 100 replications.

Evaluator		BNPWMMD + DPMINE	BNPWMMD
FID	Custom Encoder	<b>0.00063 <math>\pm</math> 0.00468</b>	0.00084 $\pm$ 0.00535
	t-SNE	<b>9.48447 <math>\pm</math> 7.05185</b>	11.10431 $\pm$ 10.8224
KID	Custom Encoder	<b>0.00069 <math>\pm</math> 0.00073</b>	0.00361 $\pm$ 0.00102
	t-SNE	<b>0.99858 <math>\pm</math> 4.16996</b>	1.011342 $\pm$ 6.66391
MMD	Custom Encoder	<b>0.00902 <math>\pm</math> 0.00236</b>	0.01172 $\pm$ 0.00703
	t-SNE	<b>0.00074 <math>\pm</math> 0.00130</b>	0.00358 $\pm$ 0.00782
MS-SSIM	–	<b>0.43465 <math>\pm</math> 0.02189</b>	0.41781 $\pm$ 0.06755

## 5 Concluding Remarks

In this work, a Bayesian bootstrap resampling framework was introduced for mutual information estimation via MINE (DPMINE), enabling variational lower-bound optimization within Bayesian nonparametric

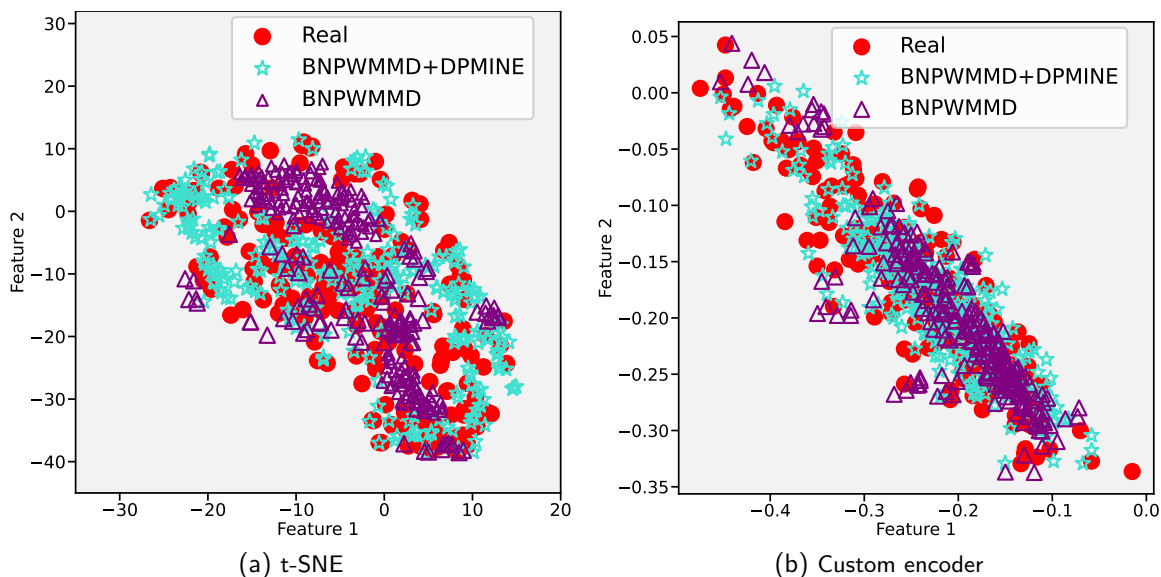


Figure 8: Scatter plots of two-dimensional feature representations in the COVID-19 example, comparing samples generated using BNPWMMD+DPMINE with those generated using BNPWMMD. Panel (a) shows the t-SNE embedding of the features, while panel (b) displays the features obtained from a custom encoder.

learning (BNPL). By replacing uniform weights with Dirichlet-resampled weights, the proposed approach naturally accounts for uncertainty in the underlying data distribution, a core feature of BNPL.

The resulting estimator, DPMINE, yields tighter and lower-variance MI estimates than standard MINE and performs competitively with existing variational bounds. Empirical results show stable behavior across batch sizes and high-dimensional settings. When incorporated into a BNP VAE-GAN model, DPMINE improves training stability and generative quality, demonstrating its practical value.

Overall, this work bridges an important gap between MI-based learning and BNPL and enables broader use of information-theoretic objectives in Bayesian nonparametric frameworks. While the focus here is on generative modeling, the proposed framework is applicable to other learning problems, including representation learning and Bayesian decision-making.

Extending DPMINE to non-IID settings and distributed or privacy-sensitive learning scenarios remains an important direction for future work. In particular, incorporating DPMINE into federated learning frameworks could strengthen privacy-preserving training by providing more robust and stable control of information flow across decentralized and heterogeneous data sources. However, the current formulation relies on an IID assumption underlying the Bayesian bootstrap. Adapting DPMINE to federated settings will therefore require new methodological developments to handle client-level heterogeneity, non-IID sampling, and communication constraints while preserving theoretical guarantees.

## Acknowledgments

Bei Jiang and Linglong Kong were partially supported by grants from the Canada CIFAR AI Chairs program, the Alberta Machine Intelligence Institute (AMII), and Natural Sciences and Engineering Council of Canada (NSERC), and Linglong Kong was also partially supported by grants from the Canada Research Chair program from NSERC. Michael Zhang was supported by the University of Hong Kong Seed Fund for PI Research #2402101367. The authors would also like to thank the anonymous reviewers for their constructive comments that improved the quality of this article.

## References

- Luai Al-Labadi, Forough Fazeli Asl, and Zahra Saberi. A test for independence via Bayesian nonparametric estimation of mutual information. *Canadian Journal of Statistics*, 50(3):1047–1070, 2022.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *International Conference on Machine Learning*, pp. 531–540. PMLR, 2018.
- Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. In *International Conference on Learning Representations*, 2018.
- Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in Neural Information Processing Systems*, 29, 2016.
- Charita Dellaporta, Jeremias Knoblauch, Theodoros Damoulas, and François-Xavier Briol. Robust Bayesian inference for simulator-based models via the MMD posterior bootstrap. In *International Conference on Artificial Intelligence and Statistics*, pp. 943–970. PMLR, 2022.
- Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- Monroe D Donsker and SR Srinivasa Varadhan. Asymptotic evaluation of certain Markov process expectations for large time. IV. *Communications on Pure and Applied Mathematics*, 36(2):183–212, 1983.
- Reuben Dorent, Polina Golland, and William M Wells Iii. Connecting jensen-shannon and kullback-leibler divergences: A new bound for representation learning. In *NeurIPS 2025*, 2025.
- Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.
- Forough Fazeli-Asl and Michael Minyi Zhang. A Bayesian non-parametric approach to generative models: Integrating variational autoencoder and generative adversarial networks using Wasserstein and maximum mean discrepancy. *arXiv preprint arXiv:2308.14048*, 2023.
- Forough Fazeli-Asl, Michael Minyi Zhang, and Lizhen Lin. A semi-bayesian nonparametric estimator of the maximum mean discrepancy measure: Applications in goodness-of-fit testing and generative adversarial networks. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=1UnlHS1FYT>.
- Thomas S Ferguson. A Bayesian analysis of some nonparametric problems. *Annals of Statistics*, 1(2): 209–230, 1973.
- Edwin Fong, Simon Lyddon, and Chris Holmes. Scalable nonparametric sampling from multimodal posteriors with the posterior bootstrap. In *International Conference on Machine Learning*, pp. 1952–1962. PMLR, 2019.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27:2672–2680, 2014.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(1):723–773, 2012.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of Wasserstein GANs. *Advances in Neural Information Processing Systems*, 30, 2017.
- R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018.

- Hemant Ishwaran and Mahmoud Zarepour. Exact and approximate sum representations for the Dirichlet process. *Canadian Journal of Statistics*, 30(2):269–283, 2002.
- Desi R Ivanova, Marvin Schmitt, and Stefan T Radev. Data-efficient variational mutual information estimation via bayesian self-consistency. In *Proc. Workshop on Bayesian Decision-Making and Uncertainty*. URL: <https://openreview.net/forum>, 2024.
- Daniel C Jones, Patrick Danaher, Youngmi Kim, Joseph M Beechem, Raphael Gottardo, and Evan W Newell. An information theoretic approach to detecting spatially varying genes. *Cell Reports Methods*, 3(6), 2023.
- Gihyun Kwon, Chihye Han, and Dae-Shik Kim. Generation of 3D brain MRI using auto-encoding generative adversarial networks. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III 22*, pp. 118–126. Springer, 2019.
- Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *International Conference on Machine Learning*, pp. 1558–1566. PMLR, 2016.
- Nunzio Alexandro Letizia, Nicola Novello, and Andrea M Tonello. Mutual information estimation via  $f$ -divergence and data derangements. *Advances in Neural Information Processing Systems*, 37:105114–105150, 2024.
- Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *International Conference on Machine Learning*, pp. 1718–1727. PMLR, 2015.
- Ruizhi Liao, Daniel Moyer, Polina Golland, and William M Wells. Demi: Discriminative estimator of mutual information. *arXiv preprint arXiv:2010.01766*, 2020.
- Simon Lyddon, Stephen Walker, and Chris C Holmes. Nonparametric learning from Bayesian models with randomized objective functions. *Advances in Neural Information Processing Systems*, 31, 2018.
- Simon P Lyddon, CC Holmes, and SG Walker. General Bayesian updating and the loss-likelihood bootstrap. *Biometrika*, 106(2):465–478, 2019.
- XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka.  $f$ -gan: Training generative neural samplers using variational divergence minimization. *Advances in Neural Information Processing Systems*, 29, 2016.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Jayaram Sethuraman. A constructive definition of Dirichlet priors. *Statistica Sinica*, pp. 639–650, 1994.
- Yuval Shalev, Amichai Painsky, and Irad Ben-Gal. Neural joint entropy estimation. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4):5488–5500, 2022.
- Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information estimators. *arXiv preprint arXiv:1910.06222*, 2019.
- Masashi Sugiyama. Machine learning with squared-loss mutual information. *Entropy*, 15(1):80–112, 2012.
- Yao-Hung Hubert Tsai, Han Zhao, Makoto Yamada, Louis-Philippe Morency, and Russ R Salakhutdinov. Neural methods for point-wise dependency estimation. *Advances in Neural Information Processing Systems*, 33:62–72, 2020.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using  $t$ -SNE. *Journal of Machine Learning Research*, 9(11), 2008.

