IMPROVING DISTRIBUTION MATCHING VIA SCORE BASED PRIORS AND STRUCTURAL REGULARIZATION

Anonymous authors

Paper under double-blind review

Abstract

Distribution matching (DM) can be applied to multiple tasks including fair classification, domain adaptation and domain translation. However, traditional variational DM methods such as VAE-based methods unnecessarily bias the latent distributions towards simple priors or fail to preserve semantic structure leading to suboptimal latent representations. To address these limitations, we propose novel VAE-based DM approach which incorporates a flexible score-based prior and a semantic structure preserving regularization. For score-based priors, the key challenge is that computing the likelihood is expensive. Yet, our key insight is that computing the likelihood is unnecessary for updating the encoder and thus we prove that the necessary gradients can be computed using only one score function evaluation. Additionally, we adapted the structure preserving regularization inspired by the Gromov-Wasserstein distance, which explicitly encourages the retention of geometric structure in the latent space, even when the latent space has fewer dimensions than the observed space. Our framework further allows the integration of semantically meaningful structure from pretrained or foundation models into the latent space, ensuring that the representations preserve semantic structure that is informative and relevant to downstream tasks. We empirically demonstrate that our DM approach leads to better latent representations compared to similar methods for fair classification, domain adaptation, and domain translation tasks.

028 029 030 031

032

004

005

006

008

009 010 011

012

013

014

015

016

017

018

019

021

022

025

026

027

1 INTRODUCTION

As machine learning (ML) continues to advance, trustworthy ML systems not only require 033 impressive performance but also properties such as fairness, robustness, causality, and 034 explainability. Unfortunately, collecting more data or building bigger models, as scaling laws (Kaplan et al., 2020) propose, improve performance with larger models and datasets but don't necessary address to solve these problems. For example, historical bias or imbalanced 037 data can cause even well-trained models to produce unfair outcomes, requiring additional constraints to mitigate such biases. Distribution matching (DM), also known as distribution 038 alignment or domain-invariant representation learning, has emerged as a promising approach to address these challenges. By minimizing the divergence between latent representations, 040 distribution matching can introduce additional objectives to ML systems, enabling them to 041 learn representations that are fair, robust, and causal. This approach has been successfully 042 applied to a wide range of problems, including domain adaptation (Ganin et al., 2016; Zhao 043 et al., 2018), domain generalization (Muandet et al., 2013) causal discovery (Spirtes & Zhang, 044 2016), and fairness-aware learning (Zemel et al., 2013).

Despite the potential of distribution matching (DM) methods, they face significant challenges due to the vast number of possible mappings in the latent space. Without sufficient constraints, these methods often fail to maintain meaningful structural relationships in the learned representations from the data distribution, resulting in suboptimal latent representations for downstream tasks. A popular method for learning representations is the use of variational approaches like Variational Autoencoders (VAEs), which have been widely adopted for their stability during training and their ability to learn meaningful representations (Chen et al., 2019; Burgess et al., 2018).

653 However, VAEs typically rely on a simple prior—commonly an isotropic Gaussian distribution—over the latent space. While this assumption simplifies optimization and ensures

054 computational tractability in generative tasks, it biases the latent space, often leading to a 055 significant loss of structural information inherent in the data during transformation. This 056 loss disrupts the preservation of the data's geometric properties, which is particularly critical in unsupervised settings for learning meaningful and robust representations Chen et al. 057 (2020);Uscidda et al. (2024). Recent advancements in manifold learning have highlighted the importance of preserving intrinsic geometry of the data (Uscidda et al., 2024; Nakagawa et al., 2023; Hahm et al., 2024; Lee et al., 2022; Horan et al., 2021; Gropp et al., 2020; Chen 060 et al., 2020). Notably, Uscidda et al. (2024) and Nakagawa et al. (2023) demonstrate that 061 incorporating geometry-preserving constraints can induce disentanglement in the latent space. 062 They propose VAE frameworks that directly regularize the objective using Gromov-Monge 063 optimal transport, leveraging its ability to align latent representations with the data's inherent geometric structure. However, these approaches face significant practical challenges: the 064 simultaneous goals of preserving data geometry and matching a simple prior often result in 065 distortions within the latent space. To address this issue, Nakagawa et al. (2023) advocates 066 for the use of more expressive priors, such as meta-priors, Gaussian mixtures, and neural 067 priors, offering greater flexibility in capturing complex data distributions while preserving 068 geometric consistency. In contrast, Uscidda et al. (2024) retains the use of a simple prior 069 but focuses on learning latent representations that minimize feature distortion as effectively 070 as possible.

071 The prospect of utilizing powerful and flexible priors is particularly compelling, as they can relax the trade-off between prior matching and data geometry preservation, reducing 073 distortion and achieving better geometric consistency in the latent space. However, we 074 argue that approaches such as Gaussian mixture priors, meta-priors, or expressive neural priors (Vahdat et al., 2021; Makhzani et al., 2016; Tomczak & Welling, 2018) may suffer 075 from practical limitations, including poor scalability to high-dimensional spaces, significant 076 computational expense, or instability during training. To overcome these limitations, we 077 introduce the Score Function Substitution (SFS) trick, a novel approach that leverages a score 078 model to indirectly parameterize the prior distribution. By doing so, our method achieves a 079 balance between memory efficiency, stability during training, and geometric consistency in the latent space, providing a robust solution to the challenges faced by traditional distribution 081 matching frameworks.

- We summarize our contributions in the field of DM as follows:
 - Introduction of Score-Based Priors for Flexible Representation: We propose the Score Function Substitution (SFS) method to learn score-based priors, preserving complex data structures while enhancing the efficiency and stability compared to prior methods.
 - Semantic Structural Preserving Constraints Inspired by Gromov-Wasserstein Distance: To preserve geometry, we adopt the Gromov-Wassersteinbased constraint from Gromov Wasserstein Autoencoders (GWAE) Nakagawa et al. (2023). Specifically, we advocate for computing the cost function within the semantic space, if available, rather than the raw pixel space, as this approach is more suitable for capturing meaningful relationships in image datasets.
 - Empirical Validation: Our experiments demonstrate improved downstream task performance in fairness learning, domain adaptation, and domain translation using score-based priors and structural preservation.

⁰⁹⁶ 2 Preliminaries

Variational Alignment Upper Bound (VAUB) The paper by Gong et al. (2024)
presents a novel approach to distribution matching for learning invariant representations.
The author proposea a non-adversarial method based on Variational Autoencoders (VAEs),
called the VAE Alignment Upper Bound (VAUB). Specifically, they introduce alignment
upper bounds for distribution matching that generalize the Jensen-Shannon Divergence
(JSD) with VAE-like objectives. The author formalizea the distribution matching problem
with the following VAUB objective:

$$\operatorname{VAUB}(q(z|x,d)) = \min_{p(z)} \mathbb{E}_{q(x,z,d)} \left[-\log \frac{p(x|z,d)}{q(z|x,d)} p(z) \right] + C, \tag{1}$$

105 106

104

085

086

087

090

091

093

094

where q(z|x, d) is the probabilistic encoder, p(x|z, d) is the decoder, p(z) is the shared prior, and C is a constant independent of model parameters. The method ensures that the distribution matching loss is an upper bound of the Jensen-Shannon divergence (JSD), up to a constant. This non-adversarial approach overcomes the instability of adversarial training, offering a robust, stable alternative for distribution matching in fairness, domain adaptation, and robustness applications. Empirical results show that VAUB and its variants outperform traditional adversarial methods, particularly in cases where model invertibility and dimensionality reduction are required.

Score-based Models Score-based Models (Song et al., 2021c) are a class of diffusion 114 models that learn to generate data by denoising noisy samples through iterative refinement. 115 Rather than directly modeling the data distribution p(x), as done in many traditional 116 generative models, score-based models focus on learning the gradient of the log-probability 117 density of the target distribution, known as the score function. To learn the score function, Vincent (2011) and Song & Ermon (2019) propose training on the Denoising Score Matching 118 (DSM) objective. Essentially, data points x are perturbed with various levels of Gaussian 119 noise, resulting in noisy observations \tilde{x} . The score model is then trained to match the score 120 of the perturbed distribution. The DSM objective is defined as follows: 121

$$\mathrm{DSM} = \frac{1}{2L} \mathbb{E}_{q_{\sigma_i}(\tilde{x}|x)p_{\mathrm{data}}(x)} [\|s_{\phi}(\tilde{x},\sigma_i) - \nabla_{\tilde{x}} \log q_{\sigma_i}(\tilde{x}|x)\|_2^2],$$
(2)

where $q_{\sigma_i}(\tilde{x}|x)$ represents the perturbed data distribution of $p_{\text{data}}(x)$, and where L is the number of noise scales $\{\sigma_i\}_{i=1}^L$. When the optimal score network s_{ϕ}^* is found, $s_{\phi}^*(x) = \nabla_x \log q_{\sigma}(x)$ almost surely (Vincent (2011),Song & Ermon (2019)) and approximates $\nabla_x \log p_{\text{data}}(x)$ when the noise is small ($\sigma \approx 0$). Since score-based models learn the gradient of the distribution rather than the distribution itself, generating samples involves multiple iterative refinement steps. These steps typically leverage techniques such as Langevin dynamics, which iteratively updates the sample using the learned score function Song & Ermon (2019).

