PARAMETER SPACE REPRESENTATION LEARNING ON MIXED-TYPE DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

A significant challenge in representation learning is to capture latent semantics in data mixing continuous, discrete, and even discretized observations (called mixedtype data), encountering issues like inconsistent discoveries and redundant modeling. Recently, Bayesian flow networks (BFNs) offer a unified strategy to represent such mixed-type data in the parameter space but cannot learn low-dimensional latent semantics since BFNs assume the size of parameters being the same as that of observations. This raises a new important question: how to learn latent semantics in parameter spaces rather than in observation spaces of mixed-type data? Accordingly, we propose a novel unified parameter space representation *learning* framework, ParamReL, which extracts progressive latent semantics in parameter spaces of mixed-type data. In ParamReL, a self-encoder learns latent semantics from intermediate parameters rather than observations. The learned semantics are then integrated into BFNs to efficiently learn unified representations of mixed-type data. Additionally, a reverse-sampling procedure can empower BFNs for tasks including input reconstruction and interpolation. Extensive experiments verify the effectiveness of ParamReL in learning parameter space representations for latent interpolation, disentanglement, time-varying conditional reconstruction, and conditional generation. The code is available at https: //anonymous.4open.science/r/ICLR25-F087/README.md.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

1 INTRODUCTION

This work explores a new important question: *How to learn latent semantics in parameter spaces rather than in observation spaces of mixed-type data comprising continuous, discrete, and even discretized observations*? We propose a novel unified *parameter space representation learning* framework that utilizes the parameter spaces rather than the observation spaces for mixed-type data.

Representation learning (Bengio et al., 2013) aims to discover low-dimensional latent semantics 037 from high-dimensional observations, widely applied in areas including computer vision (Li et al., 2023; Zhao et al., 2023a; Dong et al., 2023), and data analytics (Tonekaboni et al., 2022; Oublal et al., 2024). While the main focus has been on continuous-valued data (Kim & Mnih, 2018; Chen 040 et al., 2018; Meo et al., 2024), it is more challenging to uncover semantics in discrete (Austin et al., 041 2021; Chen et al., 2023) and even discretized (Van Den Oord et al., 2017; Razavi et al., 2019) 042 data. However, existing efforts often encounter issues like inconsistent discoveries and redundant 043 modeling (Zhou et al., 2023; Krishnan et al., 2018). Recently, Bayesian flow networks (BFNs) 044 (Graves et al., 2023; Song et al., 2024; Xue et al., 2024) emerged as a promising deep generative model. BFNs use multiple steps similar to diffusion models (Ho et al., 2020; Song et al., 2021) to refine parameters of an output distribution for reconstructing observations. Accordingly, BFNs offer 046 a unified strategy to handle mixed-type data while enabling fast sampling. However, they struggle 047 to capture low-dimensional latent semantics, raising the above open question. 048

Correspondingly, we propose a novel unified *Param*eter space *Re*presentation *L*earning framework,
 ParamReL, which leverages the multi-step generative learning of BFNs for representation learning
 on mixed-type data. ParamReL tackles this by performing representation learning in the parameter
 space to extract high-level latent semantics. The key insight lies in progressively self-encoding the
 intermediate parameters of BFNs, generating low-dimensional latent semantics step by step. Specifically, ParamReL adopts an architecture similar to BFNs but with two significant innovations: (1)

a *self-encoder* encodes intermediate parameters into lower-dimensional latent semantics, capturing gradual semantic changes throughout the multi-step generation process; and (2) a *conditional de-coder*, which conditions on latent semantics and intermediate parameters, and forms the parameters of an *output distribution* for simulating observations. Additionally, ParamReL involves *a reverse-sampling procedure* customized for tasks like image reconstruction and interpolation. Variational inference method is used in learning ParamReL, where mutual information is used to promote disentangled latent semantic learning, resulting in distinct and meaningful representations.

We evaluate ParamReL in learning meaningful high-level latent semantics from both discrete and
 continuous-valued observations on benchmark data. The sampling and reverse-sampling mecha nisms of ParamReL successfully perform tasks such as latent interpolation, disentanglement, time varying conditional reconstruction, and conditional generation. Notably, the self-encoder reveals
 progressive semantics throughout flow steps, enabling ParamReL to generate semantics with im proved clarity, while maintaining high quality of sample generation.

- 067
- 068
- 069

2 UNDERSTANDING BAYESIAN FLOW NETWORKS - AN ALTERNATIVE VIEW

Bayesian Flow Networks (BFNs) (Graves et al., 2023; Song et al., 2024; Xue et al., 2024) serve as
deep generative models with a primary objective to learn an output distribution for generating observations. The distribution's parameters are learned by a neural network, which takes the posterior
parameters of observations of inputs. Here, we try to understand BFNs from an alternative parameter
ter perspective since these (posterior) parameters play a key role in BFNs. BFNs involves concepts
such as input distribution, sender distribution and receiver distribution, to introduce BFNs, making
it less accessible to readers unfamiliar with BFNs. Interested readers may refer to Appendix A.1 and
(Graves et al., 2023) for the original illustrations.

Figure 1 shows T steps of training and sample generation in BFNs, similar to diffusion models (Ho et al., 2020; Song et al., 2021). To train BFNs, we minimize the divergence between the ground-truth data distribution and the evolving output distributions over T steps. At each step $t \in \{T, ..., 1\}$, an intermediate (posterior) parameter θ_t is first updated using a Bayesian update function $h(\cdot)$ as $\theta_t = h(\theta_{t+1}, \mathbf{x}_{t+1})$, where \mathbf{x}_{t+1} is the observation at step t+1. θ_t is then fed into a neural network $\psi(\cdot)$ to form the parameters of output distributions, i.e., a decoder $p_O(\mathbf{x}_t | \psi(\theta_t))$, for model training. After training, these intermediate output distributions can be employed to simulate observations during the sample generation process, replacing the actual observations at each step t.

By working in the parameter space, BFNs can uniformly model continuous, discrete, and discretized observations. For example, BFNs can use the mean of Gaussian distributions as parameter θ to model continuous data or use the event probabilities of categorical distributions as θ to study discrete data (see detailed settings for distributions in Table 2). However, BFNs cannot produce meaningful



Figure 1: Our alternative understanding of BFNs. Each step consists of a conditional decoder $p_O(\mathbf{x}_t | \psi(\boldsymbol{\theta}_t) \text{ (in blue rectangle) and a Bayesian update function } h(\cdot) \text{ (in peach rectangle). In training$ $BFNs, dashed arrows (between conditional decoder and <math>\{\mathbf{x}_t\}_{t=1}^T$) are non-existent as $\{\mathbf{x}_t\}_{t=1}^T$ refers to observations. The dashed arrows become solid for sample generation, representing the decoder generates \mathbf{x}_t in sample generation.

087

latent semantics capturing high-level concepts in the mixed-type observations, such as hair colors in portrait images.

3 PARAMREL: PARAMETER SPACE REPRESENTATION LEARNING

- Here, we explain the framework of ParamReL and its main design mechanisms.
- 116 3.1 The ParamReL Framework

The framework and workflow of ParamReL are shown in Figure 2. ParamReL leverages the parameter space for representation learning by extracting low-dimensional latent semantics from high-dimensional mixed-type data. Different from BFNs in approximating data distribution $p(\mathbf{x}_0)$, ParamReL learns the joint distribution over observation \mathbf{x}_0 and a series of latent semantics $\{\mathbf{z}_t\}_{t=1}^T$, with $|\mathbf{z}_t| \ll |\mathbf{x}_0|, \forall t \in \{1, \ldots, T\}$. That is, ParamReL seeks to reconstruct \mathbf{x}_0 while obtaining meaningful low-dimensional latent semantics $\{\mathbf{z}_t\}_{t=1}^T$.

Building on BFNs, ParamReL consists of four main components:

- (1) A self-encoder, conditioning on the intermediate (posterior) parameters θ_t to generate progressive latent semantics \mathbf{z}_t , described in Section 3.2.
- (2) A conditional decoder, using a neural network on latent semantics z_t and intermediate parameters θ_t to form the output distribution for subsequent steps, detailed in Section 3.3.
 - (3) A sampling and reverse-sampling process, facilitating tasks such as image reconstruction and interpolation, outlined in Section 3.4.
 - (4) A training and testing procedure, as discussed in Section 3.5, optimizing latent semantics z_t and ensuring effective model generalization.

Together, ParamReL forms a robust framework to capture and utilize latent semantics and to improve the performance of tasks including unconditional image generation and reconstruction.

3.2 PARAMETER ENCODING THROUGH A SELF-ENCODER

The *self-encoder*, denoted as $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$, progressively encodes intermediate parameters $\boldsymbol{\theta}_t$ into low-dimensional latent semantics \mathbf{z}_t , which facilitates representation learning from high-dimensional, mixed-type data at each step t. (Baranchuk et al., 2021) has shown that upsampling





layers from a U-Net in pre-trained diffusion models (Rombach et al., 2022) may capture meaningful semantic information. Inspiring from this discovery and in training ParamReL, we adopt approaches similar to (Luo et al., 2024) to parameterize $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ (see Appendix C.1 for more details). Through $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$, the intermediate parameter $\boldsymbol{\theta}_t$ effectively encodes itself into \mathbf{z}_t , together they form $\psi(\boldsymbol{\theta}_t, \mathbf{z}_t)$ for the output distribution.

Ideally, the latent semantics z_t should provide low-dimensional semantics distinct from the intermediate parameters θ_t in BFNs but without compromising the data reconstruction process. To learn high-quality latent semantics, a smooth, learnable latent space is necessary, which is ensured by integrating the prior distribution $p(z_t)$ into a robust probabilistic framework, allowing efficient sampling of x_0 . For simplicity and efficiency, we assume $p(z_t)$ follows a Gaussian distribution.