C Villani. Optimal transport, old and new. *Grundlehren der mathematischen Wissenschaften*, 3, 2008.

Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pp. 1398–1402. Ieee, 2003.

Mahmoud Zarepour and Luai Al-Labadi. On a rapid simulation of the Dirichlet process. *Statistics & Probability Letters*, 82(5):916–924, 2012.

## Appendix

### A Notations

- $F$ : Real data distribution.
- $(\mathbf{X}_{1:n})$ : A sample of  $n$  independent and identically distributed random variables generated from  $F$ .
- $\text{DP}(a + n, H_n^*)$ : The Dirichlet process (DP) posterior with concentration parameter  $a + n$ ,  $a > 0$ , and base measure  $H_n^*$ .
- $H_n^*$ : A mixture probability measure given by  $\frac{a}{a+n}H + \frac{n}{a+n}F_{\mathbf{x}_{1:n}}$ .
- $H$ : The base measure of the DP prior  $\text{DP}(a, H)$ .
- $F_{\mathbf{x}_{1:n}}$ : The empirical probability measure based on the sample data defined by  $\frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$ .
- $F_N^{\text{Pos}}$ : The probability measure of the DP posterior approximation defined by  $\sum_{i=1}^N J_i^{\text{Pos}} \delta_{\mathbf{x}_i^{\text{Pos}}}$  with  $(\mathbf{X}_{1:N}^{\text{Pos}}) \sim H_n^*$  and  $(J_{1:N}^{\text{Pos}}) \sim \text{Dirichlet}((a + n)/N, \dots, (a + n)/N)$  Ishwaran & Zarepour (2002).
- $f_i(\mathbf{X})$ : A continuous function of a random variable  $\mathbf{X} \sim F$  for  $i = 1, 2$ .
- $\{T_\gamma\}_{\gamma \in \Gamma}$ : A set of continuous functions parameterized by a neural network on a compact domain  $\Gamma$ .
- $I_d$ : Identity matrix of size  $d \times d$ .
- “a.s.”: Standing for “almost surely” and indicates that the statements hold with probability 1.

### B Definition of MMD and Wasserstein Distances: DP and Empirical Representations

Selecting an appropriate statistical distance is crucial for effective generative model training. Here, we focus on the DP representation of two popular distances used in BNP deep learning and we will briefly mention their frequentist counterparts.

#### B.1 Maximum Mean Discrepancy Distance (Feature-Matching Comparison)

The MMD distance was initially introduced in Gretton et al. (2012) for frequentist two-sample comparisons. Recently, a DP-based version of this distance has been proposed for BNP hypothesis testing Fazeli-Asl et al. (2024). Consider a set of functions  $\{G_\omega\}_{\omega \in \Omega}$  parameterized by a neural network that can generate  $n$  IID random variables  $(\mathbf{Y}_{1:n})$ , where the likelihood function is intractable and not accessible. Let  $k(\cdot, \cdot)$  be a continuous kernel function with a feature space corresponding to a universal reproducing kernel Hilbert space Gretton et al. (2012) defined on a compact sample space  $\mathfrak{X}$  Gretton et al. (2012). Given a sample

$(\mathbf{X}_{1:n}) \stackrel{\text{iid}}{\sim} F$ , the MMD distance between  $F^{\text{Pos}}$  and  $F_{\mathbf{Y}_{1:n}}$  is approximated as:

$$\text{MMD}^2(F_N^{\text{Pos}}, F_{\mathbf{Y}_{1:n}}) := \sum_{\ell, t=1}^N J_\ell^{\text{Pos}} J_t^{\text{Pos}} k(\mathbf{X}_\ell^{\text{Pos}}, \mathbf{X}_t^{\text{Pos}}) - \frac{2}{n} \sum_{\ell=1}^N \sum_{t=1}^n J_\ell^{\text{Pos}} k(\mathbf{X}_\ell^{\text{Pos}}, \mathbf{Y}_t) + \frac{1}{n^2} \sum_{\ell, t=1}^n k(\mathbf{Y}_\ell, \mathbf{Y}_t). \quad (26)$$

In the frequentist version of the MMD distance, as defined in Gretton et al. (2012), the empirical distribution  $F_{\mathbf{X}_{1:n}}$  is considered. This is denoted as  $\text{MMD}^2(F_{\mathbf{X}_{1:n}}, F_{\mathbf{Y}_{1:n}})$ , which is obtained by replacing  $N$ ,  $J_{1:N}^{\text{Pos}}$ , and  $\mathbf{X}_{1:N}^{\text{Pos}}$  with  $n$ ,  $1/n$ , and  $\mathbf{X}_{1:n}$ , respectively, in equation 26.

## B.2 Wasserstein Distance (Overall Distribution Comparison)

The frequentist version of the Wasserstein distance is completely discussed in (Villani, 2008, Part I6). Fazeli-Asl et al. Fazeli-Asl & Zhang (2023) proposed a BNP version of this distance through its Kantorovich-Rubinstein dual representation. Let  $\{D_\theta\}_{\theta \in \Theta}$  be a parametrized family of continuous functions that all are 1-Lipschitz. Then the Wasserstein distance between  $F^{\text{Pos}}$  and  $F_{\mathbf{Y}_{1:n}}$  is approximated as:

$$\text{WS}(F_N^{\text{Pos}}, F_{\mathbf{Y}_{1:n}}) := \max_{\Theta} \sum_{i=1}^N \left( J_i^{\text{Pos}} D_\theta(\mathbf{X}_i^{\text{Pos}}) - \frac{D_\theta(\mathbf{Y}_i)}{n} \right). \quad (27)$$

By utilizing modifications in equation 27 similar to those described for the MMD distance, the empirical representation of the Wasserstein distance can also be obtained. In this section, we propose two novel representations for the MINE using the DP. These representations are based on the KL and JS divergences will be used in our BNP learning framework to maximize information during the training process.

## C BNPWMMD-GAN

The BNPWMMD-GAN architecture introduced by Fazeli-Asl & Zhang (2023) consists of four networks: an encoder  $\{E_\eta\}_{\eta \in \mathcal{H}}$ , a generator  $\{G_\omega\}_{\omega \in \Omega}$ , a code generator  $\{CG_{\omega'}\}_{\omega' \in \Omega'}$ , and a discriminator  $\{D_\theta\}_{\theta \in \Theta}$ .

Given the codes  $(\mathbf{c}_{1:N} := E_\eta(\mathbf{X}_{1:N}^{\text{Pos}}))$  and generated codes  $(\tilde{\mathbf{c}}_{1:N} := CG_{\omega'}(\boldsymbol{\varepsilon}'_{1:N}))$  along with the latent noise  $(\boldsymbol{\varepsilon}_{1:N} \in \mathbb{R}^p) \sim F_\boldsymbol{\varepsilon}$  and sub-latent noise  $(\boldsymbol{\varepsilon}'_{1:N} \in \mathbb{R}^q) \sim F_{\boldsymbol{\varepsilon}'}$ , where  $q < p$ , the BNPWMMD-GAN is trained by updating the networks' parameters according to the following hybrid objective function<sup>2</sup>:

$$(\hat{\omega}, \hat{\eta})(F^{\text{Pos}}) = \arg \min_{\Omega, \mathcal{H}} \left\{ \underbrace{-\frac{1}{N} \sum_{i=1}^N [D_\theta(G_\omega(\boldsymbol{\varepsilon}_i)) + D_\theta(G_\omega(\mathbf{c}_i)) + D_\theta(G_\omega(\tilde{\mathbf{c}}_i))] + \text{MMD}(F_N^{\text{Pos}}, F_{G_\omega(\boldsymbol{\varepsilon}_{1:N})})}_{\mathcal{I}_1(\omega, \eta)} - \underbrace{\text{MMD}(F_N^{\text{Pos}}, F_{G_\omega(\tilde{\mathbf{c}}_{1:N})})}_{\mathcal{I}_1(\omega, \eta)} + \underbrace{\text{MMD}(F_N^{\text{Pos}}, F_{G_\omega(\mathbf{c}_{1:N})})}_{\mathcal{I}_1(\omega, \eta)} + \underbrace{\text{MMD}(F_{\boldsymbol{\varepsilon}_{1:N}}, F_{\mathbf{c}_{1:N}})}_{\mathcal{I}_2(\omega, \eta)} \right\}, \quad (28a)$$

$$\hat{\theta}(F^{\text{Pos}}) = \arg \min_{\Theta} \left\{ \underbrace{\frac{1}{N} \sum_{i=1}^N [D_\theta(G_\omega(\boldsymbol{\varepsilon}_i)) + D_\theta(G_\omega(\mathbf{c}_i)) + D_\theta(G_\omega(\tilde{\mathbf{c}}_i)) - 3J_i^{\text{Pos}} D_\theta(\mathbf{X}_i^{\text{Pos}})]}_{\mathcal{J}_1(\theta)} + \lambda \underbrace{L_{\text{GP-D}}}_{\mathcal{J}_2(\theta)} \right\}, \quad (28b)$$

$$\hat{\omega}'(F^{\text{Pos}}) = \arg \min_{\Omega'} \text{MMD}(F_{\mathbf{c}_{1:N}}, F_{\tilde{\mathbf{c}}_{1:N}}). \quad (28c)$$

<sup>2</sup>Each  $F$  indexed by a sample vector in equation 28 indicates the corresponding empirical distribution.