Gromov-Wasserstein Distance The Optimal Transport (OT) problem seeks the most efficient way to transform one probability distribution into another, minimizing transport cost. Given two probability distributions μ and ν over metric spaces (X, d_X) and (Z, d_z) , the OT problem is:

$$\inf_{\in \Pi(\mu,\nu)} \mathbb{E}_{(x,z)\sim\pi}[d(x,z)] \tag{3}$$

136 where $\Pi(\mu, \nu)$ is the set of couplings with marginals μ and ν , and d(x, z) is a cost function, 137 often the Euclidean distance. The Gromov-Wasserstein (GW) distance extends OT to 138 compare distributions on different metric spaces by preserving their relative structures, not 139 absolute distances. For distributions μ and ν over spaces (X, d_X) and (Z, d_z) , the GW 140 distance is:

$$GW(\mu,\nu) = \inf_{\pi \in \Pi(\mu,\nu)} \mathbb{E}_{(x,z) \sim \pi, (x',z') \sim \pi} [|d_X(x,x') - d_Z(z,z')|^2]$$
(4)

$$= \inf_{\pi \in \Pi(\mu,\nu)} \operatorname{GWCost}(\pi(x,z))$$
(5)

3 Methodology

=

122 123

135

141 142

143 144 145

146 147

148

149 150

151 152

161

3.1 TRAINING OBJECTIVE FOR DISTRIBUTION MATCHING WITH A SCORE-BASED PRIOR By employing VAUB(Gong et al., 2024) as our distribution matching(DM) objective \mathcal{L}_{DM} ,

$$\mathcal{L}_{\rm DM} = \mathcal{L}_{\rm VAUB} = \sum_{d} \frac{1}{\beta} \mathbb{E}_{q_{\theta}} \left[-\log \frac{p_{\varphi}(x|z,d)}{q_{\theta}(z|x,d)^{\beta}} Q_{\psi}(z)^{\beta} \right], \tag{6}$$

153 where d represents the domain $\forall d \in [1, \dots, D]$ (e.g., different class datasets or modalities), 154 and $\beta \in [0, 1]$ acts as a regularizer controlling the mutual information between the latent 155 variable z and the data x. $q_{\theta}(z|x, d)$ and $p_{\varphi}(x|z, d)$ are the d-th domain probabilistic encoder 156 and decoder, respectively, and $Q_{\psi}(z)$ is a prior distribution that is invariant to domains 157 (Gong et al., 2024). For notational simplicity, we ignore the SP loss and we assume $\beta = 1$. 158 We can split the VAUB objective into three components: reconstruction loss, entropy loss, 159 and cross entropy loss.

$$\mathcal{L}_{\text{VAUB}} \triangleq \sum_{d} \left\{ \underbrace{\mathbb{E}_{q_{\theta}}[-\log p_{\varphi}(x|z,d)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q_{\theta}}[-\log q_{\theta}(z|x,d)]}_{\text{entropy term}} + \underbrace{\mathbb{E}_{q_{\theta}}[-\log Q_{\psi}(z)]}_{\text{cross entropy term}} \right\}$$
(7)

The prior distribution in the cross-entropy term aligns with the encoder's posterior but is often restricted to simple forms like Gaussians or Gaussian mixtures, which can distort the encoder's transformation function Uscidda et al. (2024). To address this, we propose an expressive, learnable prior that adaptively mitigates such distortions, better capturing the underlying data structure.

167 Learning an arbitrary probabilistic density function (PDF) is often times intractable or 168 computationally expensive as the normalization constant must be computed. Therefore, 169 instead of modeling a neural network directly on the density Q(z), we propose to indirectly parameterize the prior via its score function $\nabla_z \log Q(z)$. But, the problem is that given only 170 the score function, it is difficult to compute the log likelihood of a sample. It is well-known 171 that weighted combinations of score matching losses do not generalize well and only provide 172 an approximation to maximum-likelihood estimation (MLE). Moreover, directly optimizing 173 MLE through the flow interpretation, while theoretically feasible, becomes computationally 174 expensive in practice as it requires solving an ODE at each optimization step Song et al. 175 (2021a). Modeling an arbitrary probabilistic density function (PDF) is computationally expensive due to the intractability of the normalization constant. Therefore, instead of 176 directly modeling the density Q(z), we propose to indirectly parameterize the prior via its 177 score function $\nabla_z \log Q(z)$. While this avoids direct density estimation, the score function 178 alone makes log-likelihood computations difficult. Weighted score matching losses only 179 approximate maximum-likelihood estimation (MLE), and directly optimizing MLE using 180 the flow interpretation becomes computationally prohibitive as it requires solving an ODE 181 at each step Song et al. (2021a). Unlike VAEs, where efficient sampling from the prior is 182 critical, we demonstrate that the distribution matching objective with a score-based prior can be optimized without costly sampling or computing log-likelihood. By reformulating 183 the cross-entropy term as a gradient with respect to the encoder parameters θ , we derive an 184 equivalent expression that retains the same gradient value. This allows us to decouple score 185 function training from the encoder and compute gradients with a single evaluation of the 186 score function. We call this the Score Function Substitution (SFS) trick. 187

Proposition 1 (Score Function Substitution (SFS) Trick). If $q_{\theta}(z|x)$ is the posterior distribution parameterized by θ , and $Q_{\psi}(z)$ is the prior distribution parameterized by ψ , then the gradient of the cross entropy term can be written as:

$$\nabla_{\theta} \mathbb{E}_{z_{\theta} \sim q_{\theta}(z|x)} \left[-\log Q_{\psi}(z_{\theta}) \right] = \nabla_{\theta} \mathbb{E}_{z_{\theta} \sim q_{\theta}(z|x)} \left[-\left(\underbrace{\nabla_{\bar{z}} \log Q_{\psi}(\bar{z}) \big|_{\bar{z}=z_{\theta}}}_{constant \ w.r.t. \ \theta} \right)^{\top} z_{\theta} \right], \tag{8}$$

where the notation of z_{θ} emphasizes its dependence on θ and $\cdot|_{\bar{z}=z_{\theta}}$ denotes that while \bar{z} is equal to z_{θ} , it is treated as a constant with respect to θ .

The full proof can be seen in Appendix A. In practice, Eqn. 8 detaches posterior samples from
the computational graph, enabling efficient gradient computation without additional backpropagation dependencies. Details are provided in the next section. Following Proposition 1,
we propose the score-based prior AUB (SAUB) objective defined as follows:

$$\mathcal{L}_{\text{SAUB}} \triangleq \sum_{d} \left\{ \mathbb{E}_{z \sim q_{\theta}(z|x,d)} \left[-\log p_{\varphi}(x|z,d) + \log q_{\theta}(z|x,d) - \left(\nabla_{\bar{z}} \log Q_{\psi}(\bar{z}) \big|_{\bar{z}=z} \right)^{\top} z \right] \right\}$$
(9)

Since our new loss does not affect terms related to φ , and by Proposition 1, we have $\nabla_{\theta,\varphi} \mathcal{L}_{VAUB} = \nabla_{\theta,\varphi} \mathcal{L}_{SAUB}$. However, $\nabla_{\psi} \mathcal{L}_{VAUB}$ and $\nabla_{\psi} \mathcal{L}_{SAUB}$ are not guaranteed to be equal and are likely different.

208 209 3.1.1 Deriving an Alternating Algorithm with Learnable Score-Based Priors

210 211

191 192 193

201 202

204

Optimizing the parameters θ , φ , ψ for the VAUB objective differs from the SAUB objective, as $\nabla_{\psi} \mathcal{L}_{\text{VAUB}} \neq \nabla_{\psi} \mathcal{L}_{\text{SAUB}}$, making direct optimization intractable. Furthermore, the SAUB objective is complicated by the lack of direct access to the score function. To address this, we train the prior parameters ψ separately from the encoder θ and decoder φ . Prior work Cho et al. (2022); Gong et al. (2024) shows that aligning the prior closely with the encoder's posterior improves the variational bound. Thus, we approximate the prior's score function using a score model $S_{\psi}(\cdot)$, trained on the denoising score matching objective with latent samples. This results in two training objectives:

$$\min_{\theta,\varphi} \sum_{d} \left\{ \mathbb{E}_{z \sim q_{\theta}(z|x,d)} \left[-\log p_{\varphi}(x|z,d) + \log q_{\theta}(z|x,d) - \left(S_{\psi}(z^*,\sigma_0 \approx 0) \Big|_{z^* = (z+\sigma_0\epsilon)} \right)^{\top} z \right] \right\}$$
(10)

$$\min_{\psi} \sum_{d} \left\{ \mathbb{E}_{q_{\sigma_i}(\tilde{z}|z)q_{\theta}(z|x,d)p_{\text{data}}(x,d)} \left[\left\| S_{\psi}(\tilde{z},\sigma_i) - \nabla_{\tilde{z}} \log q_{\sigma_i}(\tilde{z}|z) \right\|_2^2 \right] \right\}.$$
(11)

Eqn. 11 is the DSM objective, where $q_{\sigma_i}(\tilde{z}|z)$ is the perturbed latent representation, and $p_{\text{data}}(x,d)$ denotes the data distribution for domain d. Eqn. 10 is our SAUB loss with a fixed score model where $\epsilon \sim \mathcal{N}(0, I)$.

During VAE training, the score model is conditioned on the smallest noise level, $\sigma_0 = \sigma_{\min}$, to approximate the true score function. As previously mentioned, the output of the score model 230 is detached to prevent gradient flow, ensuring memory-efficient optimization by focusing solely 231 on the encoder and decoder parameters without tracking the score model's computational 232 graph. After optimizing the encoder and decoder, these networks are fixed while the score 233 model is updated using Eqn. 11. Theoretically, if the score model is sufficiently trained enough to fully capture latent distribution, it could be optimized using only small noise levels. However, extensive score model updates after each VAE step are computationally expensive. 235 To mitigate this, we reduce score model updates and train with a larger maximum noise 236 level, enhancing stability when the latent representation becomes out-of-distribution (OOD). 237 The complete training process is outlined in Appendix B. We also listed the stabilization 238 and optimization techniques in Appendix C.