 $q_{\phi}(\mathbf{z}_t|\boldsymbol{\theta}_t, t)$ differs from traditional auto-encoders $q_{\phi}(\mathbf{z}|\mathbf{x}_0)$ in two key aspects:

- $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ is conditioned on the intermediate parameter $\boldsymbol{\theta}_t$, rather than being conditioned on \mathbf{x}_0 . This summarizes information from all previous steps to enable generating latent semantic \mathbf{z}_t through all the T steps.
- The self-encoder generates a step-wise semantic \mathbf{z}_t , which is tailored to the dynamic behavior of variables over time t. This series of latent semantics $\{\mathbf{z}_t\}_{t=1}^T$ are expected to exhibit progressive semantic behaviors (such as gradual changes in age, smile, or skin color) throughout the generation process (as illustrated in the right panel of Figure 13).

When observations \mathbf{x}_0 are unavailable, e.g. sample generation tasks, it is also worth noting that directly using regular auto-encoders like $q_{\phi}(\mathbf{z}|\mathbf{x}_0)$ to generate latent semantics is infeasible. They may require an additional module to generate latent semantics (Preechakul et al., 2022), while training such modules would introduce computational overhead. However, in their case, not using autoencoders $q_{\phi}(\mathbf{z}|\mathbf{x}_0)$ would lead to inefficient resource use.

186 187

188

194 195

173 174

175

176

177

178 179

3.3 CONDITIONAL DECODER

The conditional decoder refers to the output distribution $p_O(\mathbf{x}_t | \psi(\boldsymbol{\theta}_t, \mathbf{z}_t))$ which conditions on latent semantics \mathbf{z}_t and intermediate parameter $\boldsymbol{\theta}_t$ to simulate \mathbf{x}_t . The condition $\psi(\boldsymbol{\theta}_t, \mathbf{z}_t)$ explicitly incorporates \mathbf{z}_t as part of its conditioning mechanism. Following the settings in diffusion models (Ho et al., 2020; Song et al., 2021), we use the U-Net architecture with the Cross-Attention in each layer specified as

Cross-Attention
$$(\boldsymbol{\theta}_t, \mathbf{z}_t) = \operatorname{softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d}})\mathbf{V}$$
, where $\mathbf{Q} = \mathbf{W}^Q \boldsymbol{\theta}_t, \mathbf{K} = \mathbf{W}^K \mathbf{z}_t, \mathbf{V} = \mathbf{W}^V \mathbf{z}_t$

where $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V$ are the query, key and value weight matrix, respectively. See the detailed U-Net architecture in Appendix C.2.

Since z_t works together with the corresponding intermediate parameter θ_t , it is expected that z_t aligns well with the progressively structured parameter θ_t . Lower-level intermediate latent x_t (such as hair texture) is progressively incorporated. The proposed self-encoder works consistently with the conditional decoder here as both work on θ_t , see Figure 6 (b).

203 204

210

3.4 SAMPLING AND REVERSE-SAMPLING PROCESSES

After training ParamReL, the sampling and reverse-sampling processes play a crucial role in generating and reconstructing data, which is essential for tasks such as image generation and interpolation. Generating samples begins with an initial guess of the intermediate parameters θ_{T+1} . From θ_{T+1} , this sampling process sequentially generates $\mathbf{x}_T, \mathbf{x}_{T-1}, \dots, \mathbf{x}_0$. Specifically, given the parameter θ_t at each step *t*, we have:

$$\mathbf{z}_{t} \sim q_{\boldsymbol{\phi}}(\mathbf{z}_{t}|\boldsymbol{\theta}_{t}, t), \ \mathbf{x}_{t} \sim p_{O}(\mathbf{x}_{t}|\boldsymbol{\psi}(\boldsymbol{\theta}_{t}, \mathbf{z}_{t})), \ \boldsymbol{\theta}_{t-1} = h(\boldsymbol{\theta}_{t}, \mathbf{x}_{t}).$$
(1)

We use the trained encoder $q_{\phi}(\mathbf{z}_t|\boldsymbol{\theta}_t, t)$ to replace the prior $p(\mathbf{z}_t)$ of \mathbf{z}_t for improving the sampling quality. After $\boldsymbol{\theta}_0$ is obtained, a sample can be generated as $\mathbf{z}_0 \sim q_{\phi}(\mathbf{z}_0|\boldsymbol{\theta}_0, 0), \mathbf{x}_0 \sim p_{O}(\mathbf{x}_0|\psi(\boldsymbol{\theta}_0, \mathbf{z}_0))$.

However, the reverse-sampling process, which transits the observation \mathbf{x}_0 through the intermediate latents $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{T-1}$ until \mathbf{x}_T , is not as straightforward as the sampling procedure. Without a



Figure 3: Reverse-sampling process in BFNs.

clearly defined reverse-sampling process, it would be challenging to perform tasks such as image reconstruction and interpolation. In fact, by taking the inverse of the Bayesian update function $h(\cdot)$ as $\theta_t = h^{-1}(\theta_{t-1}, \mathbf{x}_{t-1})$, the intermediate latent \mathbf{x}_{t-1} can transit to \mathbf{x}_t as:

$$\boldsymbol{\theta}_t = h^{-1}(\boldsymbol{\theta}_{t-1}, \mathbf{x}_{t-1}), \ \mathbf{z}_t \sim q_{\boldsymbol{\phi}}(\mathbf{z}_t | \boldsymbol{\theta}_t, t), \ \mathbf{x}_t \sim p_{\mathrm{O}}(\mathbf{x}_t | \boldsymbol{\psi}(\boldsymbol{\theta}_t, \mathbf{z}_t)).$$
(2)

Given the straightforward definition of Bayesian update function $h(\cdot)$, its inverse operation is generally easy to derive. The details of such results can be found in Figure 14. Furthermore, this developed reverse-sampling process can be naturally extended to BFNs. Transiting \mathbf{x}_{t-1} to \mathbf{x}_t at time t can be performed as $\boldsymbol{\theta}_t = h^{-1}(\boldsymbol{\theta}_{t-1}, \mathbf{x}_{t-1})$, with \mathbf{x}_t sampled as $\mathbf{x}_t \sim p_O(\mathbf{x}_t | \boldsymbol{\psi}(\boldsymbol{\theta}_t))$. With this approach, BFNs can effectively perform tasks like image reconstruction and interpolation, which were difficult or even impossible by previous BFNs models. Figure 3 shows the reverse-sampling process of BFNs. The ParamReL version is provided in Figure 7 in Appendix A.

243 244

227 228

229 230

231

232

233 234 235

3.5 TRAINING AND TEST WITH PARAMREL

Here, we outline the process of training and testing ParamReL by focusing on optimizing Param-ReL to learn meaningful latent semantics while ensuring effective reconstruction of observations. The training process involves variational inference to approximate the joint distribution of latent variables, and a mutual information term is integrated into improving the quality of learned latent semantics by strengthening the relationship between intermediate parameters and latent semantics.

Variational Inference for Intractable Joint Distribution In ParamReL, the joint distribution over \mathbf{x}_0 , intermediate latents $\{\mathbf{x}_t\}_{t=1}^T$ and latent semantics $\{\mathbf{z}_t\}_{t=1}^T$ can be defined as $p(\mathbf{x}_0, \{\mathbf{x}_t\}_{t=1}^T, \{\mathbf{z}_t\}_{t=1}^T|-) = p_O(\mathbf{x}_0|\psi(\theta_0, \mathbf{z}_0)) \cdot \prod_{t=1}^T [p(\mathbf{z}_t)\mathbb{E}_{p_O(\mathbf{x}_t|\psi(\theta_t, \mathbf{z}_t))}[p_S(\mathbf{x}_{t-1}|\mathbf{x}_t)]]$, where the output distribution $p_O(\mathbf{x}_0|\psi(\theta_0, \mathbf{z}_0))$ at step 0 is used to model observation \mathbf{x}_0 , and $\mathbb{E}_{p_O(\mathbf{x}_t|\psi(\theta_t, \mathbf{z}_t))}[p_S(\mathbf{x}_{t-1}|\mathbf{x}_t)]$ follows the definition of BFNs to model intermediate latent \mathbf{x}_{t-1} , and $p_S(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is a noisy distribution of \mathbf{x}_t .

257 With $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ defined as the encoder for \mathbf{z}_t and $p_{\mathrm{S}}(\mathbf{x}_{t-1} | \mathbf{x}_t)$ defined as the variational distribu-258 tion for \mathbf{x}_{t-1} , the evidence lower bound (ELBO) on the marginal log-likelihood of observation \mathbf{x}_0 is 259 (see the full derivation in Appendix B):

$$\log p(\mathbf{x}_{0}) \geq -\sum_{t=1}^{T} \mathbb{E}_{p_{\mathrm{F}}(\boldsymbol{\theta}_{t}|-)} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{t}|\boldsymbol{\theta}_{t},t)} \left\{ \mathrm{KL} \left[p_{\mathrm{S}}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}\right) \parallel \mathbb{E}_{p_{\mathrm{O}}(\mathbf{x}_{t}|\psi(\boldsymbol{\theta}_{t},\mathbf{z}_{t}))} \left[p_{\mathrm{S}}(\mathbf{x}_{t-1}|\mathbf{x}_{t}) \right] \right] - \mathrm{KL} \left[q_{\boldsymbol{\phi}}\left(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t},t\right) \parallel p(\mathbf{z}_{t}) \right] + \mathbb{E}_{p_{F}(\boldsymbol{\theta}_{0}|-)q_{\boldsymbol{\phi}}(\mathbf{z}_{0}|\boldsymbol{\theta}_{0},0)} \left[\ln p_{\mathrm{O}}(\mathbf{x}_{0}|\psi(\boldsymbol{\theta}_{0},\mathbf{z}_{0})) \right] := \mathrm{ELBO.} \quad (3)$$

264

260 261 262

Maximizing ELBO is equivalent to performing amortized inference (Kingma & Welling, 2014) through encoders $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ and learning likelihood function through decoders (Zhao et al., 2019). When the encodable posterior $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ is used to infer high-level semantics \mathbf{z}_t , those intermediate latents $\{\mathbf{x}_t\}_{t=1}^T$ contain low-level information in generating the observations. In ParamReL, the parameters of the output distribution are learned through iteratively proceeding the Bayesian updating functions and a learned noise model $\psi(\boldsymbol{\theta}, \mathbf{z})$ parameterized by neural networks ψ . 270 Mutual Information Regularization Ideally, during the training phase, we want to acquire the 271 latent semantic \mathbf{z}_t by the self-encoder $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ and achieve high-quality reconstruction $\widehat{\mathbf{x}_0}$ by the 272 decoder (i.e., the output distribution $p_O(\mathbf{x}_0|\psi(\boldsymbol{\theta}_0, \mathbf{z}_0)))$). However, there exists a trade-off between 273 inference and learning (Shao et al., 2020; Wu et al., 2024) coherent in optimizing the ELBO in Eq. (3). 274 In most cases, optimizing ELBO favours fitting likelihood rather than inference (Zhao et al., 2019). Based on the rate-distortion theory (Alemi et al., 2018; Bae et al., 2023), the rate, represented by the 275 KL divergence term constrained by the encoders, compresses sufficient information to minimize the 276 distortion, or reconstruction error, while simultaneously limiting the informativeness to promote a 277 smooth latent space. 278

To remedy the insufficient representation learning during the inference stage, we want to increase the dependence between intermediate parameters θ_t and latent semantics \mathbf{z}_t by maximizing their mutual information $I(\theta_t, \mathbf{z}_t)$. We can rewrite the tractable learning object in ParamReL by adding the mutual information maximization term as $ELBO_+ = ELBO + \frac{\gamma}{T} \sum_t I_q(\theta_t; \mathbf{z}_t)$, where γ is the trade-off parameter. Considering that we cannot optimize this object directly, we can rewrite it by factorizing the rate term into mutual information and total correlation (TC), see details in Appendix B.