Simultaneous updates of parameters  $\boldsymbol{\omega}$  and  $\boldsymbol{\eta}$  are facilitated by the essential role played by the generator function, which serves as a decoder during VAE training. Here, the generator is fed with  $(\boldsymbol{\varepsilon}_{1:N}, \mathbf{c}_{1:N}, \tilde{\mathbf{c}}_{1:N})$ , and terms  $\mathcal{I}_1(\boldsymbol{\omega}, \boldsymbol{\eta})$  and  $\mathcal{J}_1(\boldsymbol{\theta})$  refer to minimizing the combined distance given by:

$$d_{\text{WMMD}}(F^{\text{Pos}}, F_{G_{\boldsymbol{\omega}}}) = \text{WS}(F^{\text{Pos}}, F_{G_{\boldsymbol{\omega}}}) + \text{MMD}(F^{\text{Pos}}, F_{G_{\boldsymbol{\omega}}}). \quad (29)$$

Since the MMD measure is an  $L^2$ -norm distance, the MMD-based terms in  $\mathcal{I}_1(\boldsymbol{\omega}, \boldsymbol{\eta})$  serves as the posterior reconstruction errors, while  $\mathcal{I}_2(\boldsymbol{\omega}, \boldsymbol{\eta})$  serves as the regularization error for approximating the variational distribution  $F_{E_{\boldsymbol{\eta}}}(\mathbf{c}_{1:N} | \mathbf{X}_{1:N}^{\text{Pos}})$ . Furthermore, the gradient penalty  $\mathcal{J}_2(\boldsymbol{\theta})$  with a positive coefficient  $\lambda$  was used to ensure training stability by forcing the 1-Lipschitz constraint on the discriminator (Gulrajani et al., 2017). Specifically, this penalty is given by

$$L_{\text{GP-D}} = \frac{1}{N} \sum_{i=1}^N (\|\nabla_{\widehat{\mathbf{X}}_i^{\text{Pos}}} D_{\boldsymbol{\theta}}(\widehat{\mathbf{X}}_i^{\text{Pos}})\|_2 - 1)^2, \quad (30)$$

where  $\widehat{\mathbf{X}}_i^{\text{Pos}} = u\tilde{\mathbf{X}}_i^{\text{Pos}} + (1-u)\mathbf{X}_i^{\text{Pos}}$  and  $\tilde{\mathbf{X}}_i^{\text{Pos}}$  represents any posterior fake sample generated by the generator.

## D Theoretical Proofs

**Theorem 1.** Consider the DP posterior representations of MINE presented in the main paper. Assume that both  $T_{\boldsymbol{\gamma}}$  and  $e^{T_{\boldsymbol{\gamma}}}$  are  $M$ -bounded for all  $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}$  (i.e.,  $|T_{\boldsymbol{\gamma}}|, |e^{T_{\boldsymbol{\gamma}}}| \leq M$ ). Given the DP posterior approximation in equation 7, for any  $N > 0$ , we have,

$$i. \lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{\text{Pos}}} \left( L_{\boldsymbol{\gamma}}^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right) \geq L_{\boldsymbol{\gamma}}^{\text{KL}}(\mathbf{X}'_1, \mathbf{X}'_2),$$

$$ii. \lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{\text{Pos}}} \left( L_{\boldsymbol{\gamma}}^{\text{DPJS}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right) \geq L_{\boldsymbol{\gamma}}^{\text{JS}}(\mathbf{X}'_1, \mathbf{X}'_2).$$

*Proof.* Recall that, for all  $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}$ ,

$$L_{\boldsymbol{\gamma}}^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) := \sum_{\ell=1}^N J_{\ell}^{\text{Pos}} T_{\boldsymbol{\gamma}}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}})) - \ln \sum_{\ell=1}^N J_{\ell}^{\text{Pos}} e^{T_{\boldsymbol{\gamma}}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}}))}, \quad (31)$$

where  $(J_{1:N}^{\text{Pos}}) \sim \text{Dirichlet}((a+n)/N, \dots, (a+n)/N)$  and  $(\mathbf{X}_{1:N}^{\text{Pos}}) \stackrel{\text{iid}}{\sim} H_n^*$ .

Applying the law of total expectation yields:

$$\mathbb{E}_{F^{\text{Pos}}} \left( L_{\boldsymbol{\gamma}}^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right) = \mathbb{E}_{H_n^*} \left( \mathbb{E}_{F_{J^{\text{Pos}}}} \left( L_{\boldsymbol{\gamma}}^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \mid \mathbf{X}_{1:N}^{\text{Pos}} \right) \right), \quad (32)$$

where  $F_{J^{\text{Pos}}}$  is the CDF corresponding to  $\mathbf{J}^{\text{Pos}} := (J_{1:N}^{\text{Pos}})$ .

Taking the expectation in equation 32 with respect to the Dirichlet weights gives:

$$\mathbb{E}_{F_{J^{\text{Pos}}}} (L_{\gamma}^{\text{DPKL}} (f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \mid \mathbf{X}_{1:N}^{\text{Pos}}) = \sum_{\ell=1}^N \mathbb{E}_{F_{J^{\text{Pos}}}} (J_{\ell}^{\text{Pos}}) T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}})) - \mathbb{E}_{F_{J^{\text{Pos}}}} \left( \ln \sum_{\ell=1}^N J_{\ell}^{\text{Pos}} e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))} \right) \quad (33a)$$

(Linearity of Expectation)

$$\geq \sum_{\ell=1}^N \mathbb{E}_{F_{J^{\text{Pos}}}} (J_{\ell}^{\text{Pos}}) T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}})) - \ln \sum_{\ell=1}^N \mathbb{E}_{F_{J^{\text{Pos}}}} (J_{\ell}^{\text{Pos}}) e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))} \quad (33b)$$

(Jensen's Inequality)

$$\geq \sum_{\ell=1}^N \frac{1}{N} T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}})) - \ln \sum_{\ell=1}^N \frac{1}{N} e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))} \quad (33c)$$

(Dirichlet Weights).

Taking the expectation in equation 33 with respect to the posterior atoms gives

$$\mathbb{E}_{H_n^*} (\mathbb{E}_{F_{J^{\text{Pos}}}} (L_{\gamma}^{\text{DPKL}} (f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \mid \mathbf{X}_{1:N}^{\text{Pos}})) \geq \sum_{\ell=1}^N \frac{1}{N} \mathbb{E}_{H_n^*} (T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}}))) - \mathbb{E}_{H_n^*} \left( \ln \sum_{\ell=1}^N \frac{1}{N} e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))} \right) \quad (34a)$$

(Linearity of Expectation)

$$\geq \sum_{\ell=1}^N \frac{1}{N} \mathbb{E}_{H_n^*} (T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}}))) - \ln \sum_{\ell=1}^N \frac{1}{N} \mathbb{E}_{H_n^*} (e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))}) \quad (34b)$$

(Jensen's Inequality)

$$= \mathbb{E}_{H_n^*} (T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}}))) - \ln \mathbb{E}_{H_n^*} (e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))}) \quad (34c)$$

(Identical Atoms).

On the other hand, as  $n$  approaches infinity, the Glivenko–Cantelli theorem implies that  $F_{\mathbf{X}_{1:n}}$  converges to  $F$  and subsequently  $H_n^*$  converges to  $F$ . This convergence implies that

$$\mathbf{X}_{\ell}^{\text{Pos}} \xrightarrow{d} \mathbf{X}_{\ell}, \quad (35)$$

where “ $\xrightarrow{d}$ ” denotes convergence in distribution and  $\mathbf{X}_{\ell}$  is a random variable distributed according to  $F$ , for  $\ell = 1, \dots, N$ . By applying the continuous mapping theorem, we have

$$T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\ell}^{\text{Pos}})) \xrightarrow{d} T_{\gamma}(f_1(\mathbf{X}_{\ell}), f_2(\mathbf{X}_{\ell})), \quad e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))} \xrightarrow{d} e^{T_{\gamma}(f_1(\mathbf{X}_{\ell}), f_2(\mathbf{X}_{\pi(\ell)}))}. \quad (36)$$

Since  $T_\gamma$  and  $e^{T_\gamma}$  are  $M$ -bounded, it implies

$$\mathbb{E}_{H_n^*} (T_\gamma(f_1(\mathbf{X}_\ell^{\text{Pos}}), f_2(\mathbf{X}_\ell^{\text{Pos}}))) \rightarrow \mathbb{E}_F (T_\gamma(f_1(\mathbf{X}_\ell), f_2(\mathbf{X}_\ell))), \quad (37)$$

$$\mathbb{E}_{H_n^*} (e^{T_\gamma(f_1(\mathbf{X}_\ell^{\text{Pos}}), f_2(\mathbf{X}_{\pi(\ell)}^{\text{Pos}}))}) \rightarrow \mathbb{E}_F (e^{T_\gamma(f_1(\mathbf{X}_\ell), f_2(\mathbf{X}_{\pi(\ell)}))}). \quad (38)$$

Finally, taking the limit in equation 34 as  $n \rightarrow \infty$  and applying the continuous mapping theorem implies