3.2 Comparison with Latent Score-Based Generative Models and Connection to Score Distillation

241 242

243 Latent Score-Based Generative Models (LSGM) Vahdat et al. (2021) provide a robust 244 framework that combines latent variable models with score-based generative modeling, 245 leveraging diffusion processes to improve data generation quality. A key innovation in LSGM is the incorporation of a learnable neural network prior. Similar to our approach, LSGM 246 replaces the traditional cross-entropy term in the Evidence Lower Bound (ELBO) with 247 terms involving the score function, approximated using a diffusion model. To elucidate the 248 relationship between LSGM and our Score Function Substition (SFS) trick, we turn to 249 the concept of Score Distillation Sampling (SDS) loss. SDS loss was introduced to stabilize 250 the training of Implicit Neural Representation (INR) model parameters by circumventing 251 the computation of the Jacobian term of the diffusion model's U-Net during optimization. Computing this Jacobian term is analogous to approximating the Hessian of the data distribution, which has been empirically shown to be unstable, particularly at low noise levels. 253 Our approach appears to mirror the application of SDS loss within the LSGM framework. 254 Both methods utilize a score model to guide optimization toward higher-density regions while 255 avoiding the computation of the U-Net's Jacobian. Remarkably, this intuition is correct 256 (Appendix E for detailed derivation and explanation). By applying the Sticking-the-Landing principle Roeder et al. (2017) directly to LSGM, we derive that the SFS trick is proportional to a distilled LSGM loss. This technique allows us to update the encoder parameters 257 258 without backpropagating through the diffusion model, thereby avoiding potential instabilities 259 associated with approximating higher-order derivatives at low noise levels. The full proof of 260 this derivation is provided in subsection E.1. 261

262 263

3.2.1 Comparative Stability: SFS vs. LSGM

We evaluate stability by computing the negative log-likelihood (NLL) of the posterior against a predefined mixture Gaussian prior. Unlike standard training, which updates encoder, decoder, and prior parameters, our approach freezes the prior and uses a score model pretrained on the defined prior, updating only the encoder and decoder. The same pre-trained score model is used for both SAUB and LSGM to ensure a fair comparison. Performance is evaluated under four minimum noise levels, $\sigma_{\min} \in 0.001, 0.01, 0.1, 0.2$, with $\sigma_{\max} = 1$ fixed. While lower noise levels should improve likelihood estimation, as the score model more

290

291 292

293

295

296

297

298

299 300

301

302

305

310 311 312



Figure 1: (a) The prior distribution is the target distribution projected onto the Z-space. (b) The reconstruction loss and negative log-likelihood are presented on a logarithmic scale for improved visualization. The experiment uses consistent hyperparameters ($\beta = 0.1$), an identical VAE architecture, and the same pretrained score model.

precisely approximates the true score function, LSGM requires backpropagation through the score model's U-Net, which causes instability at low noise levels due to inaccurate gradients. As shown in Fig. 1, when $\sigma_{\min} = 0.001$, LSGM exhibits catastrophic instability, with diverging NLL and spikes in reconstruction loss. At $\sigma_{\min} = 0.1$ and $\sigma_{\min} = 0.2$, LSGM performs worse than at $\sigma_{\min} = 0.01$, indicating that unstable gradients at lower noise levels negatively impacts prior matching. This is concerning since low noise levels, like $\sigma_{\min} = 0.01$, are commonly used in practice. In contrast, the SFS trick shows greater stability across noise levels. At $\sigma_{\min} = 0.01$, the NLL is better than at $\sigma_{\min} = 0.1$, which outperforms $\sigma_{\min} = 0.2$, suggesting that SFS ensures more reliable gradients at lower noise levels. While both LSGM and SAUB degrade at $\sigma_{\min} = 0.001$, SFS stabilizes and achieves a better NLL than LSGM at $\sigma_{\min} = 0.01$, demonstrating its robustness in handling small noise configurations.

303 3.3 SEMANTIC PRESERVATION (SP) IN LATENT REPRESENTATIONS VIA GW INSPIRED 304 CONSTRAINT

The Gromov-Wasserstein (GW) distance Section 2 is a powerful tool for preserving structural relationships between distributions in different metric spaces. Nakagawa et al. (2023) introduces the GW metric \mathcal{L}_{GW} in an autoencoding framework, and we adopt this regularization in a similar manner.

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{DM}} + \lambda_{\text{GW}} \mathcal{L}_{\text{GW}}(q_{\theta}(z|x))$$
(12)

$$\mathcal{L}_{\rm GW}(q_{\theta}(z|x)) \triangleq {\rm GWCost}(\pi = q_{\rm data}(x)q_{\theta}(z|x)) = \mathbb{E}\left[\left|d_X(x,x') - d_Z(z,z')\right|^2\right]$$
(13)

where q_{data} represents the data distribution, d_X and d_Z are the predefined metric spaces for the observed and latent spaces, respectively, and λ_{GW} controls the importance of the structural preservation loss. $\mathcal{L}_{DM}(q_{\theta}(z|x))$ represents the distribution matching objective with $q_{\theta}(z|x)$ as the encoder, and $\mathcal{L}_{GW}(q_{\theta}(z|x))$ is the structural preservation loss where q_{data} is the data distribution, d_X and d_Z are the metric spaces for the observed and latent spaces, respectively, and λ_{GW} controls the GW loss $\mathcal{L}_{GW}(q_{\theta}(z|x))$. $\mathcal{L}_{DM}(q_{\theta}(z|x))$ is the distribution matching objective with encoder $q_{\theta}(z|x)$.

Selection of Metric Space and Distance Functions The GW framework's key strength
 lies in its ability to compare distributions across diverse metric spaces, where the choice of
 metric significantly impacts comparison quality. In low-dimensional datasets like Shape3D
 (Kim & Mnih, 2018) and dSprites (Matthey et al., 2017), Euclidean pixel-level distances
 align well with semantic differences, leading prior works (Nakagawa et al., 2023; Uscidda

324 et al., 2024) to use L2 or cosine distances for isometric mappings. However, this breaks down 325 in high-dimensional data, like real-world images, which lie on lower-dimensional manifolds. 326 The curse of dimensionality causes traditional metrics, such as pixel-wise distances, to lose effectiveness as dimensionality increases. Recent advancements in vision-language models 327 like CLIP (Radford et al., 2021) have shown their ability to learn robust and expressive 328 image representations by training on diverse data distributions Fang et al. (2022). Studies 329 Yun et al. (2023) demonstrate that CLIP captures meaningful semantic relationships, even 330 learning primitive concepts. Therefore, we propose using the semantic embedding space 331 of pre-trained CLIP models as a more effective metric for computing distances between 332 datasets, which we define as the Semantic Preservation (SP) loss. For a detailed evaluation 333 of the improvements from using CLIP embeddings, please refer to the Appendix F, which includes demonstrations and additional results. In the following section, we will denote the 334 Gromov-Wasserstein constraint as GW-EP, and GW-SP to differentiate the metric space we 335 used for Gromov-Wasserstein constraint as Euclidean metric space Preservation (EP) and 336 Semantic Structural Preservation (SP) respectively. 337

4 Related Works

339 **Learnable Priors** Most variational autoencoders (VAEs) typically use simple Gaussian 340 priors due to the computational challenges of optimizing more expressive priors and the 341 lack of closed-form solutions for their objectives. Early efforts to address this, such as 342 Adversarial Autoencoders (AAEs) Makhzani et al. (2016), employed adversarial networks to 343 learn flexible priors, resulting in smoother and more complete latent manifolds. Subsequent 344 research Hoffman & Johnson (2016); Johnson et al. (2017) highlighted that simple priors can 345 lead to over-regularized and less informative latent spaces, while Tomczak & Welling (2018) empirically showed that more expressive priors improve generative quality, with significant 346 gains in log-likelihood. More recently, Latent Score-based Generative Models (LSGM) Vahdat 347 et al. (2021) introduced score-based priors, leveraging a denoising score-matching objective to 348 learn arbitrary posterior distributions. This approach enables high-quality image generation 349 while capturing the majority of the data distribution.

350 **Gromov-Wasserstein Based Learning** Gromov-Wasserstein (GW) distance has found 351 numerous applications in learning problems involving geometric and structural configuration 352 of objects or distributions. Moreover, the GW metric has been adopted for mapping functions 353 in deep neural networks. One of the key benefits of GW distance is its capacity to compare distributions with heterogeneous data and/or dimensional discrepancies. Prior works, such 354 as Truong et al. (2022); Carrasco et al. (2024), although uses GW distance as part of the 355 loss in the the objective but is focusing on calculating and minimizing the GW objective 356 in the embedding space between domains $\mathcal{L}_{OT/GW} = OT/GW(z_{src}, z_{tqt})$. On the other 357 hand, Uscidda et al. (2024) defineds the GW objective as being calculated between the data 358 dimension and the embedding dimension. 359

- 5 Experiments
- 360 361 362

363

364

338

In this section, we evaluate the effectiveness of our proposed VAUB with a score-based prior on several tasks. We conduct experiments on synthetic data, domain adaptation, multi-domain matching, fairness evaluation, and domain translation. For each experiment, we compare our methods to VAUB and other baselines and evaluate performance using various metrics.

- 365 366
- 367 368

5.1 Improving Latent Space Separation by Using Score-based Prior

369 The primary objective of this experiment is to demonstrate the performance of different prior models within the VAUB framework. Additionally, we examine the effect of varying 370 the number of samples used during training, specifically considering scenarios with limited 371 dataset availability. To achieve this, we create a synthetic nested D-shaped dataset consists 372 of two domains and two labels, as illustrated in Fig. 2. The aim is to learn a shared 373 latent representation across two domains and evaluate the degree of separation between 374 class labels within this shared latent space. Since downstream tasks rely on these shared 375 latent representations, better separation of class labels in the latent space naturally leads to 376 improved classification performance. This setup draws an analogy to domain adaptation tasks, where the quality of separation in the latent representation relative to the label space 377 plays a critical role in determining downstream classification outcomes.