4 RELATED WORK

286 287

285

Recent advances have demonstrated that diffusion models (Ho et al., 2020; Song et al., 2021) are capable of generating high-quality data. Nonetheless, compared to the autoencoder framework, the intermediate outputs in diffusion stages are high-dimensional and lack smoothness, making them unsuitable for representation learning. Contemporary research focuses on encoding a conditional latent space to acquire low-dimensional semantic representations. However, those observations-based models (Preechakul et al., 2022; Wang et al., 2023), such as VAEs and diffusion models, exhibit limitations when applied to discrete data.

295 Deep hierarchical VAEs have seen progress in capturing latent dependence structures for encoding an expressive posterior, statistically or semantically. VQVAE-based (Van Den Oord et al., 2017; 296 Razavi et al., 2019) models have local-to-global features-based explanatory hierarchies at the image 297 level, forming a codebook-based discrete posterior. In (Sønderby et al., 2016; Tomczak & Welling, 298 2018), recursive latent structures in multi-layer networks form an aggregated posterior. NVAE (Vah-299 dat & Kautz, 2020) demonstrates that depth-wise hierarchies encoded by residual networks can ap-300 proximate the posterior precisely despite using shallow networks. Unlike the observation-based 301 encoder, where the information flow between input and latent is maximized in encoding-decoding 302 pipelines in the sample space, ParamReL uses progressive encoders in the parameter space to capture 303 the dynamic semantics. 304

Pre-trained diffusion models (Rombach et al., 2022), (Baranchuk et al., 2021) have shown that the 305 upsampling features from a U-Net can capture semantic information useful for downstream tasks. 306 This discovery has sparked increasing research in leveraging these upsampling features of pre-307 trained diffusion models across various applications, including classification (Xiang et al., 2023; 308 Mukhopadhyay et al., 2023), semantic segmentation (Baranchuk et al., 2021; Zhao et al., 2023b), 309 panoptic segmentation (Xu et al., 2023), semantic correspondence (Tang et al., 2023; Zhang et al., 310 2024; Luo et al., 2024; Hedlin et al., 2024), and image editing (Tumanyan et al., 2023; Hertz et al., 311 2022). In most of these approaches, identifying the optimal denoising step and upsampling layer is crucial for achieving high predictive performance. These approaches do not suggest fundamental 312 changes to model architectures or training methodologies, leaving the specific architectural com-313 ponents and techniques for learning useful semantic representations unclear. ParamReL uses these 314 discoveries to construct efficient self-encoders. 315

316

317 5 EXPERIMENTS 318

We present two ParamReL variants operating in different parameter spaces: ParamReLd for discrete input distributions (Section 5.2), and ParamReLc for continuous input distributions (Section 5.3), respectively. We evaluate the representation learning capabilities of ParamReL in three reconstructionbased tasks: latent interpolation, disentanglement, and time-varying conditional reconstruction. Additionally, we evaluate the model for unconditional generation, where samples are generated *only from the decoder* using a given prior.

324 5.1 EVALUATION SETUP 325

326 We conduct a two-fold comparison to evaluate the performance of ParamReL variants. Firstly, we 327 compare our parameter-based models (ParamReLc and ParamReLd) with established sample-based representation learning baselines, including AE and VAE-based models such as β -VAE (Higgins 328 et al., 2017), infoVAE (Zhao et al., 2019), and diffusion-based models such as DiffAE (Preechakul et al., 2022) and InfoDiffusion (Wang et al., 2023). These models represent key advancements in 330 the field: β -VAE introduce disentanglement into VAE, infoVAE incorporates MMD for balancing 331 generation and representation, while DiffAE and InfoDiffusion explore the integration of AEs and 332 VAEs into diffusion models to learn encodable latents and disentangled representations, respectively. 333 Secondly, we compare the performance of ParamReLc and ParamReLd across various input distribu-334 tions for continuous and discrete data, respectively. The discrete datasets include binarized versions 335 of MNIST (bMNIST) (Deng, 2012), FashionMNIST (bFashionMNIST) (Xiao et al., 2017), while 336 the continuous datasets include CelebA (Liu et al., 2015), CIFAR10 (Krizhevsky & Hinton, 2009), 337 and Shapes3D (Burgess & Kim, 2018)¹. The detailed hyperparameter choices and experimental 338 configurations for each dataset are provided in Appendix C.3. This comparison allows for a detailed 339 examination of how different parameter space assumptions impact the representation learning of discrete and continuous data. 340

341

343

342 5.2 SEMANTIC REPRESENTATION OF DISCRETE DATA BY PARAMRELD

Here, we measure the quality of the learned latent semantics z_0 through the downstream classifi-344 cation tasks. Since z_0 locates at step 0, they should be *general* and *transferable* (Franceschi et al., 345 2019). Various datasets by deep classifiers are assessed to ensure their universality. Specifically, 346 following the approach in Xiao & Bamler (2023), we train a classifier on labeled test sets for each 347 ParamReL model. We allocate 80% of the dataset for training a classifier and reserve the remain-348 ing 20% for test purposes. The performance on the test set is evaluated based on AUROC. This 349 process is conducted in a 5-fold cross-validation manner, with the results reported as mean \pm one 350 standard deviation. The results are shown in Figure 4 (a). Higher AUROC suggests that the learned 351 latent semantics z_0 contain more information about data. In addition to assessing the representation 352 quality, we also compare the image reconstruction ability against baselines. From the FID values 353 in Figure 4 (a) and Figures 11, 12 in Appendix E.3, we can conclude that VAE-based models still produce blurry reconstructions, while diffusion-based and parameter-based models can build near-354 exact reconstructions. Refer to Figure 11 and Figure 12 in Appendix E.3 for the generated binary 355 images. 356

- 357
- 358 359

5.3 SEMANTIC REPRESENTATION OF CONTINUOUS DATA BY PARAMRELC

On continuous data, we evaluate ParamReLc for conditional generation, conditional reconstruction, latent interpolation, and disentanglement.

High-level Representation Learning for Conditional Generation Figure 13 (a) in Appendix E.2 demonstrates that high-level semantic information is captured by the learned latent semantics $\{\mathbf{z}_t\}_{t=1}^T$ for image generation. This is illustrated by a set of latent-sample pairs $\{\mathbf{z}_t\}_{t=1}^T, \mathbf{x}_T^{i,j} >$, where $\{\mathbf{z}_t^i\}_{t=1}^T$ are obtained by reverse-sampling from the *i*-th input image through the trained ParamReL, and $\mathbf{x}_T^{i,j}$ is the *j*-th sample from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ corresponding to the *i*-th input image. Concurrently, the low-level information, such as local attributes in images (e.g., Narrow_Eyes, Mouth_Slightly_Open, Blond_Hair), are determined by $\mathbf{x}_T^{i,j}$.

Time-varying Representation Learning for Conditional Reconstruction We design a *new* timevarying reconstruction task to evaluate the effectiveness of the progressive latent semantics learned by the self-encoder. A latent-sample pair $\langle \{\mathbf{z}_t^{\text{fixed}}\}_{t=1}^T, \mathbf{x}_T^{\text{fixed}} \rangle$ is first obtained by apply the trained ParamReL's reverse-sampling process on an image. Then, we use the latent semantics at step *t** to replace other steps' ones and "reconstruct" the image as $\mathbf{x}_t \sim p_O(\mathbf{x}_t | \boldsymbol{\psi}(\boldsymbol{\theta}_t, \mathbf{z}_{t^*}^{\text{fixed}})), \boldsymbol{\theta}_{t-1} =$ $h(\boldsymbol{\theta}_t, \mathbf{x}_t), \forall t = T, ..., 1$. In that case, the attributes vary due to the semantics evolution encoded by time-specific latent. Refer to Figure 4 (b), and Figure 13 (b) in Appendix E.2 for more explanation.

¹For the discrete version, continuous data (k-bit images) can be discretized into 2^k bins by dividing the data range [-1, 1] into k intervals, each of length 2/k.



Figure 4: Quantitative representation learning comparison over generative models on discrete data (a). ParamReL demonstrates competitive performance in capturing semantic information for classification, achieving approximately 0.84 AUROC for bFashionMNIST and 0.91 for bMNIST. Additionally, it shows robust generative capabilities, with FID values ranging from 0.5 to 0.6 for bMNIST

and around 5 for bFashionMNIST. Among the ParamReL-based models, ParamReLd with a categorical distribution is particularly effective in modelling discrete data distributions, yielding lower FID values of 0.5 for bMNIST and 4.2 for bFashionMNIST. As shown in (b), the learned semantics exhibit progressive, time-varying changes. By varying time encodes at 200, 300, 400 time steps, more attributes will be influenced in the reconstruction stage: the Wavy_hair, Brown_hair, Arched_Eyebrows attributes in the first line, the Double_Chin, Mustache, Goatee attributes in the second line and the Young, High_Cheekbones, Arched_Eyebrows attributes in the third line. Notations: [AUROC, FID]; [(•, bMNIST), (■, bFashionMNIST)]; [(-, ParamReLd),(-·-, ParamReLc)].