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}_{H_n^*} (\mathbb{E}_{F_{J^{\text{Pos}}}} (L_\gamma^{\text{DPKL}} (f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \mid \mathbf{X}_{1:N}^{\text{Pos}})) &\geq \mathbb{E}_F (T_\gamma(f_1(\mathbf{X}_\ell), f_2(\mathbf{X}_\ell))) \\ &\quad - \ln \mathbb{E}_F (e^{T_\gamma(f_1(\mathbf{X}_\ell), f_2(\mathbf{X}_{\pi(\ell)}))}) \end{aligned} \quad (39)$$

$$\begin{aligned} &= \mathbb{E}_{F_{\mathbf{X}'_1 \mathbf{X}'_2}} (T_\gamma(\mathbf{X}'_1, \mathbf{X}'_2)) \\ &\quad - \ln \mathbb{E}_{F_{\mathbf{X}'_1} \otimes F_{\mathbf{X}'_2}} (e^{T_\gamma(\mathbf{X}'_1, \mathbf{X}'_2)}) \end{aligned} \quad (40)$$

$$= L_\gamma^{\text{KL}}(\mathbf{X}'_1, \mathbf{X}'_2) \quad (41)$$

with  $\mathbf{X}'_1 = f_1(\mathbf{X}_\ell)$  and  $\mathbf{X}'_2 = f_2(\mathbf{X}_\ell)$ . Consequently, for all  $\gamma \in \Gamma$ ,

$$\lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{\text{Pos}}} [L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}}))] \geq L_\gamma^{\text{KL}}(\mathbf{X}'_1, \mathbf{X}'_2). \quad (42)$$

Moreover, since

$$\max_{\gamma \in \Gamma} L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \geq L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \quad (43)$$

for any  $\gamma \in \Gamma$ , taking expectations and limits yields

$$\lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{\text{Pos}}} \left[ \max_{\gamma \in \Gamma} L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right] \geq L_\gamma^{\text{KL}}(\mathbf{X}'_1, \mathbf{X}'_2). \quad (44)$$

Since the above inequality holds for all  $\gamma \in \Gamma$ , it follows that

$$\lim_{n \rightarrow \infty} \mathbb{E}_{F_N^{\text{Pos}}} \left[ \max_{\gamma \in \Gamma} L_\gamma^{\text{DPKL}}(f_1(\mathbf{X}_{1:N}^{\text{Pos}}), f_2(\mathbf{X}_{1:N}^{\text{Pos}})) \right] \geq \max_{\gamma \in \Gamma} L_\gamma^{\text{KL}}(\mathbf{X}'_1, \mathbf{X}'_2), \quad (45)$$

which completes the proof of part (i).

A similar method is used to prove (ii) and it is then omitted. □

## E Additional Results and Implementing Details

### E.1 DPMINE Evaluation on Complex Examples

In this section, we evaluate DPMINE against several benchmark methods on a set of complex examples adapted from Letizia et al. (2024), which are commonly used for assessing mutual information estimators.

#### E.1.1 Complex Gaussian

**Gaussian Example.** Following the same protocol used in the main experimental section, a Gaussian dimension, staircase experiment is considered to evaluate mutual information estimators as the dimensionality of dependence increases. Let

$$\mathbf{X} = (\mathbf{U}, \mathbf{V}) \in \mathbb{R}^{2k}, \quad \mathbf{U}, \mathbf{V} \in \mathbb{R}^k,$$

with  $k \in \{1, 2, 3, 4, 5\}$ , where each value of  $k$  defines one staircase level.

For a fixed correlation parameter  $\rho = 0.95$ , the variables are generated according to the linear Gaussian relationship

$$\mathbf{V} = \rho \mathbf{U} + \sqrt{1 - \rho^2} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I}_k),$$

with  $\mathbf{U} \sim N(\mathbf{0}, \mathbf{I}_k)$  independent of  $\boldsymbol{\varepsilon}$ .

The true MI at each staircase level is given by

$$\text{MI}(\mathbf{U}, \mathbf{V}) = -\frac{k}{2} \log(1 - \rho^2),$$

and is used as a reference for comparing estimator performance in this example.

The cubic, half-cube, and asinh transformations introduced below induce heavy-tailed marginals while preserving the MI of this Gaussian example, since each is an invertible transformation applied componentwise.

**Cubic Example.** Let

$$\mathbf{U} \sim N(\mathbf{0}, \mathbf{I}_k),$$

and  $\mathbf{V}$  is generated through a nonlinear transformation of a correlated Gaussian,

$$\mathbf{V} = (\rho \mathbf{U} + \sqrt{1 - \rho^2} \boldsymbol{\varepsilon})^3, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I}_k),$$

where the correlation parameter is fixed to  $\rho = 0.95$ .

**Half-Cube Example.** Let  $\mathbf{U} \sim N(\mathbf{0}, \mathbf{I}_k)$ . The variable  $\mathbf{V}$  is first generated through a linear Gaussian relationship

$$\mathbf{Z} = \rho \mathbf{U} + \sqrt{1 - \rho^2} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I}_k),$$

with  $\rho = 0.95$ , and is then transformed elementwise using a half-cube (tail-stretching) nonlinearity,

$$\mathbf{V} = \text{sign}(\mathbf{Z}) |\mathbf{Z}|^{3/2}.$$

**Asinh Example.** Let  $\mathbf{U} \in \mathbb{R}^k$  be sampled from a standard multivariate Gaussian distribution,  $\mathbf{U} \sim N(\mathbf{0}, \mathbf{I}_k)$ . The variable  $\mathbf{V}$  is generated by first forming a Gaussian vector

$$\mathbf{Z} = \rho \mathbf{U} + \sqrt{1 - \rho^2} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I}_k),$$

and then applying the elementwise inverse hyperbolic sine (asinh) transformation,

$$\mathbf{V} = \text{asinh}(\mathbf{Z}),$$

where  $\text{asinh}(\cdot)$  denotes the inverse hyperbolic sine function and  $\rho = 0.95$  is fixed across all staircase levels.

### E.1.2 Non-Gaussian Distributions

**Student Example.** To examine estimator behavior under heavy-tailed dependence, a Student staircase example is considered. For each value of  $k$ , the joint vector  $(\mathbf{U}, \mathbf{V})$  is sampled from a centered multivariate Student- $t$  distribution with the correlation structure

$$(\mathbf{U}, \mathbf{V}) \sim t_\nu \left( \mathbf{0}, \begin{pmatrix} \mathbf{I}_k & \rho \mathbf{I}_k \\ \rho \mathbf{I}_k & \mathbf{I}_k \end{pmatrix} \right),$$

where the degrees of freedom are fixed to  $\nu = 2$  and the correlation parameter is set to  $\rho = 0.95$ .

The true MI at each staircase level is approximated via Monte Carlo integration of the corresponding Student- $t$  densities and serves as a reference throughout training.

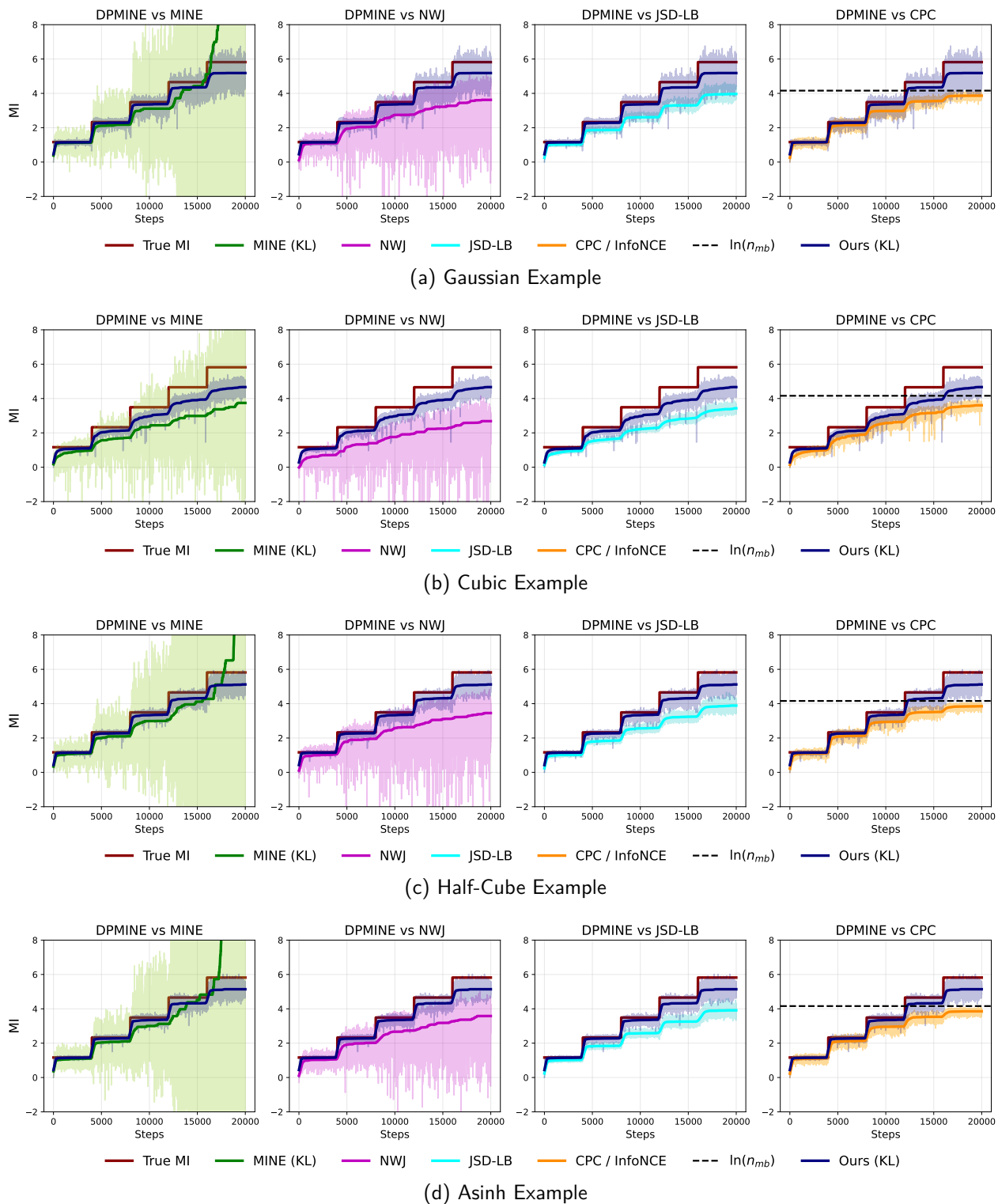
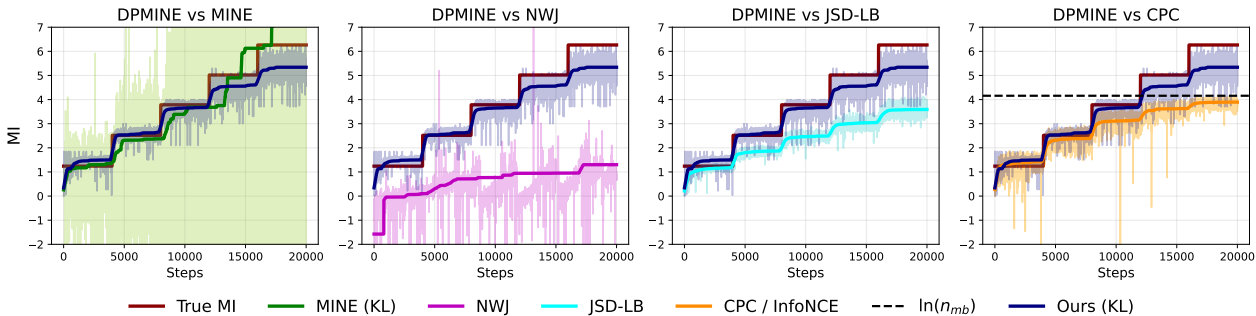