381 382

384

386 387

388

393

394

396

397

398



Figure 2: The dataset consists of two domains: Domain 1 (left nested 'D-shaped') and Domain 2 (right flipped 'D-shaped'). In each domain, the outer 'D' corresponds to Label 1, and the inner 'D' to Label 2. The shared latent spaces are visualized for models trained with varying data sizes (n = 20, 100, 500 samples) using Gaussian(Kingma et al., 2019), Mixture of Gaussians(Gong et al., 2024), Vampprior(Tomczak & Welling, 2018), LSGM,(Song et al., 2021c) and our score-based model (columns). Legends follow the format D{domain_index}_L{label_index}

399 In this experiment, we control the total 400 number of data samples generated for the 401 dataset, and compare the model's perfor-402 mance using five types of priors: Gaussian prior, Mixture of Gasussian Prior(MoG), 403 Vampprior, and a score-based prior trained 404 with LSGM, and ours (SFS method). Con-405 sidering the strong relations between point-406 wise distance and the label information of 407 the dataset, we use GW-EP to compute the 408 constraint loss in both in the data domain 409 and the latent domain. This helps to better visually reflect the underlying structure and 410 separations in the latent space. As shown 411 in Fig. 3, this performance improvement 412 is evident in the latent space: the nested 413 D structure is well-preserved under trans-414 formation with the two score-based prior 415 method (LSGM and ours), resulting in well-416 separated latent representations across different classes. This holds consistently true 417 for varying numbers of data points, from as 418 low as 20 samples to higher counts. On the



Figure 3: This figure shows label separation in the latent space under varying sample sizes and prior configurations, quantified by AU-ROC scores from the prediction of support vector classifier. Higher scores indicate better separation. Details of the metric are described in the appendix.

419 other hand, the Gaussian prior, MoG and Vamprior only achieves 90% of separation in the 420 latent space when the number of data samples is sufficiently large (n = 100 for MoG and 421 Vampprior prior and n = 20 for Gaussian prior), allowing the inner and outer classes to have a 422 classifier bound supported by enough data points as shown in Fig. 3. This finding is especially 423 relevant for real-world datasets, where the original data dimensionality can easily reach up to tens of thousands; while in this experiment, we worked with only a two-dimensional dataset, 424 yet the Gaussian, MoG and Vampprior required more than hundreds of samples to achieve 425 effective latent separation, whereas the score-based prior (LGSM and SFS) succeeded with 426 as few as 20 samples. 427

- 428
- 429 430

431

5.2 IMPROVING THE TRADEOFF BETWEEN ACCURACY AND PARITY ON FAIRNESS REPRESENTATION LEARNING

432 For this experiment, we apply our model to the well-433 known Adult dataset, derived from the 1994 census, 434 which contains 30K training samples and 15K test 435 samples. The target task is to predict whether an individual's income exceeds \$50K, with gender (a binary 436 attribute in this case) considered as the protected 437 attribute. We adopt the same preprocessing steps 438 in Zhao et al. (2020), and the encoder and classifier 439 architectures are consistent with those in Gupta et al. 440 (2021). We adapt GW-EP as our constraint loss con-441 sidering the lack of semantic models in tabular dataset such as Adult dataset. Please refer to Appendix H for 442 more detailed architecture setup. For comparison, we 443 benchmark our model against three non-adversarial 444 models FCRL(Gupta et al., 2021), CVIB(Moyer et al., 445 2018), VAUB(Gong et al., 2024) and one adversarial 446 model LAFTR-DP(Madras et al., 2018) and one ex-447 tra baseline 'Unfair Classifier' which is obtained to 448 serve as a baseline, computed by training the classifier directly on the original dataset. 449



Figure 4: Demographic Parity gap (Δ_{DP}) vs. Accuracy trade-off for UCI Adult dataset. Lower Δ_{DP} is better, and higher **Accuracy** is better.

450 As illustrated in Fig. 4, our method not only retains the advantages of the SAUB method, 451 achieving near-zero demographic parity (DP) gap while maintaining accuracy, but it also 452 improves accuracy across the board under the same DP gap comparing to other methods. We attribute this improvement largely to the introduction of the score-based prior, which 453 potentially allows for better semantic preservation in the latent space, enhancing both 454 accuracy and fairness. 455

- 456
- 5.3DOMAIN ADAPTATION
- 457 458

459 We evaluate our method on the MNIST-USPS domain adaptation task, transferring 460 knowledge from the labeled MNIST (70,000 461 images) to the unlabeled USPS (9,298 im-462 ages) without using target labels. We com-463 pare our SAUB method (with and with-464 out structure-preserving constraints) against 465 baseline DA methods: ADDA (Zhao et al., 2018), DANN (Ganin et al., 2016), and 466

Model	MNIST to USPS	(%) USPS to MNIST $(%)$
ADDA	89.4	90.1
DANN	77.1	73
VAUB	40.7	45.3
Ours w/o GW	88.1	85.54
Ours w/ GW-EP	91.4	92.7
Ours w/ GW-SP	96.1	97.4

Table 1: Domain adaptation accuracy (%) for MNIST to USPS and USPS to MNIST tasks.

VAUB (Gong et al., 2024). All methods use the same encoder and classifier architec-467 ture for fairness, with structure-preserving constraints applied using L2 distance in Euclidean 468 space(GW-EP) and CLIP embedding(GW-SP). 469

As shown in Table 1, our method outperforms the baselines in both directions. Unlike ADDA 470 and DANN, which require joint classifier and encoder training, our approach allows for 471 classifier training after the encoder is learned, simplifying domain adaptation. Additionally, 472 the inclusion of a decoder enables our model to naturally adapt to domain translation 473 tasks, as demonstrated in Fig. 14. We additionally conduct novel experiments to assess 474 the generalizability and robustness of our model with limited source-labeled data, detailed 475 in Appendix D. Additionally, image translation results between MNIST and USPS are 476 presented in Appendix J.

477 478

479

DOMAIN TRANSLATION 5.4

480 We conduct domain translation experiments on the CelebA dataset, translating images of 481 females with blonde hair to black hair and vice versa. We compare three settings: GW loss 482 in semantic space, GW loss in Euclidean space, and no GW loss. This comparison shows 483 that GW loss in the semantic space better preserves semantic features, while Euclidean 484 space GW loss is less effective in high-dimensional settings. We want to note that achieving state-of-the-art image translation performance is not the primary objective of our work; 485 instead, this experiment demonstrates our model's versatility across tasks.

100					
486	Task/Model	Top-1 (%)	Top-5 (%)	Top-10 (%)	Top-20 (%)
487 488	Black-to-Blonde Hair	5.0 ± 1.4	14.6 ± 9.4	24.4 ± 4.0	40.0 ± 2.5
489	GW-EP	3.0 ± 1.4 4.0 ± 1.0	14.0 ± 2.4 11.6 ± 2.2	24.4 ± 4.0 22.0 ± 2.9	40.0 ± 3.5 35.0 ± 2.6
490	GW-SP	9.0 ± 1.6	27.8 ± 3.1	39.2 ± 4.2	59.0 ± 2.9
491	Blonde-to-Black Hair	04115	100100	100 1 00	004000
492	NO GW CW FD	3.4 ± 1.7 2.0 ± 0.7	10.8 ± 3.3 0.2 \pm 1.8	19.0 ± 2.9 15.8 \pm 2.6	33.4 ± 3.9
493	GW-EF GW-SP	$egin{array}{c} 2.0 \pm 0.7 \ 4.8 \pm 2.3 \end{array}$	9.2 ± 1.0 18.8 ± 3.4	15.8 ± 2.0 28.6 ± 4.1	$egin{array}{c} 30.4 \pm 3.1 \ 46.2 \pm 2.5 \end{array}$
494					

Table 2: Top-k retrieval accuracy (%) for semantic preservation experiments. Bold values indicate the best performance for each metric.

For quantitative evaluation of semantic preservation, we utilize image retrieval accuracy as our metric. The models, trained for 1,500 epochs, translate images from a domain of 100 females with black hair to a domain of 100 females with blonde hair and vice versa. For each translated image, we compute the cosine similarity with all translated images in the target domain using CLIP embeddings To ensure fairness, we use a different pretrained CLIP model for evaluation and for training GW-SP for more information see Appendix H. This process is repeated five times with randomly selected datasets to account for variability in the data. The experiment aims to measure how well the translated images preserve their semantic content. We compute the top-k accuracy, where the task is to retrieve the correct translated image from the set of all translated images. This bidirectional evaluation black-to-blonde and blonde-to-black ensures robustness and highlights the model's ability to maintain semantic consistency during translation. The results show that applying GW-EP harms performance in high-dimensional datasets due to poor distance scaling. In contrast, GW-SP in semantic space consistently improves accuracy. Notably, GW-EP performs worse than no GW loss. The domain translation images in Appendix L confirm that models with semantic space GW loss better preserve semantic features like hairstyle, smile, and facial structure, demonstrating its advantage. For additional experiments, we provide image translations between male and female subjects on the FairFace dataset in Appendix K for interested readers.



(a) Black to Blonde Hair Female Translation



(b) Blonde to Black Hair Female Translation

Figure 5: All models use the same architecture. Refer to Appendix H for details on the neural network and CLIP model. Applying GW loss in the CLIP semantic space shows superior semantic preservation in both (a) and (b). The samples are selectively chosen to represent diverse variations; random samples are in Appendix L.

6 DISCUSSION AND CONCLUSION

In conclusion, we introduce score-based priors and structure-preserving constraints to address the limitations of traditional distribution matching methods. Our approach uses score models to capture complex data distributions while maintaining geometric consistency.
By applying Gromov-Wasserstein constraints in the semantic CLIP embedding space, we preserve meaningful relationships without the computational cost of expressive priors. Our experiments demonstrate improved performance in tasks like fairness learning, domain adaptation, and domain translation.

References 541

546

547

554

542	Christopher P. Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume
543	Desjardins, and Alexander Lerchner. Understanding disentangling in β -vae, 2018. URL
544	https://arxiv.org/abs/1804.03599.