Table 1: Comparison of representation learning algorithms on continuous data by disentanglement performance (mean \pm std) and classification. The quantitative results for each algorithm are averaged over five trials. Notations: Modeling on data space \mathcal{D} , parameter space \mathcal{P} . Prior distributions: Gaussian *g*, Categorical *c*, Delta *d*. \uparrow : higher better, \downarrow : lower better. Color: **Top-1**, Top-2.

Prio	r Prior	Methods	CelebA				Shapes3D		CIFAR-10	
on	type	Methous	$\mathcal{TAD}\uparrow$	$ATTRS\uparrow$	$\mathcal{FID}\downarrow$	$AUROC\uparrow$	$\mathcal{DCI}\uparrow$	$AUROC\uparrow$	$\mathcal{FID}\downarrow$	$AUROC\uparrow$
	-	AE	0.042 ± 0.004	1.0 ± 0.0	90.4±1.8	0.759 ±0.003	0.219 ±0.001	0.796±0.007	169.4±2.4	0.721±0.001
	g	VAE Kingma & Welling (2014)	0.000 ± 0.000	0.0 ± 0.0	94.3±2.8	0.770 ±0.002	0.276 ±0.001	0.799±0.002	177.2±3.2	0.743±0.002
Ð	g	β -VAE Burgess et al. (2017)	0.088 ±0.051	1.6 ±0.8	99.8±2.4	0.699 ±0.001	0.281 ±0.001	0.801±0.001	183.3±3.1	0.769±0.003
D	g	InfoVAE Zhao et al. (2019)	0.000 ± 0.000	0.0 ± 0.0	77.8±1.6	0.757 ±0.003	0.134 ±0.001	0.829±0.003	160.7±2.5	0.814±0.006
	g	DiffAE Preechakul et al. (2022)	0.155 ±0.010	2.0 ± 0.0	22.7±2.1	0.799 ±0.002	0.196 ±0.001	0.899±0.001	32.1±1.1	0.859±0.002
	g	InfoDiffusion Wang et al. (2023)	0.299 ± 0.006	3.0 ± 0.0	23.8±1.6	0.848 ± 0.001	0.342 ± 0.002	0.882 ± 0.001	32.4±1.8	0.886 ± 0.004
D	с	ParamReL $(\gamma = 1, \lambda = 0.01)$	0.261 ±0.01	5.0 ±0.0	22.6±1.2	0.846 ±0.009	0.477 ±0.002	0.901±0.007	31.8±1.1	0.892±0.004
,	d	ParamReL $(\gamma = 0.9, \lambda = 0.01)$	0.302 ±0.005	<u>4.0 ±0.0</u>	<u>22.1±1.6</u>	0.850 ±0.006	0.567 ±0.005	0.902±0.001	<u>31.2±1.1</u>	0.901±0.001
	d	$ParamReL (\gamma = 1, \lambda = 0.01)$	0.368 ±0.005	3.0 ± 0.0	21.6±1.1	0.865±0.004	<u>0.485 ±0.009</u>	0.931±0.001	31.1±1.1	0.911±0.002

Smooth Representation Learning for Latent Interpolation Latent space interpolation (Goodfellow et al., 2014; Higgins et al., 2017) is commonly used to validate the smoothness, continuity, and semantic coherence of the learned latent semantics in generative models. Typically, two samples are embedded into the latent space, and interpolating between the latent variables generates interpolated representations. The reconstructed outputs produced by the sampling process reveal the semantic richness of the latent space. Demonstration of the image interpolation is detailed in Appendix D.1.

424 As shown in Figure 14 in Appendix E.3, ParamReL achieves near-exact reconstruction, in contrast 425 to the downgraded performance of VAE variants such as (a) vanilla VAE, and (b) β -VAE. Compared 426 with diffusion models (c) DiffAE and (d) InfoDiffusion, ParamReL characterizes a smoother and 427 more consistent latent space with high-quality samples.

Disentanglement We perform latent traversals on the FFHQ and CelebA datasets to evaluate the disentanglement properties of our trained ParamReL, as illustrated in Figure 5 and Figure 15 in Appendix E.3. In this process, we modify one dimension of the learned latent semantics $\{z_t\}_{t=1}^T$ each step, and replace it with M evenly distributed numbers within a standardized range (e.g., -3to +3), while keeping the other dimensions fixed. After decoding these adjusted latent semantics, (a) Mustache (a) Mustache (b) Brown_Hair (c) Brown_Hair (c) Brown_Hair

(c) Eyeglasses

Figure 5: Disentanglement of ParamReL on FFHQ-128. The interpretable traversal directions are displayed by traversing the encodings ranging from [-3, 3].

we evaluate the generated samples for changes in specific attributes. Successful disentanglement 451 is verified when manipulating one single dimension alters only one distinguishable attribute, such 452 as age, while leaving all other attributes unchanged. As shown in Figure 5 and Figure 15 in the 453 Appendix, ParamReL effectively isolates and controls individual data attributes in both FFHQ and 454 CelebA. For example, on FFHQ, manipulating latent dimensions controls attributes like Mustache, 455 Brown Hair, and Eyeqlasses, while other attributes remain constant. Similarly, on CelebA, 456 attributes such as Smiling, Pale Skin, and Big Nose are independently manipulated without 457 affecting others. 458

To provide a thorough and unbiased quantitative assessment of disentanglement, we utilize two met-459 rics: 1) Disentanglement, Completeness, and Informativeness (DCI) (Eastwood & Williams, 2018), 460 which is a prediction-based indicator; and 2) Total AUROC Difference (TAD) (Yeats et al., 2022), 461 an intervention-based criterion. Additionally, we report the generation quality in Appendix E.3 and 462 conclude that ParamReL achieves near-exact reconstruction on CelebA (Figure 16 (a)), Shapes3D 463 (Figure 16 (b)), and CIFAR-10 (Figure 16 (c)). Both the qualitative latent traversal results and 464 the quantitative disentanglement metrics show that ParamReL effectively learns disentangled rep-465 resentations, with visual traversals closely aligning with the attributes that the latent semantics are 466 intended to capture.

467 468

469

445 446 447

448

449 450

6 CONCLUSION AND LIMITATIONS

470 In this work, we introduce ParamReL, a novel unified parameter space representation learning 471 framework, as a unified strategy to handle continuous, discrete and even discretized data. Unlike 472 traditional encoder methods that map observations into static latent semantics, ParamReL employs a 473 self-encoder to derive progressively structured latent semantics from intermediate parameters at each 474 step of the generation process. This allows for more effective representation learning across different 475 data types. Our experiments on tasks including latent interpolation, disentanglement, time-varying 476 conditional reconstruction, and conditional generation validate the effectiveness of ParamReL. The results demonstrate its superior ability to extract meaningful high-level semantics, leading to unified 477 representations and a clear semantic understanding of the underlying data. 478

While ParamReL shows promising results, our experiments reveal areas for potential expansion. (1)
The precision variables, which play a key role in the sampling process, could be further optimized to
reduce computational time and to improve efficiency. This was observed during the sampling stages
where slight inefficiencies in parameter updates are detected. (2) We noticed that employing a standard U-Net architecture without pre-training may limit the performance of ParamReL, particularly
in tasks involving complex data. Therefore, exploring the integration of a pre-trained U-Net model
into ParamReL could provide a significant boost in accuracy and representation quality. We will
investigate these in the future work.

486 REFERENCES 487

492

521

524

525

- Alexander Alemi, Ben Poole, Ian Fischer, Joshua Dillon, Rif A Saurous, and Kevin Murphy. Fixing 488 a Broken ELBO. ICML, 2018. 489
- 490 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured 491 Denoising Diffusion Models in Discrete State-Spaces. NeurIPS, 2021.
- Juhan Bae, Michael R Zhang, Michael Ruan, Eric Wang, So Hasegawa, Jimmy Ba, and Roger 493 Grosse. Multi-rate VAE: Train Once, Get the Full Rate-Distortion Curve. ICLR, 2023. 494
- 495 Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khrulkov, and Artem Babenko. Label-496 Efficient Segmentation with Diffusion Models. arXiv preprint arXiv:2112.03126, 2021.
- 497 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation Learning: A Review and New 498 Perspectives. TPAMI, 35(8):1798-1828, 2013. 499
- 500 Chris Burgess and Hyunjik Kim. 3d shapes dataset. https://github.com/deepmind/3dshapes-dataset/, 501 2018.
- Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -VAE. NeurIPS, 2017. 504
- 505 Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. Isolating Sources of Disentanglement in Variational Autoencoders. NeurIPS, 2018. 506
- 507 Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog Bits: Generating Discrete Data using 508 Diffusion Models with Self-Conditioning. ICLR, 2023. 509
- Li Deng. The MINST Database of Handwritten Digit Images for Machine Learning Research. IEEE 510 *Signal Processing Magazine*, 29(6):141–142, 2012. 511
- 512 Prafulla Dhariwal and Alexander Nichol. Diffusion Models Beat GANs on Image Synthesis. 513 NeurIPS, 2021. 514
- Sixun Dong, Huazhang Hu, Dongze Lian, Weixin Luo, Yicheng Qian, and Shenghua Gao. Weakly 515 Supervised Video Representation Learning with Unaligned Text for Sequential Videos. CVPR, 516 2023. 517
- 518 Cian Eastwood and Christopher KI Williams. A Framework for the Quantitative Evaluation of 519 Disentangled Representations. *ICLR*, 2018.
- Babak Esmaeili, Robin Walters, Heiko Zimmermann, and Jan-Willem van de Meent. Topological Obstructions and How to Avoid Them. NeurIPS, 2023. 522
- 523 B Everett. An Introduction to Latent Variable Models. Springer Science & Business Media, 2013.
 - Jean-Yves Franceschi, Aymeric Dieuleveut, and Martin Jaggi. Unsupervised Scalable Representation Learning for Multivariate Time Series. NeurIPS, 2019.
- 527 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, 528 Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. NeurIPS, 2014.
- Alex Graves, Rupesh Kumar Srivastava, Timothy Atkinson, and Faustino Gomez. Bayesian Flow 530 Networks. arXiv preprint arXiv:2308.07037, 2023. 531
- 532 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016.
- 534 Eric Hedlin, Gopal Sharma, Shweta Mahajan, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi, 535 and Kwang Moo Yi. Unsupervised Semantic Correspondence Using Stable Diffusion. NeurIPS, 536 2024.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 538 Prompt-to-Prompt Image Editing with Cross Attention Control. arXiv preprint arXiv:2208.01626, 2022.