Figure 9: Complex Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE uses mini-batch size 64, with  $F \sim DP(a_{\text{MAP}}, H)$  and  $H = N(\bar{\mathbf{X}}_{\text{mb}}, S_{\mathbf{X}_{\text{mb}}})$ . Experiments are run for  $k = 1, \dots, 5$ , using 4000 steps per dimension.

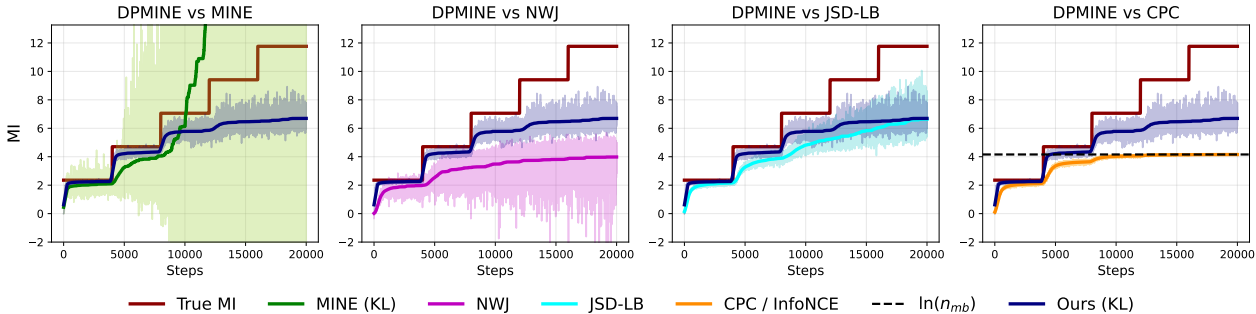
**Uniform Example.** A uniform-additive staircase example is used to study estimator behavior under bounded, non-Gaussian noise. At staircase level  $k$ , we consider a pair  $(\mathbf{U}, \mathbf{V})$  with variable  $\mathbf{U}$  is sampled uniformly on  $[0, 1]^k$ , and  $\mathbf{V}$  is generated by adding independent uniform noise,

$$\mathbf{V} = \mathbf{U} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \text{Unif}(-\varepsilon_0, \varepsilon_0)^k,$$

with  $\varepsilon_0 = 0.05$  fixed across all staircase levels. Since the dependence is additive and independent across dimensions, the MI scales linearly with the dimension. The one-dimensional MI is approximated via numerical integration and scaled by  $k$  to obtain the true MI at each staircase level.



(a) Student Example



(b) Uniform Example

Figure 10: Complex non-Gaussian setting with  $\mathbf{X} = (\mathbf{U}, \mathbf{V}) \sim F$ . BNP-MINE uses mini-batch size 64, with  $F \sim \text{DP}(a_{\text{MAP}}, H)$  and  $H = N(\bar{\mathbf{X}}_{\text{mb}}, S_{\mathbf{X}_{\text{mb}}})$ . Experiments are run for  $k = 1, \dots, 5$ , using 4000 steps per dimension.

Figures 9 and 10 compare DPMINE with several variational MI estimators across Gaussian and non-Gaussian settings. Across all examples, DPMINE tracks the true MI more closely while exhibiting substantially reduced variance compared to MINE, NWJ, JSD-LB, and CPC/InfoNCE. Competing estimators either suffer from high variability or exhibit systematic underestimation and saturation effects, particularly as MI increases. In contrast, DPMINE shows stable and monotonic convergence toward the true MI across iterations, demonstrating improved robustness and reliability in both low- and high-MI regimes.

## E.2 Application of DPMINE (KL) in VAE-GAN

### E.2.1 Toy Example

**Stanford Bunny Dataset:** We use the Stanford Bunny dataset, available at <https://graphics.stanford.edu/data/3Dscanrep/>, to investigate the impact of DPMINE on an additional toy example. The dataset consists of point clouds representing the Stanford Bunny, a renowned 3D model provided by the Stanford University Computer Graphics Laboratory. These point clouds, captured with the Cyberware 3030 MS scanner and stored in PLY files (Polygon File Format) developed at Stanford, represent spatial locations on the object’s surface. Point clouds provide a 3D spatial representation of the object, enabling

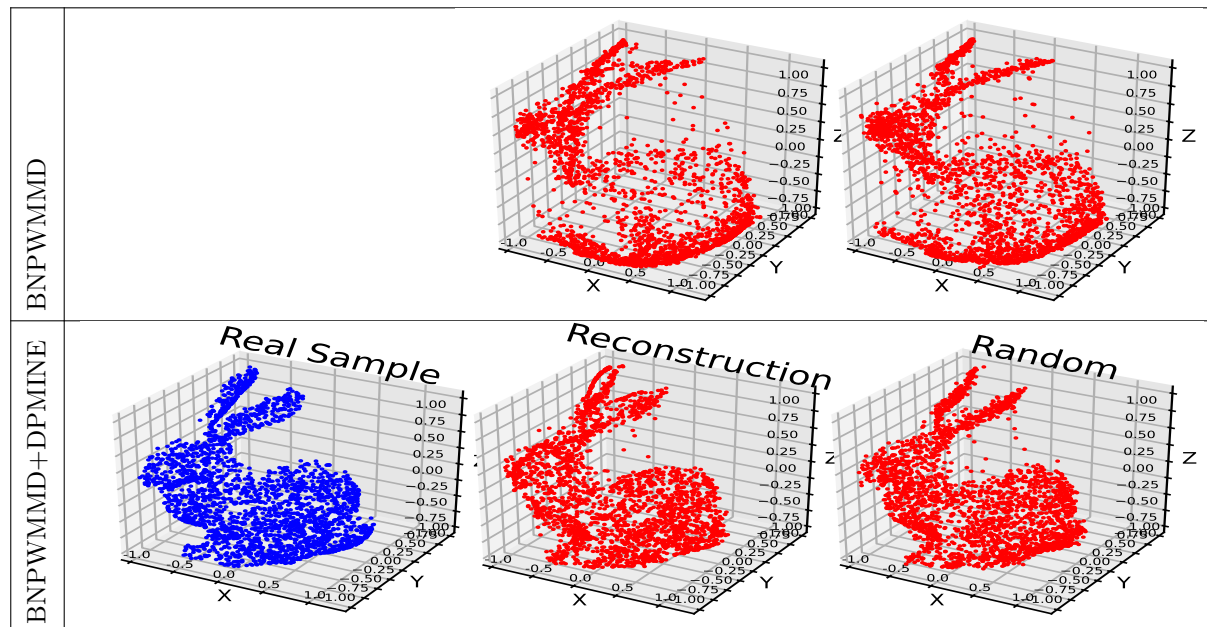


Figure 11: 2500 samples of the Stanford Bunny generated using BNPWMMD+DPMINE after 5000 epochs, compared with BNPWMMD alone.

detailed visualization and analysis. For our study, we use three point clouds from this dataset and apply filtering with the *pyoints* Python library to extract points in 3D space.

We randomly sampled 5000 points from the available 43,188 points as the training dataset and implemented our model on them by feeding the model with noise inputs to generate 2500 random samples and with encoded real inputs to obtain reconstruction samples. Figure 11 illustrates the significant impact of DPMINE on the performance of the BNP VAE-GAN in data generation.

## E.2.2 Real Examples

**Covid-19 dataset:** The red border in Figure 12 indicates the corresponding slices depicted in Figure 7 of the experimental findings discussed in the main paper. The results indicate that BNPWMMD+DPMINE yields improved sharpness and diversity in the generated 3D slices relative to the baseline methods. Figure 13 further shows reconstructed samples that exhibit increased similarity to the training data under the BNPWMMD+DPMINE model.