- 545 Xavier Aramayo Carrasco, Maksim Nekrashevich, Petr Mokrov, Evgeny Burnaev, and Alexander Korotin. Uncovering challenges of solving the continuous gromov-wasserstein problem, 2024. URL https://arxiv.org/abs/2303.05978. 548
- Nutan Chen, Alexej Klushyn, Francesco Ferroni, Justin Bayer, and Patrick Van Der Smagt. 549 Learning flat latent manifolds with vaes. arXiv preprint arXiv:2002.04881, 2020. 550
- 551 Ricky T. Q. Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of 552 disentanglement in variational autoencoders, 2019. URL https://arxiv.org/abs/1802. 553 04942.
- Wonwoong Cho, Ziyu Gong, and David I. Inouye. Cooperative distribution alignment via 555 jsd upper bound. In Neural Information Processing Systems (NeurIPS), dec 2022. 556
- Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishaal Shankar, Achal Dave, 558 and Ludwig Schmidt. Data determines distributional robustness in contrastive language image pre-training (clip), 2022. URL https://arxiv.org/abs/2205.01397. 559
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François 561 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural 562 networks. The journal of machine learning research, 17(1):2096–2030, 2016. 563
- Ziyu Gong, Ben Usman, Han Zhao, and David I. Inouye. Towards practical non-adversarial distribution matching. In International Conference on Artificial Intelligence and Statistics 565 (AISTATS), May 2024.566
- 567 Amos Gropp, Matan Atzmon, and Yaron Lipman. Isometric autoencoders. arXiv preprint 568 arXiv:2006.09289, 2020.
- 569 Umang Gupta, Aaron M Ferber, Bistra Dilkina, and Greg Ver Steeg. Controllable guarantees 570 for fair outcomes via contrastive information estimation. In Proceedings of the AAAI 571 Conference on Artificial Intelligence, volume 35, pp. 7610–7619, 2021. 572
- Jaehoon Hahm, Junho Lee, Sunghyun Kim, and Joonseok Lee. Isometric representation 573 learning for disentangled latent space of diffusion models. arXiv preprint arXiv:2407.11451, 574 2024.575
- 576 Matthew D Hoffman and Matthew J Johnson. Elbo surgery: yet another way to carve up 577 the variational evidence lower bound. In Workshop in Advances in Approximate Bayesian Inference, NIPS, volume 1, 2016. 578
- 579 Daniella Horan, Eitan Richardson, and Yair Weiss. When is unsupervised disentanglement 580 possible? Advances in Neural Information Processing Systems, 34:5150–5161, 2021. 581
- 582 Matthew J. Johnson, David Duvenaud, Alexander B. Wiltschko, Sandeep R. Datta, and Ryan P. Adams. Composing graphical models with neural networks for structured repre-583 sentations and fast inference, 2017. 584
- 585 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon 586 Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. URL https://arxiv.org/abs/2001.08361.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas 589 Krause (eds.), Proceedings of the 35th International Conference on Machine Learning, 590 volume 80 of Proceedings of Machine Learning Research, pp. 2649–2658. PMLR, 10–15 Jul 2018. URL https://proceedings.mlr.press/v80/kim18b.html. 592
- Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. 593 Foundations and Trends® in Machine Learning, 12(4):307-392, 2019.

602

614

- Yonghyeon Lee, Sangwoong Yoon, MinJun Son, and Frank C Park. Regularized autoen coders for isometric representation learning. In International Conference on Learning Representations, 2022.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially
 fair and transferable representations. In *International Conference on Machine Learning*,
 pp. 3384–3393. PMLR, 2018.
 - Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders, 2016.
- Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Dis entanglement testing sprites dataset. https://github.com/deepmind/dsprites-dataset/,
 2017.
- Daniel Moyer, Shuyang Gao, Rob Brekelmans, Aram Galstyan, and Greg Ver Steeg. Invariant representations without adversarial training. Advances in Neural Information Processing Systems, 31, 2018.
- Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp. 10–18, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL https://proceedings.mlr.press/v28/muandet13.html.
- ⁶¹⁵ Nao Nakagawa, Ren Togo, Takahiro Ogawa, and Miki Haseyama. Gromov-wasserstein autoencoders, 2023. URL https://arxiv.org/abs/2209.07007.
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini
 Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL https://arxiv.org/abs/2103.00020.
- Geoffrey Roeder, Yuhuai Wu, and David Duvenaud. Sticking the landing: Simple, lower-variance gradient estimators for variational inference, 2017. URL https://arxiv.org/abs/1703.09194.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022.
 URL https://arxiv.org/abs/2010.02502.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. Advances in neural information processing systems, 32, 2019.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. Advances in neural information processing systems, 34: 1415–1428, 2021a.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models, 2021b. URL https://arxiv.org/abs/2101.09258.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
 Ben Poole. Score-based generative modeling through stochastic differential equations. In
 International Conference on Learning Representations, 2021c. URL https://openreview.
 net/forum?id=PxTIG12RRHS.
- Peter Spirtes and Kun Zhang. Causal discovery and inference: concepts and recent methodological advances. In *Applied informatics*, volume 3, pp. 1–28. Springer, 2016.
- Jakub Tomczak and Max Welling. Vae with a vampprior. In International conference on artificial intelligence and statistics, pp. 1214–1223. PMLR, 2018.
- Thanh-Dat Truong, Naga Venkata Sai Raviteja Chappa, Xuan Bac Nguyen, Ngan Le, Ashley
 Dowling, and Khoa Luu. Otadapt: Optimal transport-based approach for unsupervised
 domain adaptation, 2022. URL https://arxiv.org/abs/2205.10738.

Théo Uscidda, Luca Eyring, Karsten Roth, Fabian Theis, Zeynep Akata, and Marco Cuturi. Disentangled representation learning with the gromov-monge gap, 2024. URL https: //arxiv.org/abs/2407.07829. Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. Advances in neural information processing systems, 34:11287–11302, 2021. Pascal Vincent. A connection between score matching and denoising autoencoders. Neural Computation, 23(7):1661–1674, 2011. doi: 10.1162/NECO_a_00142. Tian Yun, Usha Bhalla, Ellie Pavlick, and Chen Sun. Do vision-language pretrained models learn composable primitive concepts?, 2023. URL https://arxiv.org/abs/2203.17271. Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In Sanjoy Dasgupta and David McAllester (eds.), Proceedings of the 30th International Conference on Machine Learning, volume 28 of Proceedings of Machine Learning Research, pp. 325–333, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL https://proceedings.mlr.press/v28/zemel13.html. Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J Gordon. Adversarial multiple source domain adaptation. Advances in neural information processing systems, 31, 2018. Han Zhao, Amanda Coston, Tameem Adel, and Geoffrey J. Gordon. Conditional learning of fair representations. In International Conference on Learning Representations, 2020. URL https://openreview.net/forum?id=HkeklONFPr.

702 703	С	ONTENTS	
704			
705	1	Introduction	1
706			
707	9	Droliminarios	9
708	4	1 Tellinna les	4
709	•		•
710	3	Methodology	3
711		3.1 Training Objective for Distribution Matching with a Score-based Prior	3
712		3.1.1 Deriving an Alternating Algorithm with Learnable Score-Based Priors	4
713			_
714 715		Score Distillation	5
716		3.2.1 Comparative Stability: SFS vs. LSGM	5
717		3.3 Semantic Preservation (SP) in Latent Representations via CW Inspired Con-	
718 719		straint	6
720			_
721	4	Related Works	7
722			
723	5	Experiments	7
724 725		5.1 Improving Latent Space Separation by Using Score-based Prior	7
726		5.2 Improving the Tradeoff between Accuracy and Parity on Fairness Representa- tion Learning	8
728		5.3 Domain Adaptation	9
729			0
730		5.4 Domain Translation	9
731 732	6	Discussion and Conclusion	10
733			
734	Ta	able of Contents	14
735			
736 737	Α	Proof of Proposition 1	16
738 739	В	Pseudo-code for learning VAUB with Score-Based Prior	16
740 741	С	Stabilization and Optimization Techniques	16
742			
743	D	Limited Source Label for Domain Adaptation	17
744 745		D.1 Results	17
746	-		
747	Ľ	More Detailed Discussion of Gradient Comparison Between LSGM and SFS Trick	18
748			10
749		E.1 Proof: SFS trick is proportional to a distilled LSGM loss	19
750			
751	F	Choices of different metric spaces in different dataset	21
752			
754	G	Multi-Domain Distribution Matching Setting	22
755	н	Detailed Architecture of the model	22

756	H.1 Fairness Representation Learning	. 22
757 758	H.2 Separation Metric for Synthetic Dataset	22
759	H.3 Domain Adaptation VAE Model	. 22
760	H.4 Domain Translation VAE Model	23
761	H.5 Domain Adaptation Classifier	. 23
763 764	H.6 Pretrained CLIP	23
765 I 766	More Synthetic Dataset Results	23
767 768 J	Image translation between MNIST and USPS	30
769 770 K	K FairFace Image Translation	30
771	K.1 Handpicked samples	. 30
772 773	K.2 Random Samples	. 30
774 775 L 776	Additional Random Image Translations on CelebA	31

810 A PROOF OF PROPOSITION 1

812 Proposition 1 (Score Function Substitution (SFS) Trick) 813

If $q_{\theta}(z|x)$ is the posterior distribution parameterized by θ , and $Q_{\psi}(z)$ is the prior distribution parameterized by ψ , then the *gradient* of the cross entropy term can be written as:

$$\nabla_{\theta} \mathbb{E}_{z_{\theta} \sim q_{\theta}(z|x)} \left[-\log Q_{\psi}(z_{\theta}) \right] = \nabla_{\theta} \mathbb{E}_{z_{\theta} \sim q_{\theta}(z|x)} \left[-z_{\theta}^{T} \underbrace{\nabla_{\bar{z}} \log Q_{\psi}(\bar{z})}_{\text{constant w.r.t. } \theta} \right], \tag{14}$$

where the notation of z_{θ} emphasizes its dependence on θ and $|_{\bar{z}=z_{\theta}}$ denotes that while \bar{z} is equal to z_{θ} , it is treated as a constant with respect to θ .