540	Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M Botyinick,
541	Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning basic visual concepts with a
542	constrained variational framework. ICLR, 2017.
543	
544	Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. <i>NeurIPS</i> ,
545	2020.
546 547	Geonho Hwang, Jaewoong Choi, Hyunsoo Cho, and Myungjoo Kang. MAGANet: Achieving Com-
548	binatorial Generalization by Wodening a Group Action. Tewl., 2025.
549 550	Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou. Variational Deep Embedding: An Unsupervised and Generative Approach to Clustering. <i>arXiv preprint</i>
551 552	arXiv:1611.05148, 2016.
553	Hyunjik Kim and Andriy Mnih. Disentangling by Factorising. <i>ICML</i> , 2018.
554	Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. ICLR, 2014.
555 556 557	Rahul Krishnan, Dawen Liang, and Matthew Hoffman. On the challenges of learning with inference networks on sparse, high-dimensional data. <i>AISTATS</i> , 2018.
558 559	Alex Krizhevsky and Geoffrey Hinton. Learning Multiple Layers of Features from Tiny Images. Technical report, University of Toronto, 2009.
560 561	James Kwok and Ryan P Adams. Priors for Diversity in Generative Latent Variable Models. <i>NeurIPS</i> , 2012.
562	
563	Tianhong Li, Huiwen Chang, Shlok Kumar Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan.
564	MAGE: MAsked Generative Encoder to Unify Representation Learning and image synthesis.
505	CVPR, 2023.
567 568	Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep Learning Face Attributes in the Wild. <i>ICCV</i> , 2015.
569 570	Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The Variational Fair Autoencoder. <i>ICLR</i> , 2016.
571 572 573	Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. Diffusion Hyperfeatures: Searching through Time and Space for Semantic Correspondence. <i>NeurIPS</i> , 2024.
574 575	Cristian Meo, Louis Mahon, Anirudh Goyal, and Justin Dauwels. αTC-VAE: On the relationship between Disentanglement and Diversity. <i>ICLR</i> , 2024.
576 577 578 579	Soumik Mukhopadhyay, Matthew Gwilliam, Yosuke Yamaguchi, Vatsal Agarwal, Namitha Padman- abhan, Archana Swaminathan, Tianyi Zhou, and Abhinav Shrivastava. Do Text-free Diffusion Models Learn Discriminative Visual Representations? <i>arXiv preprint arXiv:2311.17921</i> , 2023.
580 581 582	Khalid Oublal, Said Ladjal, David Benhaiem, Emmanuel LE BORGNE, and François Roueff. Disentangling Time Series Representations via Contrastive Independence-of-Support on l-Variational Inference. <i>ICLR</i> , 2024.
583 584 585	Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Diffusion Autoencoders: Toward a Meaningful and Decodable Representation. <i>CVPR</i> , 2022.
586 587	Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating Diverse High-Fidelity Images with VQ-VAE-2. <i>NeurIPS</i> , 2019.
วชช 589 590	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- Resolution Image Synthesis with Latent Diffusion Models. <i>CVPR</i> , 2022.
591 592 593	Huajie Shao, Shuochao Yao, Dachun Sun, Aston Zhang, Shengzhong Liu, Dongxin Liu, Jun Wang, and Tarek Abdelzaher. ControlVAE: Controllable Variational Autoencoder. <i>ICML</i> , 2020.

Ken Shoemake. Animating rotation with quaternion curves. SIGGRAPH, 1985.

594 595 596	Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder Variational Autoencoders. <i>NeurIPS</i> , 2016.
597 598	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. <i>ICLR</i> , 2021.
599 600	Yuxuan Song, Jingjing Gong, Hao Zhou, Mingyue Zheng, Jingjing Liu, and Wei-Ying Ma. Unified Generative Modeling of 3D Molecules with Bayesian Flow Networks. <i>ICLR</i> , 2024.
602 603 604	Hiroshi Takahashi, Tomoharu Iwata, Atsutoshi Kumagai, Sekitoshi Kanai, Masanori Yamada, Yuuki Yamanaka, and Hisashi Kashima. Learning Optimal Priors for Task-Invariant Representations in Variational Autoencoders. <i>KDD</i> , 2022.
605 606 607	Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent Correspondence from Image Diffusion. <i>NeurIPS</i> , 2023.
608	Jakub Tomczak and Max Welling. VAE with a Vampprior. AISTATS, 2018.
609 610 611	Sana Tonekaboni, Chun-Liang Li, Sercan O Arik, Anna Goldenberg, and Tomas Pfister. Decoupling Local and Global Representations of Time Series. <i>AISTATS</i> , 2022.
612 613	Michael Tschannen, Olivier Bachem, and Mario Lucic. Recent Advances in Autoencoder-Based Representation Learning. <i>NeurIPS Workshop</i> , 2018.
614 615 616	Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation. <i>CVPR</i> , 2023.
617	Arash Vahdat and Jan Kautz. NVAE: A Deep Hierarchical Variational Autoencoder. NeurIPS, 2020.
618 619 620	Aaron Van Den Oord, Oriol Vinyals, et al. Neural Discrete Representation Learning. <i>NeurIPS</i> , 2017.
621 622 623	Yingheng Wang, Yair Schiff, Aaron Gokaslan, Weishen Pan, Fei Wang, Christopher De Sa, and Volodymyr Kuleshov. InfoDiffusion: Representation Learning Using Information Maximizing Diffusion Models. <i>ICML</i> , 2023.
624 625 626	Zhangkai Wu, Longbing Cao, and Lei Qi. eVAE: Evolutionary Variational Autoencoder. TNNLS, 2024.
627 628	Weilai Xiang, Hongyu Yang, Di Huang, and Yunhong Wang. Denoising Diffusion Autoencoders are Unified self-supervised Learners. <i>ICCV</i> , 2023.
629 630 631	Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MINST: a Novel Image Dataset for Bench- marking Machine Learning Algorithms. <i>arXiv preprint arXiv:1708.07747</i> , 2017.
632 633	Tim Z Xiao and Robert Bamler. Trading Information between Latents in Hierarchical Variational Autoencoders. <i>ICLR</i> , 2023.
635 636	Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open- Vocabulary Panoptic Segmentation with Text-to-Image Diffusion Models. <i>CVPR</i> , 2023.
637 638 639 640	Jie Xu, Yazhou Ren, Huayi Tang, Xiaorong Pu, Xiaofeng Zhu, Ming Zeng, and Lifang He. Multi-VAE: Learning Disentangled View-Common and View-Peculiar Visual Representations for Multi-View Clustering. <i>ICCV</i> , 2021.
641 642 643	Kaiwen Xue, Yuhao Zhou, Shen Nie, Xu Min, Xiaolu Zhang, Jun Zhou, and Chongxuan Li. Uni- fying Bayesian Flow Networks and Diffusion Models through Stochastic Differential Equations. <i>ICML</i> , 2024.
644 645 646	Tao Yang, Xuanchi Ren, Yuwang Wang, Wenjun Zeng, and Nanning Zheng. Towards Building A Group-based Unsupervised Representation Disentanglement Framework. <i>ICLR</i> , 2022.
647	Eric Yeats, Frank Liu, David Womble, and Hai Li. NashAE: Disentangling Representations through

Adversarial Covariance Minimization. ECCV, 2022.

648 649 650	Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A Tale of Two Features: Stable Diffusion Complements DINO for Zero-Shot Semantic Correspondence. <i>NeurIPS</i> , 2024.
652 653	Haojie Zhao, Dong Wang, and Huchuan Lu. Representation Learning for Visual Object Tracking by Masked Appearance Transfer. <i>CVPR</i> , 2023a.
654 655	Shengjia Zhao, Jiaming Song, and Stefano Ermon. InfoVAE: Balancing Learning and Inference in Variational Autoencoders. <i>AAAI</i> , 2019.
657 658	Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text- to-image diffusion models for visual perception. <i>ICCV</i> , 2023b.
659 660	Mingyuan Zhou, Tianqi Chen, Zhendong Wang, and Huangjie Zheng. Beta Diffusion. <i>NeurIPS</i> , 2023.
660	
662	
664	
665	
666	
667	
668	
669	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
691	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
606	
697	
698	
699	
700	
701	

1	Introduction
2	Understanding Bayesian Flow Networks - An Alternative View
3	ParamReL: Parameter Space Representation Learning
	3.1 The ParamReL Framework
	3.2 Parameter Encoding through A Self-encoder
	3.3 Conditional Decoder
	3.4 Sampling and Reverse-sampling Processes
	3.5 Training and Test with ParamReL
4	Related Work
5	Experiments
	5.1 Evaluation Setup
	5.2 Semantic Representation of Discrete Data by ParamReLd
	5.3 Semantic Representation of Continuous Data by ParamReLc
6	Conclusion and Limitations
A	Preliminaries
	A.1 Bayesian Flow Networks
	A.2 Bayesian Flow Distribution
	A.3 Generative Latent Variable Models for Representation Learning
	A.4 Illustration of Parameter Space Optimization
B	Proofs
	B.1 Derivation of ELBO for ParamReL
	B.2 Mutual Information Learning
С	Technical Details and Experimental Setup
	C.1 Encoder Architecture
	C.2 BFN Architecture
	C.3 Hyperparameters
D	Experiment details
	D.1 Interpolation
E	Additioanl results
	E.1 Sensitivity Analysis

A PRELIMINARIES

A.1 BAYESIAN FLOW NETWORKS

In Graves et al. (2023), BFNs assume two types of distributions: a simple *input distribution* $P_{\rm I}(\cdot)$ representing the initial belief about observations and an *output distribution* $P_{\rm O}(\cdot)$ simulating the observation distribution. The parameters of input distribution are first updated through a Bayesian inference scheme and then passed into a neural network $\psi(\cdot)$ to form the parameters of output distributions. The main objective of BFNs is to minimize the divergence between the ground-truth data distribution and the output distribution, ensuring that the output distribution closely approximates the ground-truth data distribution.