**BraTS 2018 Dataset:** The BRATS 2018 dataset, a benchmark resource in medical imaging available at <https://www.med.upenn.edu/sbia/brats2018/data.html>, is employed for training models in the generation of brain tumor MRI scans. For the experiments, data from 210 subjects labeled as “HGG” (High-Grade Glioma) are used, focusing on patients with aggressive brain tumors. Each subject’s MRI data includes four distinct imaging modalities: T1-weighted (T1), T1-weighted with contrast enhancement (T1ce), T2-weighted (T2), and Fluid Attenuated Inversion Recovery (FLAIR). Notably, the FLAIR modality, which is particularly effective for highlighting edema and tumor boundaries, is used for the experiments, providing essential insights for synthetic MRI generation.

The results presented in Figures 14-17 and Table 3 reinforce our previous findings on the lung dataset, demonstrating similarly strong performance when applying the methodology to the BraTS dataset. This cross-domain evaluation further validates the model’s robustness and adaptability to diverse data types.

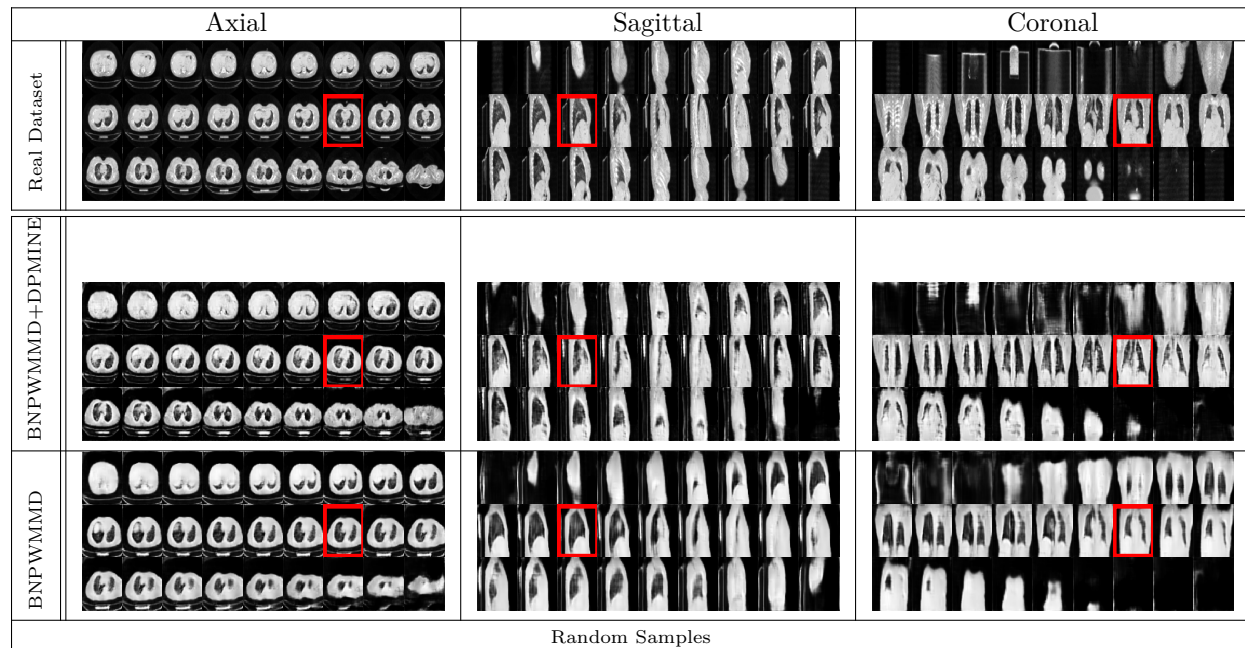


Figure 12: 27 slices along each axis of a 3D sample randomly generated using BNPWMMD+DPMINE after 7500 epochs, compared with BNPWMMD, for the COVID-19 example.

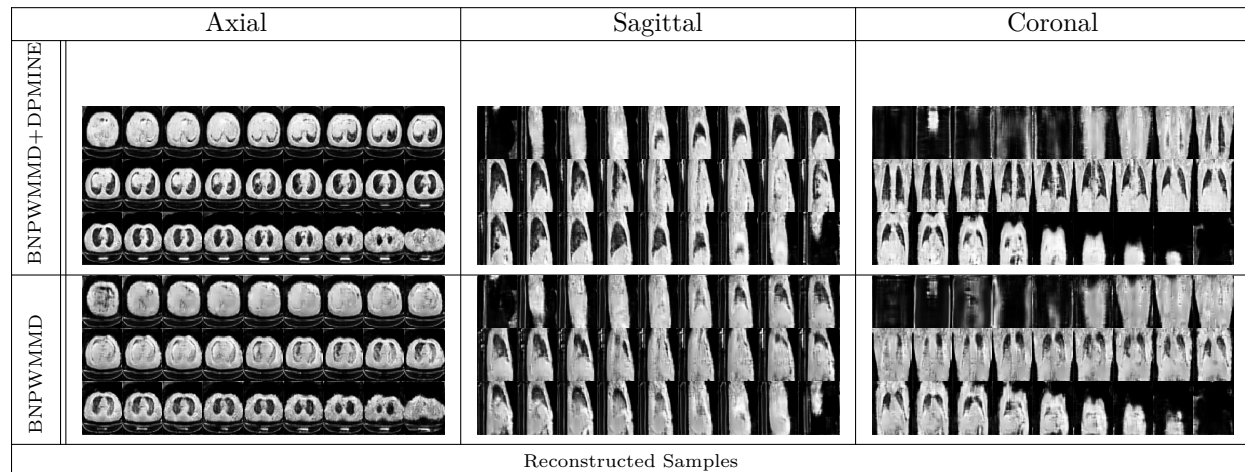


Figure 13: 27 slices along each axis of a 3D sample reconstructed using BNPWMMD+DPMINE after 7500 epochs, compared with BNPWMMD, for the COVID-19 example.

### E.3 Implementing Details

#### E.3.1 Calculation of Evaluation Scores

Considering the real and generated features,  $\mathbf{f}_r := (\text{Feature}_{1:200,r}, \text{Feature}_{21:200,r})$  and  $\mathbf{f}_g := (\text{Feature}_{1:200,g}, \text{Feature}_{21:200,g})$ , the FID and KID metrics are calculated using the following Eqs., respectively Bińkowski et al. (2018):

$$\text{FID}(\mathbf{f}_r, \mathbf{f}_g) = \|\boldsymbol{\mu}_{\mathbf{f}_r} - \boldsymbol{\mu}_{\mathbf{f}_g}\|^2 + \text{Tr}\left(\boldsymbol{\Sigma}_{\mathbf{f}_r} + \boldsymbol{\Sigma}_{\mathbf{f}_g} - 2\sqrt{\boldsymbol{\Sigma}_{\mathbf{f}_r}\boldsymbol{\Sigma}_{\mathbf{f}_g}}\right), \quad (46)$$

$$\text{KID}(\mathbf{f}_r, \mathbf{f}_g) = \text{MMD}^2(\mathbf{f}_r, \mathbf{f}_g), \quad (47)$$

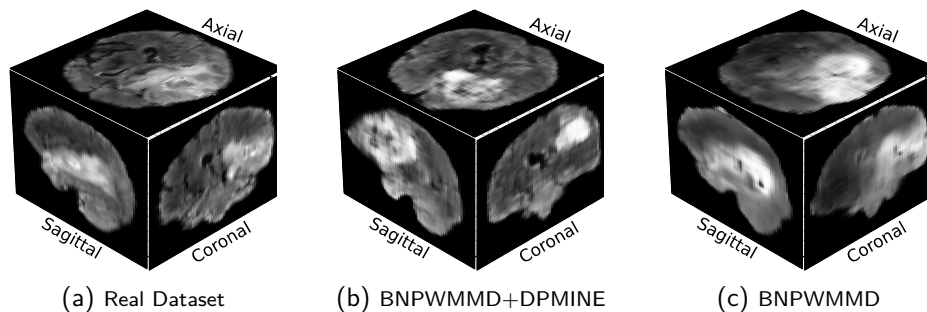


Figure 14: A 3D visualization of a real sample and a randomly generated sample using BNPWMMMD+DPMINE after 7500 epochs, compared with BNPWMMMD, in the BraTS18 example.

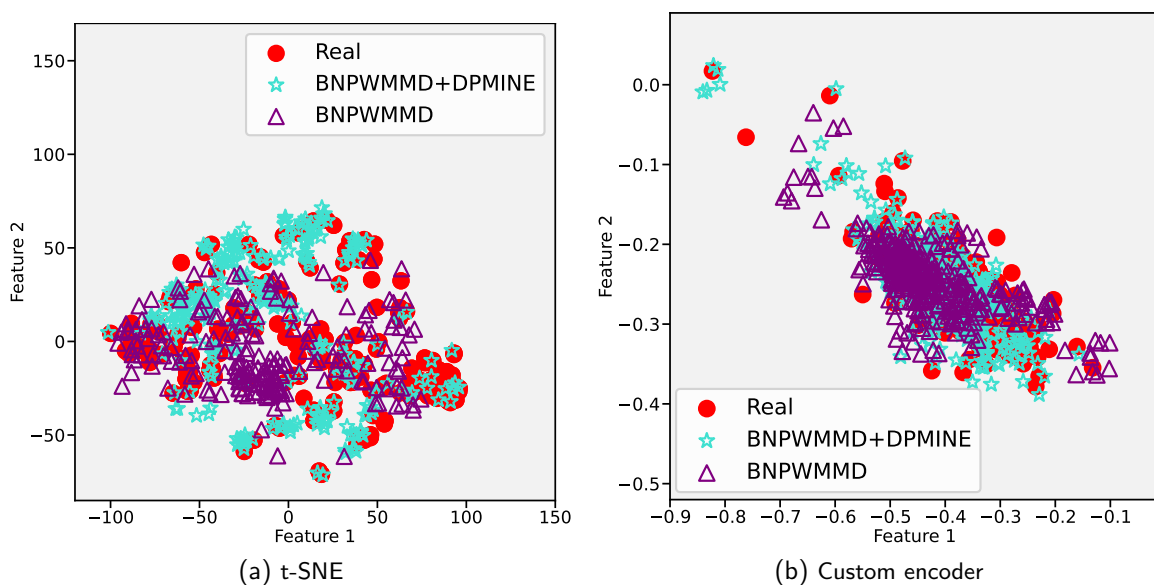


Figure 15: Scatter plots of two-dimensional feature representations in the BraTS18 example, comparing samples generated using BNPWMMMD+DPMINE with those generated using BNPWMMMD. Panel (a) shows the t-SNE embedding of the features, while panel (b) displays the features obtained from a custom encoder.