Proof.

816 817 818

819

820 821 822

823

824

825 826

827

828

829

830 831

832

833 834

835

836

837 838

839

840

841

842 843

844

845

846 847

848 849

850 851

852 853

854 855

856

857

858

859

 $\nabla_{\theta} \mathbb{E}_{z_{\theta} \sim q_{\theta}(z|x)} \left[-\log Q_{\psi}(z_{\theta}) \right]$ (15) $= \nabla_{\theta} \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[-\log Q_{\psi}(g_{\theta}(\epsilon)) \right]$ (Reparameterization trick: $z_{\theta} = g_{\theta}(\epsilon)$) (16) $= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[\nabla_{\theta} \left(-\log Q_{\psi}(g_{\theta}(\epsilon)) \right) \right]$ (17) $= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[\frac{\partial g_{\theta}(\epsilon)}{\partial \theta}^{\top} \frac{\partial \log Q_{\psi}(\bar{z})}{\partial \bar{z}} \Big|_{\bar{z}=g_{\theta}(\epsilon)} \right]$ (Chain rule: differentiating at $q_{\theta}(\epsilon)$) (18) $= \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[\nabla_{\theta} g_{\theta}(\epsilon)^{\top} \frac{\partial \log Q_{\psi}(\bar{z})}{\partial \bar{z}} \Big|_{\bar{z} = g_{\theta}(\epsilon)} \right]$ (Simplify notation) (19) $= \mathbb{E}_{\epsilon \sim p(\epsilon)} \nabla_{\theta} \left[\left(\frac{\partial \log Q_{\psi}(\bar{z})}{\partial \bar{z}} \Big|_{\bar{z} = g_{\theta}(\epsilon)} \right)^{\top} g_{\theta}(\epsilon) \right]$ (Move ∇_{θ} outside) (20) $= \nabla_{\theta} \mathbb{E}_{\epsilon \sim p(\epsilon)} \left[- \left(\nabla_{\bar{z}} \log Q_{\psi}(g_{\theta}(\epsilon)) \Big|_{\bar{z} = g_{\theta}(\epsilon)} \right)^{\top} g_{\theta}(\epsilon) \right]$ (Gradient applied to parts dependent on θ) (21) $= \nabla_{\theta} \mathbb{E}_{z_{\theta} \sim q_{\theta}(z|x)} \left[- \left(\nabla_{\bar{z}} \log Q_{\psi}(z_{\theta}) \Big|_{\bar{z}=z_{\theta}} \right)^{\top} z_{\theta} \right]$ (Change back to z_{θ} after pulling out gradient) (22)

B PSEUDO-CODE FOR LEARNING VAUB WITH SCORE-BASED PRIOR

See Alg. 1.

C STABILIZATION AND OPTIMIZATION TECHNIQUES

Several factors, such as interactions between the encoder, decoder, and score model, as well as the iterative nature of the optimization process, can introduce instability. To mitigate these issues, we implemented stabilization and optimization techniques to ensure smooth and robust training.

Batch Normalization on Encoder Output (Without Affine Learning) Applying
batch normalization to the encoder's mean output without affine transformations facilitates
smooth transitions in the latent space, acting as a soft distribution matching mechanism.
By centering the mean and mitigating large shifts, it prevents disjoint distributions, allowing
the score model to keep up with the encoder's updates. This regularization ensures the

-	Alg	orithm 1 Training VAUB with Score-based Prior (Alternating Optimization)
5	T	The Determinant of the second state of the sec
;	Inp	ut: Data x, domain a, parameters $\{\theta_d, \varphi_d, \psi\}$, hyperparameters: noise levels
,	$\{\sigma_{\mathrm{m}}\}$	$\{n, \sigma_{\max}\}$, number of loops L for score model update
}		
)	1:	Initialize: Parameters of Encoders θ , Decoders φ , and Score model ψ
	2:	while not converged do
	3:	Step 1: Update Encoder and Decoder parameters $\{\theta, \varphi\}$
	4:	Draw $x, d \sim p_{\text{data}}(x, d)$
	5:	Draw $z \sim q_{\theta}(z x,d)$
	6:	Calculate score by computing $S_{\psi}(z^*, \sigma = \sigma_{\min})$ using z^* , which is detached from the
		computational graph
	7:	Compute the following objective in Eqn. 10:
	8:	Perform gradient descent to minimize the objective and update $\{\theta, \varphi\}$
	9:	Step 2: Update Score Model parameters ψ
	10:	for $loop = 1$ to L do \triangleright Number of loops for score model update
	11:	Draw $x, d \sim p_{\text{data}}(x, d)$
	12:	Draw $z \sim q_{\theta}(z x,d)$
	13:	Draw perturbed latent variable $\tilde{z} \sim q_{\sigma_i}(\tilde{z} z)$, where $\sigma_i \in [\sigma_{\min}, \sigma_{\max}]$
	14:	Compute the DSM loss for the score model in Eqn. 11:
	15:	Perform gradient descent to minimize the DSM objective and update ψ
	16:	end for
	17:	Repeat alternating optimization steps until convergence.
	18:	end while

latent space remains within regions where the score model is trained, enhancing stabilityand reducing the risk of divergence.

Gaussian Score Function for Undefined Regions: To further stabilize training, we incorporate a small Gaussian score function into the score model to handle regions beyond the defined domain of the score function (i.e., outside the maximum noise level, σ_{max}). Inspired by the mixture neural score function in LSGMs Vahdat et al. (2021), this approach blends score functions to address out-of-distribution latent samples. The Gaussian score ensures smooth transitions and prevents instability in poorly defined areas of the latent space, maintaining robustness even in undertrained regions of the score model.

Weight Initialization and Hyperparameter Tuning: We observed that the initialization of weights significantly impacts the stability and convergence of our model. Poor
initialization can lead to bad alignment. Therefore, gridsearch was used to find an optimal
weight scale.

900 901

907

909

D LIMITED SOURCE LABEL FOR DOMAIN ADAPTATION

We introduce, to the best of our knowledge, a novel downstream task setup where there is limited labeled data in the source domain (i.e., 1%, 5%, 10%) and no supervision in the target domain. We apply this setup to the MNIST-to-USPS domain adaptation task. The objective is to determine how well our model with and without structural preservation can generalize with limited source supervision.

908 D.1 Results

As shown in Figure Fig. 6, our method without the SP constraint (which is entirely unsuper-910 vised in the source domain) demonstrates remarkable sample efficiency. With as little as 911 0.04% of the dataset (roughly two images per class), our method achieves an accuracy of 912 around 40%. By increasing the labeled data to just 0.1% (about five images per class), the 913 accuracy surpasses 73%. When we introduce the structural preservation constraint, which 914 allows the model to transfer knowledge from a pretrained model, we observe a significant 915 improvement in performance. With only 0.2% of the labeled data, the model's accuracy 916 approaches the performance of models trained on the full dataset. This boost in performance shows the effectiveness of incorporating semantic information into the latent space, allowing 917 the model to generalize better with minimal supervision.

The performance gap between models with and without the structural preservation (SP) constraint becomes more evident through UMAP visualizations of the latent space (Figure Fig. 6). While both methods achieve distribution matching and show label separation, the model without SP struggles to distinguish structurally similar digits, such as "4" and "9". In contrast, with the SP constraint, the latent space exhibits clearer, distinct separations, even for similar digits. The semantic structure injected by the SP constraint leads to more robust and meaningful representations, helping the model better differentiate between challenging classes. This highlights the effectiveness of the SP constraint in refining latent space organization.



Figure 6: (a) MNIST to USPS (LSDA). (b) USPS to MNIST (LSDA). (c) UMAP with SP. (d) UMAP without SP. All labeled data is randomly selected from the source dataset and tested on the target dataset, with results averaged over 10 trials. Both (a) and (b) demonstrate that with SP loss, the model is more robust to limited data. This is further supported by the corresponding UMAP visualizations, where (c) shows larger separation between classes compared to (d), reflecting better class distinction.

E More Detailed Discussion of Gradient Comparison Between LSGM and SFS Trick

Below, we detail the encoder and decoder optimization objectives for LSGM:

$$\min_{\theta,\varphi} \mathbb{E}_{q_{\theta}(z_{0}|x)} \left[-\log p_{\varphi}(x|z_{0}) \right] + \mathbb{E}_{q_{\theta}(z_{0}|x)} \left[\log q_{\theta}(z_{0}|x) \right] + \mathbb{E}_{t,\epsilon,q(z_{t}|z_{0}),q_{\theta}(z_{0}|x)} \left[\frac{w(t)}{2} \|\epsilon - \epsilon_{\psi}(z_{t},t)\|_{2}^{2} \right]$$

967 where w(t) is a weighting function, $\epsilon_{\psi}(\cdot)$ represents a diffusion model, and $\epsilon \sim \mathcal{N}(0, I)$. Similar 968 to our loss objective (refer to Eqn. 10), LSGM substitutes the traditional cross-entropy term 969 with a learnable neural network prior. Specifically, the final term in the Evidence Lower 970 Bound (ELBO) is replaced with a weighted denoising score matching objective.