Following the notations in diffusion models, we denote x_0 as the observations. There are T reverse steps in BFNs which gradually reveals the information of x_0 through $\{x_T, x_{T-1}, \ldots, x_1\}$ to the input distribution². At each step t, \mathbf{x}_t is first noised through a sender distribution $p_{\rm S}(\hat{\mathbf{x}}_t | \mathbf{x}_t; \alpha_t)$, with α_t denoting the precision. Combined with input distribution $p_{I}(\mathbf{x}_t; \boldsymbol{\theta}_{t+1})$, the posterior dis-tribution of \mathbf{x}_t is obtained as $p(\mathbf{x}_t; h(\boldsymbol{\theta}_{t+1}, \hat{x}_t, \alpha_t)) \propto p_{\mathrm{I}}(\mathbf{x}_t; \boldsymbol{\theta}_{t+1}) p_{\mathrm{S}}(\hat{x}_t | \mathbf{x}_t; \alpha_t)$, where $\boldsymbol{\theta}_t = \mathbf{x}_t$ $h(\theta_{t+1}, \hat{x}_t, \alpha_t)$ is the Bayesian update function. By feeding this intermediate (posterior) param-eter θ_t into a neural network $\psi(\cdot)$, \mathbf{x}_t 's output distribution $p_O(\cdot)$ is parameterized as $p_O(\mathbf{x}_t; \psi(\boldsymbol{\theta}_t))$. Finally, a *receiver distribution* $p_{\rm R}(\cdot)$ is defined as the expectation of the sender distribution with respect to the output distribution, i.e., $p_{\mathrm{R}}(\widehat{x}_t; \psi(\boldsymbol{\theta}_t), \alpha_t) := \mathbb{E}_{p_{\mathrm{O}}(\mathbf{x}_t; \psi(\boldsymbol{\theta}_t))}[p_{\mathrm{S}}(\widehat{x}_t | \mathbf{x}_t; \alpha_t)]$. See Fig-ure 6 (a) for a visualization of the relationships between these distributions.

In BFNs, the joint distribution over the observation \mathbf{x}_0 and the intermediates $\{\mathbf{x}_t\}_t$ is defined as $p(\mathbf{x}_0, \{\mathbf{x}_t\}_t | -) := p_O(\mathbf{x}_0; \psi(\boldsymbol{\theta}_0)) \prod_{t=1}^T p_R(\hat{x}_t; \psi(\boldsymbol{\theta}_t), \alpha_t)$. This intractable joint distribution can be approximated under the variational inference framework as follows:

$$\log p(\mathbf{x}_{0}) \geq \mathbb{E}_{p_{\mathrm{F}}(\boldsymbol{\theta}_{1:T}|-)p_{\mathrm{S}}(\{\mathbf{x}_{t}\}_{t}|-)} \left[\log \frac{p_{\mathrm{O}}(\mathbf{x}_{0};\psi(\boldsymbol{\theta}_{0}))\prod_{t=1}^{T}p_{\mathrm{R}}(\hat{x}_{t};\psi(\boldsymbol{\theta}_{t}),\alpha_{t})}{\prod_{t=1}^{T}p_{\mathrm{S}}(\hat{x}_{t} \mid \mathbf{x}_{t};\alpha_{t})}\right] = -\sum_{t=1}^{T} \underbrace{\mathbb{E}_{p_{F}(\boldsymbol{\theta}_{t}|-)}\mathrm{KL}\left[p_{\mathrm{S}}\left(\hat{x}_{t} \mid \mathbf{x}_{0};\alpha_{T:t}\right) \parallel p_{\mathrm{R}}\left(\hat{x}_{t};\psi(\boldsymbol{\theta}_{t}),\alpha_{t}\right)\right]}_{\mathcal{L}_{t}^{\mathrm{R}}(\mathbf{x})} + \underbrace{\mathbb{E}_{p_{F}(\boldsymbol{\theta}_{0}|-)}\ln p_{\mathrm{O}}(\mathbf{x}_{0};\psi(\boldsymbol{\theta}_{0}))}_{\mathcal{L}^{\mathrm{D}}(\mathbf{x})},$$
(4)

where $p_{\rm F}(\theta_t|-)$ is the distribution of θ_t (see Appendix A.2 for a detailed calculation). Maximizing Eq. 4 equals minimizing the discrepancy $\mathcal{L}_t^{\rm R}(\mathbf{x})$ between the sender and receiver distributions and penalizing Distortion $\mathcal{L}^{\rm D}(\mathbf{x})$ to maximize the likelihood distribution over data.

Table 2: Examples of detailed distribution formats in BFNs. $\theta_{t+1} = \{\mu_{t+1}, \rho_{t+1}^{-1}\}$). cate: categorical distribution.

Data type	$p_{\mathrm{I}}(\mathbf{x}_t oldsymbol{ heta}_{t+1})$	$p_{\mathrm{S}}(\widehat{x}_t \mathbf{x}_t; lpha_t)$	$\boldsymbol{\theta}_t = h(\boldsymbol{\theta}_{t+1}, \widehat{x}_t, \alpha_t)$		
Continuous data	$\mathcal{N}(\mathbf{x}_t; \mu_{t+1}, \rho_{t+1}^{-1})$	$\mathcal{N}(\widehat{x}_t; \mathbf{x}, \alpha_t^{-1})$	$\mu_t = \frac{\alpha_t \widehat{x}_t + \rho_{t+1} \mu_{t+1}}{\alpha_t + \rho_{t+1}}$		
Discrete data	$\operatorname{Cat}(\mathbf{x}_t; \frac{1}{K} \cdot 1)$	$\mathcal{N}(\widehat{x}_t; \alpha_t K \mathbf{e}_{\mathbf{x}_t} - \alpha_t, \alpha_t K \mathbf{I})$	$\boldsymbol{\theta}_t = rac{e^{\widehat{x}_t} \boldsymbol{\theta}_{t+1}}{\sum_k e^{\mathbf{x}_{t-1,k}} \theta_{t+1,k}}$		
Data type	$p_{\mathrm{O}}(\mathbf{x}_t \boldsymbol{\theta}_t)$	$p_{ ext{R}}(\widehat{x}_t \psi(oldsymbol{ heta}_t),lpha_t)$			
Continuous data	$\delta(\mathbf{x}_t - \psi(\boldsymbol{\theta}_t))$	$\mathcal{N}(\widehat{x}_t;\psi(oldsymbol{ heta}_t),lpha_t^{-1})$			
Discrete data	$\operatorname{Cat}(\operatorname{softmax}(\psi(\boldsymbol{\theta}_t)))$	$\sum_{k} p_O(k; \psi(\boldsymbol{\theta}_t)) \mathcal{N}(\widehat{x}_t; \alpha_t K \mathbf{e}_k - \alpha_t, \alpha_t R)$			

²It is noted that the index t is used reversely in Graves et al. (2023). We make such changes to be consistent with the diffusion model settings Ho et al. (2020); Song et al. (2021).



Figure 7: The reverse-sampling process in ParamReL.

A.2 BAYESIAN FLOW DISTRIBUTION

833 834

835 836

837

838

839 840 841

842 843

844

Bayesian flow distribution $p_{\rm F}(\cdot | \mathbf{x}; t)$ is the marginal distribution over input parameters at time t, given prior distribution, accuracy schedule α and Bayesian update distribution $p_U(\cdot | \boldsymbol{\theta}, \mathbf{x}; \alpha)$, as follows:

$$p_{\rm F}(\boldsymbol{\theta} \mid \mathbf{x}; t) = p_U(\boldsymbol{\theta} \mid \boldsymbol{\theta}_0, \mathbf{x}; \beta(t)).$$
(5)

A.3 GENERATIVE LATENT VARIABLE MODELS FOR REPRESENTATION LEARNING

Latent Variable Models (LVMs) Everett (2013) which aim at learning the joint distribution $p(\mathbf{x}, \mathbf{z})$ over data \mathbf{x} and latent variables \mathbf{z} present efficient ways for uncovering hidden semantics. In LVMs, the joint distribution $p(\mathbf{x}, \mathbf{z})$ is usually decomposed as $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$, where $p(\mathbf{z})$ represents prior knowledge for inference Tschannen et al. (2018), thus facilitating learning the conditional distribution $p(\mathbf{x} | \mathbf{z})$. Among LVMs, Variational AutoEncoders (VAEs) Kingma & Welling (2014) and diffusion models Ho et al. (2020); Song et al. (2021) are two representative approaches Kwok & Adams (2012).

In VAEs, latent variables z is obtained through an *encoder network* $q_{\phi}(\mathbf{z} | \mathbf{x})$, whereas observations are reconstructed through a *decoder network* $p_{\theta}(\mathbf{x} | \mathbf{z})$, with ϕ and θ being the encoder and decoder parameters.

The dimensions of z are usually much smaller than those of x, denoted as $|z| \ll |x|$, such that redundant information is effectively removed and the most semantically meaningful factors are abstracted Louizos et al. (2016). VAEs are popular for downstream tasks like disentanglement Higgins et al. (2017); Yang et al. (2022); Hwang et al. (2023); Esmaeili et al. (2023), classification Takahashi et al. (2022); Tonekaboni et al. (2022), and clustering Jiang et al. (2016); Xu et al. (2021).

860 On the other hand, diffusion models Ho et al. (2020); Song et al. (2021) first use T diffusion steps 861 to transform observation \mathbf{x} into a white noise \mathbf{x}_T and then use T denoising steps to reconstruct the 862 observation. Diffusion models have obtained impressive performance in the fidelity and diversity 863 of generation tasks. However, they might be unable to obtain meaningful latent semantics since the 864 dimensions of \mathbf{x} and \mathbf{x}_T are the same as $|\mathbf{x}| = |\mathbf{x}_T|$. Preechakul et al. (2022); Wang et al. (2023) have attempted to integrate a decodable auxiliary variable z to enable diffusion models to obtain
 low-dimensional latent semantics. However, they have not overcome issues like the slow training
 speed inherent to the diffusion and reverse processes.