where  $\mu_{f_r}$  and  $\mu_{f_g}$ , and  $\Sigma_{f_r}$  and  $\Sigma_{f_g}$  are the mean vector and covariance matrix of  $f_r$  and  $f_g$ , respectively. The empirical MMD metric in Gretton et al. (2012) is used to compute  $\text{MMD}^2(\cdot, \cdot)$ , with a polynomial kernel  $k(f_r, f_g) = (0.5f_r^T f_g + 1)^\nu$  and degree  $\nu$ . The “Tr” denotes matrix trace, “T” denotes matrix transpose, and “ $\|\cdot\|$ ” denotes Euclidean norm. We use the provided code at [https://torchmetrics.readthedocs.io/en/v0.8.2/image/kernel\\_inception\\_distance.html](https://torchmetrics.readthedocs.io/en/v0.8.2/image/kernel_inception_distance.html) to compute KID with  $\nu = 3$ , and also use the provided code at <https://pytorch.org/ignite/generated/ignite.metrics.FID.html> to calculate FID.

Additionally, we calculate MS-SSIM using the available codes at [https://torchmetrics.readthedocs.io/en/v0.8.2/image/multi\\_scale\\_structural\\_similarity.html](https://torchmetrics.readthedocs.io/en/v0.8.2/image/multi_scale_structural_similarity.html).

### E.3.2 Network Setting

This section summarizes the network architectures used in all experiments. The encoder, generator, and discriminator architectures listed in Tables 4, 6, and 7 follow standard designs commonly adopted in 3D generative modeling and are consistent with those used in the  $\alpha$ -WGAN setting (Kwon et al., 2019). The BNPWMMMD model Fazeli-Asl & Zhang (2023) builds on these components and additionally includes a code

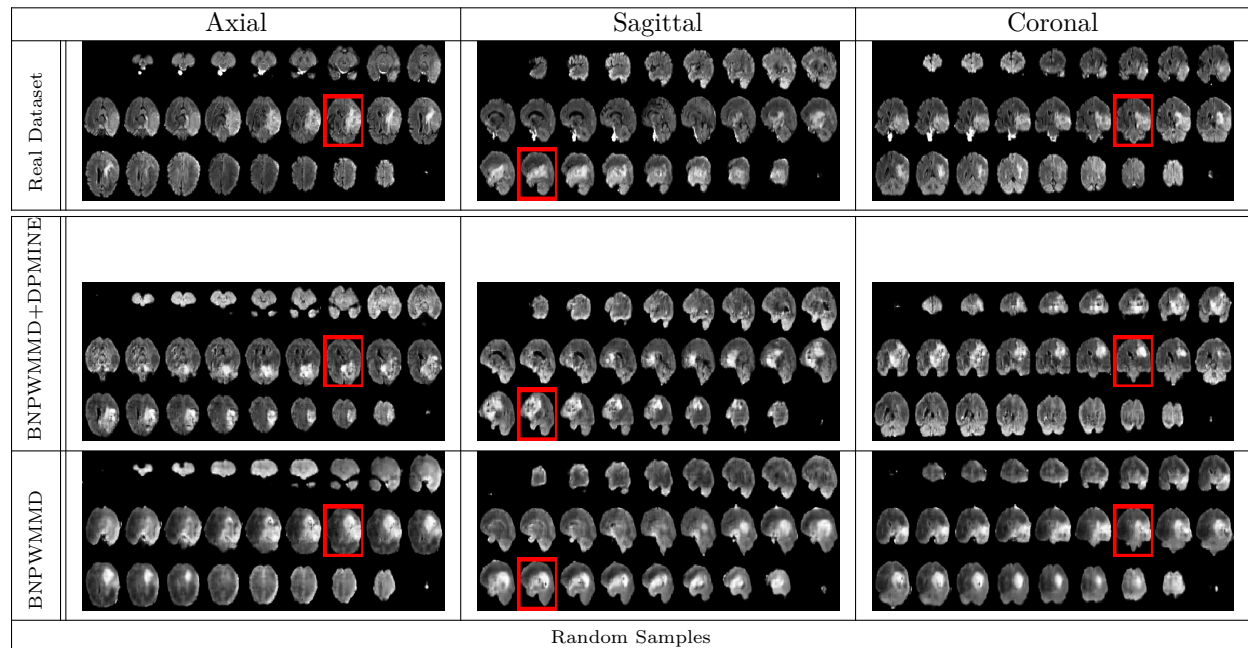


Figure 16: 27 slices along each axis of a 3D sample randomly generated using BNPWMMD+DPMINE after 7500 epochs, compared with BNPWMMD, for the BraTS18 example.

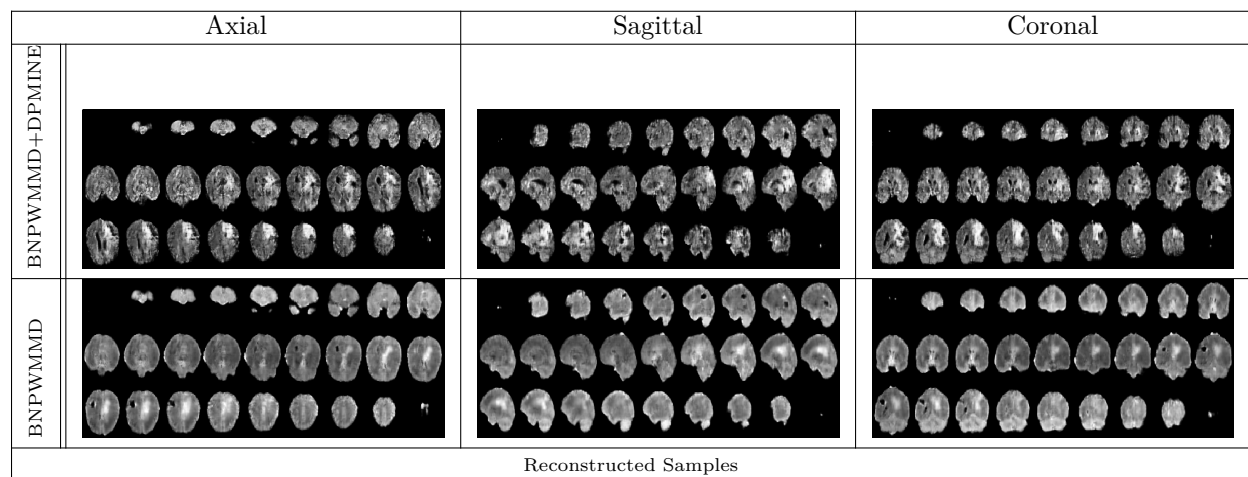


Figure 17: 27 slices along each axis of a 3D sample reconstructed using BNPWMMD+DPMINE after 7500 epochs, compared with BNPWMMD, for the BraTS18 example.

generator, detailed in Table 5. Figure 18 illustrates the overall workflow of the BNPWMMD+DPMINE framework.

Additionally, Table 8 presents a simple architecture used for updating the parameters  $\{T_\gamma\}_{\gamma \in \Gamma}$  in the DPMINE and MINE calculations Belghazi et al. (2018).

All architectures have been implemented in PyTorch using the Adam optimizer and a learning rate of 0.0002 on an NVIDIA Tesla V100-SXM2 with 4 GPUs, each with 32GB of RAM. Our code requires 48-72 hours to provide results for the COVID-19 and BraTS18 examples. For the coil experiments, our code takes about 4-6 hours.

Table 3: Comparison of statistical scores (mean  $\pm 2\sigma$ ) for the BraTS18 example ver 100 replications.

Evaluator		BNPWMMD + DPMINE	BNPWMMD
FID	Custom Encoder	<b>0.00197 <math>\pm</math> 0.00015</b>	0.00264 $\pm$ 0.00021
	t-SNE	<b>25.93170 <math>\pm</math> 14.68720</b>	34.13418 $\pm$ 16.00405
KID	Custom Encoder	<b>0.00038 <math>\pm</math> 0.00004</b>	0.00067 $\pm$ 0.00007
	t-SNE	<b>8.81000 <math>\pm</math> 7.18880</b>	9.63544 $\pm$ 8.95555
MMD	Custom Encoder	<b>0.00075 <math>\pm</math> 0.00078</b>	0.00092 $\pm$ 0.00087
	t-SNE	<b>0.19887 <math>\pm</math> 0.06760</b>	0.27312 $\pm$ 0.07383
MS-SSIM	-	<b>0.75007 <math>\pm</math> 0.10722</b>	0.74998 $\pm$ 0.12031

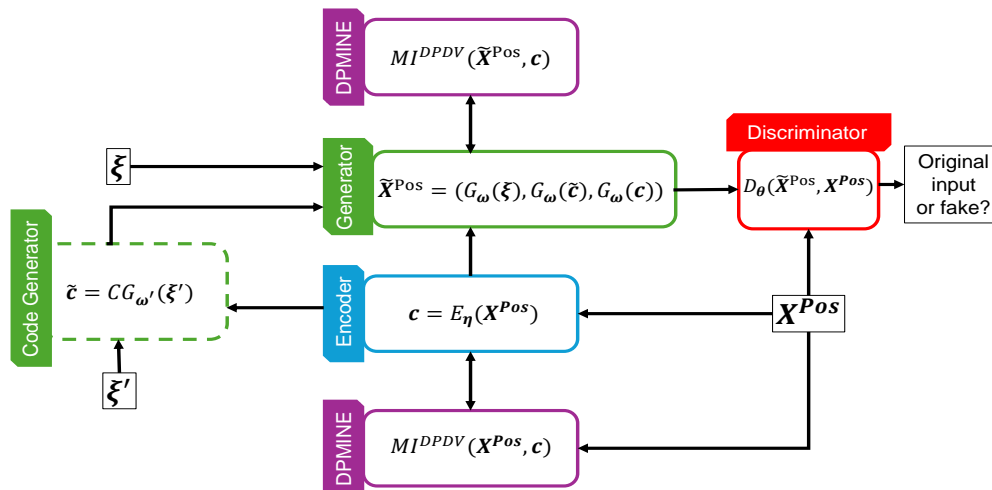


Figure 18: A general flowchart to illustrate the actual operation process of BNPWMMD+DPMINE.