971 We first adapt notations used in our objective for easy readability during comparison. Diffusion model can approximate the denoising score function by rewritting $\epsilon_{\psi}(z_t, t) =$

973 $\sigma_t S_{\psi}(z + \sigma_t \epsilon, \sigma_t)$ Song et al. (2022). To streamline discussion and avoid repetition, we 973 will refer to the final term of this formulation as the **LSGM objective**, we can write the 974 cross-entropy term of LSGM as below with the weighting function as $w(t) = g(t)^2/\sigma_t^2$ which 975 maximizes the likelihood between the encoder posterior and the prior where $g(\cdot)$ is the 976 diffusion coefficient typically proportional to the variance scheduling function Song et al. 977 (2021b)Vahdat et al. (2021).

$$\mathcal{L}_{\text{LSGM}} = \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z), q_{\theta}(z|x)} \left[\frac{w(t)}{2} \| \epsilon - \epsilon_{\psi}(z_t, t) \|_2^2 \right]$$
(23)

$$= \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z), q_{\theta}(z|x)} \left[\frac{g(t)^2}{2} \left\| \frac{\epsilon}{\sigma_t} - S_{\psi}(\tilde{z} = z + \sigma_t \epsilon, \sigma_t) \right\|_2^2 \right]$$
(24)

During encoder updates, the gradient computation for the last term with respect to the encoder parameters is expressed as:

$$\nabla_{\theta} L_{\text{LSGM}} = \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z), q_{\theta}(z|x)} \left[g(t)^2 \left(\frac{\epsilon}{\sigma_t} - S_{\psi}(\tilde{z}, \sigma_t) \right)^\top \frac{\partial S_{\psi}(\tilde{z}, \sigma_t)}{\partial \tilde{z}} \frac{\partial \tilde{z}}{\partial \theta} \right].$$
(25)

This framework requires computing the Jacobian term $\frac{\partial S_{\psi}(\bar{z},\sigma_t)}{\partial z_t}$, which is both computationally expensive and memory-intensive. To mitigate this, the Score Function Substitution (SFS) trick eliminates the need for Jacobian computation by detaching the latent input z^* in the score function from the encoder parameters. The resulting gradient is expressed as:

$$\nabla_{\theta} L_{\rm SFS} = -\mathbb{E}_{z \sim q_{\theta}(z|x,d)} \left[\frac{\partial z}{\partial \theta}^{\top} \left(S_{\psi}(z^*, \sigma \approx 0) \Big|_{z^*=z} \right) \right].$$
(26)

997 This modification provides significant advantages, reducing memory usage by bypassing the computational graph of the diffusion model's U-NET and enhancing stability. Poole et al. (2022) highlighted that the Jacobian computation approximates the Hessian of the dataset distribution, which is particularly unstable at low noise levels. Our empirical results in Fig. 1 confirm these findings, demonstrating improved stability with our loss objective compared to LSGM.

E.1 PROOF: SFS TRICK IS PROPORTIONAL TO A DISTILLED LSGM LOSS

To demonstrate that applying the Sticking-the-Landing principle Roeder et al. (2017) to LSGM yields the SFS trick, we begin by expressing Eqn. 25 in its score function form:

$$\nabla_{\theta} L_{\text{LSGM}} = g(t)^2 \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z), q_{\theta}(z|x)} \left[\left(\frac{\epsilon}{\sigma_t} - \nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z) \right)^\top \frac{\partial (\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z))}{\partial \tilde{z}} \frac{\partial \tilde{z}}{\partial \theta} \right].$$
(27)

1020 For clarity, we decompose Eqn. 27 into three components:

•
$$A = \left(\frac{\epsilon}{\sigma_t} - \nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z)\right)^{\mathsf{T}}$$

•
$$B = \frac{\partial(\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z))}{\partial \tilde{z}}$$

•
$$C = \frac{\partial \tilde{z}}{\partial \theta}$$

We first compute the expectation $\mathbb{E}[A^{\top}BC]$:

$$\mathbb{E}[A^{\top}BC] = \mathbb{E}\left[\frac{\epsilon}{\sigma_t}^{\top}BC\right] - \mathbb{E}\left[\nabla_{\tilde{z}}\log p_{\psi}(\tilde{z}|z)^{\top}BC\right]$$
(28)

$$= \mathbb{E}\left[\frac{\epsilon}{\sigma_t}^{\top}\right] \mathbb{E}[BC] - \mathbb{E}\left[\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z)^{\top} \frac{\partial(\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z))}{\partial \tilde{z}}C\right]$$
(29)

$$= -\mathbb{E}_{q_{\theta}(z|x)} \left[\mathbb{E}_{q_{\sigma_{t}}(\tilde{z}|z)} \left[\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z)^{\top} \frac{\partial (\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z))}{\partial \tilde{z}} \right] C \right].$$
(30)

1036 Here, Eqn. 28 follows from substituting the definition of A, Eqn. 29 separates the expectation 1037 terms because $\frac{\epsilon}{\sigma_t}^{\top}$ is independent of BC, and Eqn. 30 eliminates the first term since 1038 $\mathbb{E}\left[\frac{\epsilon}{\sigma_t}^{\top}\right] = 0.$

Next, we evaluate the expectation $\mathbb{E}[A^{\top}C]$:

$$\mathbb{E}[A^{\top}C] = \mathbb{E}\left[\left(\frac{\epsilon}{\sigma_t} - \nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z)\right)^{\top}C\right]$$
(31)

$$= -\mathbb{E}_{q_{\theta}(z|x)} \left[\mathbb{E}_{q_{\sigma_{t}}(\tilde{z}|z)} \left[\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z)^{\top} \right] C \right].$$
(32)

1047 In Eqn. 32, the term $\frac{\epsilon}{\sigma_t}$ is removed due to its independence from C and its zero expectation. 1048 Now, we assume the score model perfectly predicts the noisy latent representation of the 1049 encoder, i.e., $q_{\sigma_t}(\tilde{z}|z) = p_{\psi}(\tilde{z}|z)$, to compute $\mathbb{E}[A^\top BC]$ and $\mathbb{E}[A^\top C]$.

1051 For $\mathbb{E}[A^{\top}BC]$, considering only the inner expectation, we note that $p_{\psi}(\tilde{z}|z)$ is conditional Gaussian:

$$\mathbb{E}_{q_{\sigma_t}(\tilde{z}|z)} \left[\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z)^\top \frac{\partial (\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z))}{\partial \tilde{z}} \right]$$
(33)

$$= \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z)} \left[\left(\frac{\tilde{z} - z}{\sigma^2} \right)^\top \nabla_{\tilde{z}} \left(\frac{\tilde{z} - z}{\sigma^2} \right) \right]$$
(34)

$$= \mathbb{E}_{p_{\psi}(\tilde{z}|z)} \left[\left(\frac{\tilde{z} - z}{\sigma^2} \right)^{\prime} \left(\frac{1}{\sigma^2} \right) \right]$$
(35)

$$=0.$$
 (36)

1064 In Eqn. 35, the substitution $q_{\sigma_t}(\tilde{z}|z) = p_{\psi}(\tilde{z}|z)$ simplifies the expectation to zero. Similarly, 1065 for the inner expectation of $\mathbb{E}[A^{\top}C]$:

$$\mathbb{E}_{q_{\sigma_t}(\tilde{z}|z)}\left[\nabla_{\tilde{z}}\log p_{\psi}(\tilde{z}|z)\right] = \mathbb{E}_{p_{\psi}(\tilde{z}|z)}\left[\nabla_{\tilde{z}}\log p_{\psi}(\tilde{z}|z)\right]$$
(37)

$$= \int p_{\psi}(\tilde{z}|z) \frac{\nabla_{\tilde{z}} p_{\psi}(\tilde{z}|z)}{p_{\psi}(\tilde{z}|z)} d\tilde{z}$$
(38)

$$=\nabla_{\tilde{z}}\int p_{\psi}(\tilde{z}|z)d\tilde{z} \tag{39}$$

$$=0.$$
 (40)

1076 Thus, when the score model accurately predicts the posterior score function, removing the B1077 term from $\mathbb{E}[A^{\top}BC]$ introduces no gradient bias. Consequently, applying the Sticking-the-Landing methodology eliminates the Hessian term, reducing variance. The result is:

$$\nabla_{\theta} \mathcal{L}_{\text{Distilled}-\text{LSGM}} = g(t)^2 \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z), q_{\theta}(z|x)} \left[\left(\frac{\epsilon}{\sigma_t} - \nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z) \right)^\top \frac{\partial \tilde{z}}{\partial \theta} \right].$$
(41)

Sticking-the-Landing can be applied once more to Eqn. 41 by removing $\frac{\epsilon}{\sigma_t}$ and constraining on small noise levels ($\sigma_t \approx 0$), the gradient becomes proportional to that of the SFS trick:

$$\nabla_{\theta} L_{\text{SFS}} = -\mathbb{E}_{z \sim q_{\theta}(z|x,d)} \left[\frac{\partial z}{\partial \theta}^{\top} \left(S_{\psi}(z^*, \sigma \approx 0) \Big|_{z^*=z} \right) \right] \propto -g(t)^2 \mathbb{E}_{q_{\sigma_t}(\tilde{z}|z), q_{\theta}(z|x)} \left[\left(\nabla_{\tilde{z}} \log p_{\psi}(\tilde{z}|z) \right)^{\top} \frac{\partial \tilde{z}}{\partial \theta} \right].$$

$$(42)$$

\mathbf{F} Choices of different metric spaces in different dataset



Figure 7: Histogram of the pairwise distance between data samples within a class and between different classes for three datasets: MNIST, USPS, and CelebA. The amount of separation of two histogram is computed by using the AUROC score which being measured by a binary classifier to distinguish between with-in class results and between-class results. The class considered in MNIST and USPS is the digits, and in CelebA is hair color.

From the graph, we observe that for the MNIST and USPS datasets, both the Euclidean pixel space metric and the semantic space metric can effectively separate data pairs into within-class or between-class categories. However, the semantic space metric demonstrates a higher AUROC separation score, indicating that it provides a more reliable metric for distinguishing between these pair types.

In contrast, for the CelebA dataset, relying solely on pixel-based Euclidean distances struggles to differentiate whether the paired distances belong to within-class or between-class data

pairs. By employing a semantic metric, such as the one derived from CLIP, a clear distinction emerges, underscoring its utility.