A.4 ILLUSTRATION OF PARAMETER SPACE OPTIMIZATION

Figure 8 illustrates the optimal data distribution learned in the parameter space. The plot presents stochastic parameter trajectories for the input distribution mean (indicated by white lines) overlaid on a Bayesian flow distribution logarithmic heatmap.



Figure 8: This figure illustrates optimization in the parameter space after t iterations.

B PROOFS

B.1 DERIVATION OF ELBO FOR PARAMREL

We derive the ELBO of ParamReL defined in Eq. (3).

$$\begin{split} \log p(\mathbf{x}_{0}) \\ &= \log \int_{\{\mathbf{z}_{t}\}_{t}} \int_{\{\mathbf{x}_{t}\}_{t}} p\left(\mathbf{x}_{0}, \{\mathbf{x}_{t}\}_{t}, \{\mathbf{z}_{t}\}_{t} \mid \boldsymbol{\theta}_{0}, \alpha\right) d\{\mathbf{z}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} \\ &= \log \int_{\{\boldsymbol{\theta}_{t}\}_{t}} \int_{\{\mathbf{z}_{t}\}_{t}} \int_{\{\mathbf{x}_{t}\}_{t}} p(\{\boldsymbol{\theta}_{t}\}_{t}| -) p_{O}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{0}, \mathbf{z}_{0})) \prod_{t=T}^{1} p(\mathbf{z}_{t}) \mathbb{E}_{p_{O}(\mathbf{x}_{t}; \psi(\boldsymbol{\theta}_{t}, \mathbf{z}_{t}))} [p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t})] \\ &\quad d\{\mathbf{z}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} d\{\boldsymbol{\theta}_{t}\}_{t} \\ &= \log \int_{\{\mathbf{z}_{t}\}_{t}} \int_{\{\mathbf{x}_{t}\}_{t}} \int_{\{\boldsymbol{\theta}_{t}\}_{t}} p(\{\boldsymbol{\theta}_{t}\}_{t}| -) \frac{p_{O}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{0}, \mathbf{z}_{0})) \prod_{t=T}^{1} p(\mathbf{z}_{t}) \mathbb{E}_{p_{O}(\mathbf{x}_{t}; \psi(\boldsymbol{\theta}_{t}, \mathbf{z}_{t}))} [p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t})] \\ &\quad \cdot \prod_{t=1}^{T} p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t}) q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}, t) d\{\mathbf{z}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} d\{\boldsymbol{\theta}_{t}\}_{t} \\ &\geq \mathbb{E}_{\prod_{t=1}^{T} p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t}) q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}, t) d\{\mathbf{z}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} d\{\boldsymbol{\theta}_{t}\}_{t} \\ &\geq \mathbb{E}_{\prod_{t=1}^{T} p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t}) q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}, t) p(\boldsymbol{\theta}_{t} \mid -) \left[\log \frac{p_{O}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{0}, \mathbf{z}_{0})) \prod_{t=T}^{1} p(\mathbf{z}_{t}) \mathbb{E}_{p_{O}(\mathbf{x}_{t}; \psi(\boldsymbol{\theta}_{t}, \mathbf{z}_{t})) [p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t})] \\ &\quad \cdot \prod_{t=1}^{T} p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t}) q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}, t) d\{\mathbf{z}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} d\{\mathbf{x}_{t}\}_{t} d\{\boldsymbol{\theta}_{t}\}_{t} \\ &\geq \mathbb{E}_{\prod_{t=1}^{T} p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t}) q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}, t) \left[\log \frac{p_{O}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{0}, \mathbf{z}_{0})) \prod_{t=T}^{1} p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}; \alpha_{t}) q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}, t)} \right] \\ &= \sum_{t=1}^{T} \mathbb{E}_{p_{F}(\boldsymbol{\theta}_{t} \mid -) \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{t})} \left\{ \mathbb{E}_{p_{S}(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}; \alpha_{T}; t)} \left[\log \frac{p_{O}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{t}, \mathbf{z}_{t}), \alpha_{t}} \right] \right] \\ &\quad -\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}) \left[\log \frac{q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}}{p(\mathbf{z}_{t})} \right] \right\} + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{0}, \boldsymbol{\theta}_{0}} [\ln p_{O}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{0}, \mathbf{z}_{0}))] \\ &\quad -\mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}) \left[\log \frac{q_{\boldsymbol{\phi}}(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t}}{p(\mathbf{z}_{t}$$

B.2 MUTUAL INFORMATION LEARNING

$$\mathcal{L}_{\text{ParamReL}^{+}} = -\sum_{t=1}^{T} \mathbb{E}_{p_{\text{F}}(\boldsymbol{\theta}_{t}|-)} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{t})} \left\{ \text{KL} \left[p_{\text{S}} \left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}; \alpha_{T:t} \right) \parallel p_{\text{R}} \left(\mathbf{x}_{t-1}; \psi(\boldsymbol{\theta}_{t}, \mathbf{z}_{t}), \alpha_{t} \right) \right] - \frac{1-\gamma}{T} \text{KL} \left[q_{\boldsymbol{\phi}} \left(\mathbf{z}_{t} \mid \boldsymbol{\theta}_{t} \right) \parallel p(\mathbf{z}) \right] - \frac{\gamma+\lambda-1}{T} \text{KL} \left[q_{\boldsymbol{\phi}} \left(\mathbf{z}_{t} \right) \parallel p(\mathbf{z}) \right] \right\} + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{0},\boldsymbol{\theta}_{0})} \left[\ln p_{\text{O}}(\mathbf{x}_{0}; \psi(\boldsymbol{\theta}_{0}, \mathbf{z}_{0})) \right].$$
(7)

Unlike the rest of the terms that can be optimized directly using reparameterization tricks, the TC term cannot be directly optimized due to intractable marginal distribution $q_{\phi}(\mathbf{z}_t)$. Here, we follow the guidance in Zhao et al. (2019) to replace the TC term with any strict divergence D, where $D(q_{\phi}(\mathbf{z})||p(\mathbf{z})) = 0$ iff $q_{\phi}(\mathbf{z}) = p(\mathbf{z})$. We implement the Maximum-Mean Discrepancy (MMD) Zhao et al. (2019) from the divergence family. MMD is a statistical measure that quantifies the difference between two probability distributions by comparing their mean embeddings in a high-dimensional feature space. By defining the kernel function $\kappa(\cdot, \cdot)$, D_{MMD} is denoted as:

$$D_{\text{MMD}}\left(q(\cdot) \| p(\cdot)\right) = \mathbb{E}_{p(\mathbf{z}), p(\mathbf{z}')}\left[\kappa\left(\mathbf{z}, \mathbf{z}'\right)\right] - 2\mathbb{E}_{q(\mathbf{z}), p(\mathbf{z}')}\left[\kappa\left(\mathbf{z}, \mathbf{z}'\right)\right] + \mathbb{E}_{q(\mathbf{z}), q(\mathbf{z}')}\left[\kappa\left(\mathbf{z}, \mathbf{z}'\right)\right].$$
(8)

C TECHNICAL DETAILS AND EXPERIMENTAL SETUP

C.1 ENCODER ARCHITECTURE

In our proposed encoder architecture, the self-encoder $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t)$ also conditions on step (t + 1)'s upsampling layers $\{\mathbf{u}_{t+1,l}\}_{l=1}^L$, where L is the number of layers in the U-Net architecture. For the l-th upsampling layer $\mathbf{u}_{t+1,l}$ at step t + 1, we upsample it to the size of \mathbf{x}_t , update by the Bayesian update function, and pass through a bottleneck layer $B_l(\cdot)$ (He et al., 2016) to the low-dimensional size. As a result, the self-encoder is defined as $q_{\phi}(\mathbf{z}_t | \boldsymbol{\theta}_t, t) = \mathcal{N}\left(\mathbf{z}_t; g_{\mu}(\boldsymbol{\theta}_t, \{\mathbf{u}_{t+1,l}\}_{l=1}^L, t), g_{\sigma}(\boldsymbol{\theta}_t, \{\mathbf{u}_{t+1,l}\}_{l=1}^L, t)^2\right)$, where $g_{\mu}(\cdot), g_{\sigma}(\cdot)$ use the same structure as:

$$g_{\mu}(\boldsymbol{\theta}_{t}, \{\mathbf{u}_{t+1,l}\}_{l=1}^{L}, t), g_{\sigma}(\boldsymbol{\theta}_{t}, \{\mathbf{u}_{t+1,l}\}_{l=1}^{L}, t) := \sum_{l=0}^{L} \omega_{l} \cdot B_{l}(h(\mathbf{x}_{t}, \mathbf{u}_{t+1,l})) + \omega_{L+1} \cdot B_{L+1}(\boldsymbol{\theta}_{t})$$

where ω_l is the mixing weight of the *l*-th layer.

C.2 BFN ARCHITECTURE

Similar to the diffusion-based representation learning model, we update the U-Net architecture based on Residual Blocks and Attention Modules. However, unlike previous approaches Ho et al. (2020); Song et al. (2021); Preechakul et al. (2022); Wang et al. (2023), we use shallower layers in the upper and down modules while incorporating an additional attention mechanism in the bottleneck module to achieve significant representations. Figure 9 illustrates the specific structural differences.

C.3 HYPERPARAMETERS

Table 3 presents the hyperparameter settings for training ParamReL. Different bin values are provided for various continuous datasets. All models are trained for 50 epochs. "Channel mult" denotes the channel shapes in each ResNet block within the U-Net architecture.

- D EXPERIMENT DETAILS
- 969 D.1 INTERPOLATION
- The latent space interpolation can be described as follows. Firstly, we noise source images to generate latent pairs by sender distribution, $\langle \mathbf{x}_1^1, \mathbf{x}_1^2 \rangle$, where $\mathbf{x}_1^1 \sim q(\cdot | \mathbf{x}_N^1)$ and $\mathbf{x}_1^2 \sim q(\cdot | \mathbf{x}_N^2)$.