## F Limitations

The BNP framework proposed in this paper assumes that the training data follows the IID assumption. However, this assumption restricts the model’s applicability to certain privacy-preserving techniques, such as federated learning, which often rely on non-IID assumptions at the local device level. While this limitation is acknowledged, it is beyond the scope of the current work to address it. In future research, we plan to extend the proposed BNP framework to incorporate non-IID assumptions and apply it to federated learning. This extension aims to mitigate issues such as slow convergence and unfair predictions that may arise when applying a global model to unseen data.

## G Broad Impacts and Safeguards

One potential positive impact of this paper is the advancement of artificial intelligence (AI) techniques in healthcare. The proposed model offers a solution to the challenge of accessing high-quality, diverse medical datasets for clinical decision-making. This can lead to improved accuracy and effectiveness in diagnosing and treating diseases, ultimately benefiting patients and healthcare providers. Additionally, the incorporation of a BNP procedure helps uncover underlying structures in the data and reduce overfitting, enhancing the reliability of the generated samples.

However, the use of AI in healthcare also raises concerns about data privacy and security. It is important to prioritize patient privacy and data protection by implementing strict regulations and protocols to ensure that patient data is anonymized and protected. Additionally, there is a risk of bias in the generated samples, which could lead to disparities in healthcare outcomes if the models are not properly validated and tested on diverse populations. Therefore, transparency and accountability in the development and deployment of these AI models are necessary. Rigorous testing and validation on diverse datasets should be conducted to mitigate

Table 4: Encoder: A 3D Convolutional Network Architecture for the COVID-19 Dataset used in BNPWMMD-GAN.

Layer	Input dimension	Output dimension	Kernel Size	Stride	Padding	Activation Function
Convolution	$1 \times 64 \times 64 \times 64$ (Data dimension)	$64 \times 32 \times 32 \times 32$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
Convolution	$64 \times 32 \times 32 \times 32$	$128 \times 16 \times 16 \times 16$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
BatchNorm	$128 \times 16 \times 16 \times 16$	$128 \times 16 \times 16 \times 16$	-	-	-	-
Convolution	$128 \times 16 \times 16 \times 16$	$256 \times 8 \times 8 \times 8$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
BatchNorm	$256 \times 8 \times 8 \times 8$	$256 \times 8 \times 8 \times 8$	-	-	-	-
Convolution	$256 \times 8 \times 8 \times 8$	$512 \times 4 \times 4 \times 4$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
BatchNorm	$512 \times 4 \times 4 \times 4$	$512 \times 4 \times 4 \times 4$	-	-	-	-
Convolution	$512 \times 4 \times 4 \times 4$	$1000 \times 1 \times 1 \times 1$	$4 \times 4 \times 4$	2	1	-

Table 5: Code-Generator: A 3D Convolutional Network Architecture for the COVID-19 Dataset used in BNPWMMD-GAN.

Layer	Input dimension	Output dimension	Kernel Size	Stride	Padding	Activation Function
Convolution	100 (Sub-latent dimension)	$16 \times 2 \times 5 \times 10$	$3 \times 3 \times 3$	1	1	ReLU
BatchNorm2	$16 \times 2 \times 5 \times 10$	$16 \times 2 \times 5 \times 10$	-	-	-	-
Max-pooling	$16 \times 2 \times 5 \times 10$	$16 \times 2 \times 5 \times 5$	$1 \times 1 \times 2$	$1 \times 1 \times 2$	-	-
Convolution	$16 \times 2 \times 5 \times 5$	$32 \times 2 \times 5 \times 5$	$3 \times 3 \times 3$	1	1	ReLU
BatchNorm2	$32 \times 2 \times 5 \times 5$	$32 \times 2 \times 5 \times 5$	-	-	-	-
Max-pooling	$32 \times 2 \times 5 \times 5$	$32 \times 2 \times 5 \times 2$	$1 \times 1 \times 2$	$1 \times 1 \times 2$	-	-
Convolution	$32 \times 2 \times 5 \times 2$	$64 \times 2 \times 5 \times 2$	$3 \times 3 \times 3$	1	1	ReLU
BatchNorm3	$64 \times 2 \times 5 \times 2$	$64 \times 2 \times 5 \times 2$	-	-	-	-
Max-pooling	$64 \times 2 \times 5 \times 2$	$64 \times 2 \times 5 \times 1$	$1 \times 1 \times 2$	$1 \times 1 \times 2$	-	-
Flatten	$64 \times 2 \times 5 \times 1$	640	-	-	-	-
Fully-connected	640	1000	-	-	-	-

bias and ensure fairness. Collaboration between healthcare professionals, AI researchers, and policymakers is essential to establish guidelines and regulations that prioritize patient privacy, data protection, and equitable healthcare outcomes.

Table 6: Generator: A 3D Convolutional Network Architecture for the COVID-19 Dataset used in BNPWMMD-GAN.

Layer	Input dimension	Output dimension	Kernel Size	Stride	Padding	Activation Function
Transposed convolution	1000 (Latent dimension)	$512 \times 4 \times 4 \times 4$	$4 \times 4 \times 4$	1	0	ReLU
BatchNorm	$512 \times 4 \times 4 \times 4$	$512 \times 4 \times 4 \times 4$	-	-	-	-
Upscale	$512 \times 4 \times 4 \times 4$	$512 \times 4 \times 4 \times 4$	-	-	-	-
Convolution	$512 \times 4 \times 4 \times 4$	$256 \times 8 \times 8 \times 8$	$3 \times 3 \times 3$	1	1	ReLU
BatchNorm	$256 \times 8 \times 8 \times 8$	$256 \times 8 \times 8 \times 8$	-	-	-	-
Upscale	$256 \times 8 \times 8 \times 8$	$256 \times 8 \times 8 \times 8$	-	-	-	-
Convolution	$256 \times 8 \times 8 \times 8$	$128 \times 16 \times 16 \times 16$	$3 \times 3 \times 3$	1	1	ReLU
BatchNorm	$128 \times 16 \times 16 \times 16$	$128 \times 16 \times 16 \times 16$	-	-	-	-
Upscale	$128 \times 16 \times 16 \times 16$	$128 \times 16 \times 16 \times 16$	-	-	-	-
Convolution	$128 \times 16 \times 16 \times 16$	$64 \times 32 \times 32 \times 32$	$3 \times 3 \times 3$	1	1	ReLU
BatchNorm	$64 \times 32 \times 32 \times 32$	$64 \times 32 \times 32 \times 32$	-	-	-	-
Upscale	$64 \times 32 \times 32 \times 32$	$64 \times 32 \times 32 \times 32$	-	-	-	-
Convolution	$64 \times 32 \times 32 \times 32$	$1 \times 64 \times 64 \times 64$	$3 \times 3 \times 3$	1	1	Tanh

Table 7: Discriminator: A 3D Convolutional Network Architecture for the COVID-19 Dataset used in BNPWMMD-GAN.

Layer	Input dimension	Output dimension	Kernel Size	Stride	Padding	Activation Function
Convolution	$1 \times 64 \times 64 \times 64$ (Data dimension)	$64 \times 32 \times 32 \times 32$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
Convolution	$64 \times 32 \times 32 \times 32$	$128 \times 16 \times 16 \times 16$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
BatchNorm	$128 \times 16 \times 16 \times 16$	$128 \times 16 \times 16 \times 16$	-	-	-	-
Convolution	$128 \times 16 \times 16 \times 16$	$256 \times 8 \times 8 \times 8$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
BatchNorm	$256 \times 8 \times 8 \times 8$	$256 \times 8 \times 8 \times 8$	-	-	-	-
Convolution	$256 \times 8 \times 8 \times 8$	$512 \times 4 \times 4 \times 4$	$4 \times 4 \times 4$	2	1	Leaky ReLU (negative slope=0.2)
BatchNorm	$512 \times 4 \times 4 \times 4$	$512 \times 4 \times 4 \times 4$	-	-	-	-
Convolution	$512 \times 4 \times 4 \times 4$	$1 \times 1 \times 1 \times 1$	$4 \times 4 \times 4$	2	1	Sigmoid

Table 8: Architecture of the shared  $T_\gamma$  network employed by DPMINE and all benchmark mutual information estimators.

Layer	Input Dimension	Output Dimension	Activation Function
concatenation	$1000 + 64^3$ (Latent dimension + Data dimension)	263144	-
Linear	263144	400	ReLU
Linear	400	400	ReLU
Linear	400	400	ReLU
Linear	400	1	ReLU