Figure 9: U-Net comparisons of InfoDiffusion (a) and ParamReL (b). We apply the Attention module in the bottleneck layer, shallower than the InfoDiffusion's U-Net.

Table 3: Hyperparameters for training Bayesian Flow Networks, U-Net architecture, training protocols, and devices. The training configuration of ParamReL is based on Preechakul et al. (2022); Dhariwal & Nichol (2021).

	Hyperparameter	CelebA	Shapes3D	CIFAR-10	FFHQ			
	Encoder base channels	64	64	64	128			
Encoder	Encoder attention resolution	[16]	[16]	[16]	[16]			
	Encoder channel multipliers	[1,2,4,8,8]	[1,1,2,3,4,4]	[1,1,2,3,4,4]	[1,1,2,3,4,4			
	Latent code z dimension	512	512	512	512			
	Base channels	64	64	64	128			
	Channel multipliers	[1,2,4,8]	[1,1,2,3,4]	[1,1,2,3,4]	[1,1,2,3,4]			
Decoder	Attention resolution	[16]	[16]	[16]	[16]			
	Images trained	130M	130M	130M	130M			
	Batch size	128	128	128	128			
Decoder	Learning rate		1	e-4)				
	Optimizer		Adam (no v	weight decay)				
	EMA rate	0.9999						
	Training T	1000						
	Diffusion loss	MSE with noise prediction ϵ						
	Diffusion var.	Not important for DDIM						
Device	GPU	H100	H100	H100	H100			

Then, we implement two methods from Shoemake (1985) to generate four interpolated latent pairs $\bar{\mathbf{x}}_{1:4}$, i.e., linear interpolation, and spherical interpolation:

$$\bar{\mathbf{x}}_{i} = (1 - \lambda)\mathbf{x}_{0}^{1} + \lambda \mathbf{x}_{0}^{2},$$

$$\bar{\mathbf{x}}_{i} = \frac{\sin((1 - \alpha)\theta)}{\sin(\theta)}\mathbf{x}_{0}^{1} + \frac{\sin(\alpha\theta)}{\sin(\theta)}\mathbf{x}_{0}^{1},$$
(9)

1024 where λ is the scale coefficient, $\alpha \in [0,1]$ denotes the interpolation steps, and $\theta = 1025$ $\operatorname{arccos}\left(\frac{(\mathbf{x}_{0}^{1})^{\top}\mathbf{x}_{0}^{2}}{\|\mathbf{x}_{0}^{1}\|\|\mathbf{x}_{0}^{2}\|}\right)$ is the angle between \mathbf{x}_{0}^{1} and \mathbf{x}_{0}^{2} .

Table 4: Comparison of representation learning algorithms on continuous data by disentanglement performance (mean \pm std) and classification. The quantitative results for each algorithm averaged over five trials. (Modeling on data space \mathcal{D} , parameter space \mathcal{P} ; Prior distribution specify: Gaussian q, Categorical c, Delta d; \uparrow higher is better, \downarrow lower is better; [Top-1, Top-2, Top-3]).

Prior	Prior	Mathada	CelebA				3DShapes		CIFAR-10	
on	type	Wiethous	$\mathcal{TAD}\uparrow$	$ATTRS\uparrow$	$\mathcal{FID}\downarrow$	$AUROC\uparrow$	$\mathcal{DCI}\uparrow$	$AUROC\uparrow$	$\mathcal{FID}\downarrow$	$AUROC\uparrow$
	-	AE	0.042 ±0.004	1.0 ±0.0	90.4±1.8	0.759 ± 0.003	0.219 ±0.001	0.796±0.007	169.4±2.4	0.721±0.001
	g	VAE Kingma & Welling (2014)	0.000 ± 0.000	0.0 ±0.0	94.3±2.8	0.770 ±0.002	0.276 ±0.001	0.799±0.002	177.2±3.2	0.743±0.002
\mathcal{D}	g	β -VAE Burgess et al. (2017)	0.088 ±0.051	1.6 ±0.8	99.8±2.4	0.699 ±0.001	0.281 ±0.001	0.801±0.001	183.3±3.1	0.769±0.003
2	g	InfoVAE Zhao et al. (2019)	0.000 ± 0.000	0.0 ±0.0	77.8±1.6	0.757 ±0.003	0.134 ± 0.001	0.829±0.003	160.7±2.5	0.814±0.006
	g	DiffAE Preechakul et al. (2022)	0.155 ±0.010	2.0 ±0.0	22.7±2.1	0.799 ±0.002	0.196 ±0.001	0.899±0.001	32.1±1.1	0.859±0.002
	g	InfoDiffusion Wang et al. (2023)	0.299 ±0.006	3.0 ±0.0	23.8±1.6	0.848 ±0.001	0.342 ±0.002	0.882±0.001	32.4±1.8	0.886±0.004
	c	ParamReL	0.221+0.032	30+00	23 8+1 7	0.841 ±0.006	0 453 +0 002	0.871+0.007	33 6+2 3	0.857+0.005
	C	$(\gamma = 0.9, \lambda = 0.1)$	0.22110.052	5.0 ±0.0	23.011.7	0.041 ±0.000	0.455 ±0.002	0.07110.007	55.012.5	0.00720.000
	c	ParamReL	0.286 ± 0.001	4.0 ± 0.0	24.7±1.3	0.848 ± 0.002	0.477 ±0.002	0.892±0.006	33.2±0.6	0.871±0.002
-		$(\gamma = 0.9, \lambda = 0.01)$								
P	с	ParamkeL	0.256 ±0.008	3.0 ±0.0	22.5±1.2	0.839 ±0.003	0.417 ±0.002	0.891±0.001	31.9±1.1	0.868±0.003
		$(\gamma = 1, \lambda = 0.1)$								
	c	$(\alpha = 1, \lambda = 0.01)$	0.261 ±0.01	5.0 ±0.0	22.6±1.2	0.846 ±0.009	0.477 ±0.002	0.901±0.007	31.8±1.1	0.892±0.004
		P_{a}								
	d	$(\alpha = 0.0, \lambda = 0.1)$	0.299 ±0.005	3.0 ±0.0	24.1±1.1	0.844 ±0.012	0.482 ±0.001	0.891±0.002	34.7±0.9	0.882±0.005
		$(\gamma = 0.3, \lambda = 0.1)$								
	d	$(\gamma - 0.9, \lambda - 0.01)$	0.302 ±0.005	4.0 ±0.0	22.1±1.6	0.850 ± 0.116	0.567 ±0.005	0.902±0.001	31.2±1.1	0.901±0.001
		ParamReL								
	d	$(\gamma = 1, \lambda = 0, 1)$	0.287 ±0.005	3.0 ±0.0	23.6±1.7	0.821 ±0.006	0.441 ±0.008	0.887±0.002	32.8±2.1	0.877±0.002
		ParamReL								
	d	$(\gamma = 1, \lambda = 0, 01)$	0.368 ±0.005	3.0 ± 0.0	21.6±1.1	0.865±0.004	0.485 ± 0.009	0.931±0.001	31.1±1.1	0.911±0.002

Ε ADDITIOANL RESULTS

E.1 SENSITIVITY ANALYSIS

The coefficient in the Eq. 7 will regulate the information flow under the variational bottleneck guid-ance Burgess et al. (2017); Shao et al. (2020); Wu et al. (2024), resulting in the tradeoff between generation and representation learning.

Figure 10 depicts the generation and representation tradeoff in discrete datasets under the different coefficient sets ($\gamma = \{0.9, 1\}, \lambda = \{0.01, 0.1\}$). When disentanglement pressure is applied (), the AUROC increases.

Table 4 depicts the generation and representation tradeoff in continous datasets under the different coefficient sets ($\gamma = \{0.9, 1\}, \gamma = \{0.1, 0.01\}$). ParamReL consistently scores highest on average, with moderate variance.



Figure 10: Effect of γ , and λ by different representation learning metrics over ParamReLd and ParamReLc. Notations: [AUROC, FID]; [(•, bMNIST), (**□**, bFashionMNIST)]; [(-, ParamReLd), $(-\cdot -, ParamReLc)$].

E.2 LOW RESOLUTION REPRESENTATION LEARNING

We illustrate the representation learning ability in CelebA for high-level representation learning in Figure 13 (a), time-varying representation learning in Figure 13 (b), latent interpolation in Figure 14 and disentanglement in Figure 15.

1080 E.3 UNCONDITIONAL GENERATION

Figure 12 illustrates the unconditional generation quality on bMNIST. Images sampled from VAEbased model are blurry, as shown in Figure 12 (b). We implement two sampling strategies in the Diffusion-based model Wang et al. (2023), and both can only sample grey-scale images. Figure 12 (c) is sampled from the DDIM sampler, and Figure 12 (d) is sampled from a two-phased sampling procedure: form timesteps T to T/2, denoise and sample using a pre-trained vanilla denoising diffusion model. For timesteps ranging from T/2 to 0, proceed with sampling utilizing the InfoDiffusion model. Figure 12 (e) is images generated from our ParamReLc model. We can conclude that ParamReL can be sampled from the discrete distribution where the image value is binarized.



(a) : Generated samples of ParamReL on binaryMNIST (b) : Generated samples of ParamReL on binaryFashionMNIST

Figure 11: Samples reconstructed from our trained ParamReL on dataset Binary-MNIST.







1178

1179 Figure 14: Comparisons of latent space interpolation among sample-based models and parameter-1180 based models on dataset CelebA. Only our ParamReL model (e) can learn a continuous, smooth latent space while ensuring near-exact image reconstruction. Specifically, while sample-based gen-1181 erative models can learn a continuous but unsmooth latent space, this leads to incomplete recon-1182 structions. For example, in (a-d), the attribute of eyeglasses is frequently omitted. Moreover, VAEs 1183 (a,b) tend to produce blurry images. Additionally, it is observable that sample-based models often 1184 compromise reconstruction in favour of representation learning, as evidenced by the failure of dif-1185 fusion model variants (c-d) to accurately reconstruct background characters in imageB. 1186



Figure 15: Disentanglement of ParamReL on CelebA. The interpretable traversal directions are displayed by traversing the encodings ranging from [-3, 3].



Figure 16: Generated samples trained ParamReL on CelebA (a), Shapes3D (b), CIFAR-10 (c).