These observations highlight that while pixel space metrics like Euclidean distance may be useful for certain datasets, semantic distance metrics, when available, often offer superior performance and may even be essential for datasets with more complex structures or features.

1141 G MULTI-DOMAIN DISTRIBUTION MATCHING SETTING

We train SAUB with SP on three different MNIST rotation angles: 0°, 30°, 60°. The top row is the ground truth image, the second row is the reconstruction, the third row is translation to MNIST 30°, and last row is translation to MNIST 60° in Fig. 8. Qualitatively most of the stylistic and semantic features are preserved with the correct rotation.

1147

1151

1140

1142

- 1148 H DETAILED ARCHITECTURE OF THE MODEL
- 1150
 - H.1 FAIRNESS REPRESENTATION LEARNING

The encoder is a 3-layer MLP with hidden dimension 64, and latent dimension 8 with ReLU layers connecting in between. The classifier is a 3-layer MLP with hidden dimension 64 with ReLU layers connecting in between.

11561157H.2Separation Metric for Synthetic Dataset

The classifier is trained by a support vector where hyperparmeters are chosen from the list
'C': [0.1, 1, 10, 100], 'gamma': [1, 0.1, 0.01, 0.001] with 5-fold cross validation. Error plot is generated from 5 runs.

1162

1164

- 1163 H.3 DOMAIN ADAPTATION VAE MODEL
- **1165** ENCODER ARCHITECTURE
- 1166

1169

1170

1171

1172

1173

1174

1175

1167 The encoder compresses the input image $\mathbf{x} \in \mathbb{R}^{1 \times 28 \times 28}$ into a latent representation. The 1168 architecture consists of the following layers:

- Conv2D: 4×4 , stride 2, 16 channels (input size $28 \times 28 \rightarrow 14 \times 14$).
 - **Residual Block:** 16 channels.
 - Conv2D: 4×4 , stride 2, 64 channels (input size $14 \times 14 \rightarrow 7 \times 7$).
- **Residual Block:** 64 channels.
- Conv2D: 3×3 , stride 2, $2 \times$ latent size channels (input size $7 \times 7 \rightarrow 4 \times 4$).
- **Residual Block:** $2 \times$ latent size channels.
 - Conv2D: 4×4 , stride 1, $2 \times$ latent size channels (output size $4 \times 4 \rightarrow 1 \times 1$).
 - Split into two branches for μ and $\log \sigma^2$, each with latent size channels.
- 1176 1177

1179

1186

1187

1178 DECODER ARCHITECTURE

1180 The decoder reconstructs the input image $\mathbf{x}' \in \mathbb{R}^{1 \times 28 \times 28}$ from the latent representation. The 1181 architecture consists of the following layers:

- 1182
 Reshape: Latent vector reshaped to size (latent size, 1, 1).
 Residual Block: latent size channels.
 - ConvTranspose2D: 4 × 4, stride 1, 64 channels (output size 1 × 1 → 4 × 4).
- 1184 ConvTranspose2D: 4 × 4, str
 1185 Residual Block: 64 channels.
 - ConvTranspose2D: 4×4 , stride 2, 16 channels (output size $4 \times 4 \rightarrow 8 \times 8$).
 - **Residual Block:** 16 channels.
 - **ConvTranspose2D:** 4×4 , stride 4, 1 channel (output size $8 \times 8 \rightarrow 28 \times 28$).

1188 1189	H.4 Domain Translation VAE Model					
1190	Encoder Architecture					
1191 1192 1193	The encoder compresses the input image $\mathbf{x} \in \mathbb{R}^{3 \times 64 \times 64}$ into a latent representation. The architecture consists of the following layers:					
1194 1195 1196 1197 1198 1199 1200 1201 1202	 Conv2D: 3 × 3, stride 2, 64 channels (input size 64 × 64 → 32 × 32). Residual Block: 64 channels. Conv2D: 3 × 3, stride 2, 128 channels (input size 32 × 32 → 16 × 16). Residual Block: 128 channels. Conv2D: 3 × 3, stride 2, 256 channels (input size 16 × 16 → 8 × 8). Residual Block: 256 channels. Conv2D: 3 × 3, stride 2, 2 × latent size channels (input size 8 × 8 → 4 × 4). Residual Block: 2 × latent size channels. Split into two branches for µ and log σ², each with latent size channels. 					
1203	Decoder Architecture					
1205 1206	The decoder reconstructs the input image $\mathbf{x}' \in \mathbb{R}^{3 \times 64 \times 64}$ from the latent representation. The architecture consists of the following layers:					
1207 1208 1209 1210 1211 1212 1213 1214 1215 1216	 Reshape: Latent vector reshaped to size (latent size, 4, 4). Residual Block: latent size channels. ConvTranspose2D: 3 × 3, stride 2, 256 channels (output size 4 × 4 → 8 × 8). Residual Block: 256 channels. ConvTranspose2D: 3 × 3, stride 2, 128 channels (output size 8 × 8 → 16 × 16). Residual Block: 128 channels. ConvTranspose2D: 3 × 3, stride 2, 64 channels (output size 16 × 16 → 32 × 32). Residual Block: 64 channels. ConvTranspose2D: 3 × 3, stride 2, 3 channels (output size 32 × 32 → 64 × 64). Sigmoid Activation: To map outputs to the range [0, 1]. 					
1217	H.5 Domain Adaptation Classifier					
1218 1219	Classifier consists of 2 linear layers and a ReLU activation function.					
1220	H.6 PRETRAINED CLIP					
1222 1223	For this work, we utilized pretrained CLIP models from the OpenCLIP repository. Specifically:					
1224	• ViT-H-14-378-quickgelu on dfn5b dataset was employed for training the GW-					
1226 1227	 ViT-L-14-quickgelu on dfn2b dataset was used for evaluation on the Image Retrieval task. 					
1228 1229 1230 1231 1232 1233 1234 1235	I More Synthetic Dataset Results					
1236 1237 1238 1239 1240 1241						



Figure 8: Multi-domain adaptation: MNIST images rotated at various angles.



Figure 9: This figures show the translated dataset, reconstructed dataset, as well as the latent space under sample size 20.



Figure 10: This figures show the translated dataset, reconstructed dataset, as well as the latent space under sample size 50.



Figure 11: This figures show the translated dataset, reconstructed dataset, as well as the latent space under sample size 100.



Figure 12: This figures show the translated dataset, reconstructed dataset, as well as the latent space under sample size 200.



Figure 13: This figures show the translated dataset, reconstructed dataset, as well as the latent space under sample size 500.

MNIST to USPS Sure O I J <thJ</th> J J J <t

$\rm J$ $\,$ Image translation between $\rm MNIST$ and $\rm USPS$

Figure 14: MNIST to USPS translated image trained with SP.

¹⁵⁸³ K FAIRFACE IMAGE TRANSLATION

This experimental setting is conducted in a fully unsupervised manner without SP loss. We compare our proposed score-based prior (SAUB) with a multi-Gaussian-based learning prior (VAUB) to evaluate their effectiveness.

K.1 HANDPICKED SAMPLES



(a) Male to Female translation

Female to Male translation

Figure 15: In this experiment, both models are trained in an unsupervised manner (i.e., SAUB is trained without GW-SP loss). SAUB clearly exhibits superior semantic preservation in both (a) and (b), particularly with respect to features such as skin color, race, and age. Notably, SAUB makes minimal adjustments when altering gender, while VAUB struggles to retain the identity of the original data. (These samples are handpicked to illustrate the trend.)

1618 K.2 RANDOM SAMPLES

In Fig. 16, we show completely random samples from the FairFace dataset.



Figure 16: Random samples from the FairFace experiment using our method. Top three rows translate from male to female and the bottom three rows translate from female to male. First row is original, second is reconstructed, and third is translated.

L Additional Random Image Translations on CelebA

Examples of random image translations between black hair and blonde hair are presented in Fig. 17 and Fig. 18

1674	GW-SP (semantic space)	QW-E (pixel space)	w/o GW
1675	Ground Truth Recon Translation Gro	sund Truth Recon Translation Ground	Futh Recon Translation
1676			
1677			
1678			
1679	122. 23. 63. 63.	25. 25. 10 12	E. 20. 165/
1680	EL EL EL		
1681			
1682			
1683			
1684			
1685			
1686			A DECK
1687			
1600			
1000	Read And And		The second second
1609			
1090			
1691			
1692			
1693			
1694			
1695			
1696			
1697	- 10 1 00 And 1		
1698	3 . 9 . 9.	3 . 3 . 3 . 3	
1699			
1700			
1701			
1702	nnal		
1703			
1704			
1705			
1706			
1707			
1708			
1709			
1710		🔊 🔊 📖 🖉	
1711			ALANS
1712			
1713			
1714			
1715			
1716			
1717			
1718			
1719			N (
1720			
1721			
1722	ELEPEN	EL EL EL TE	
1723			
1724			
1725			
1726			
1727			

Figure 17: Random Samples from Black to Blonde Hair Female

1728	GW-SP (semantic space)	GW-E (pixel space)	w/o GW
1729			
1730	60. 60. 60	(a). (a). (a)	a. (a), (a), (a)
1731			
1732			
1733			
1734			
1735			
1736			
1737			
1738			
1739			
1740			
1741			
1742			
1743			
1744		ion tor is	
1745			
1746	00 00 00	A	
1747			
1748		10-10-F	- m-
1749			
1750			
1751			and the second
1752	A ISA ISA		
1753	AAA	AAA	
1754	1-0 (-1 1-1)	1-17 1-17 1-1	A-PK A-PK A-P
1755			
1756			
1757			
1758	Nonia		
1759	60 60	60 60 6	
1760	Ve S N		
1761			
1762			
1763			
1764		60 00 0	
1765			
1766			
1767	20 60 00	20 00 00 0	A REAL AREA CE
1768			
1769			
1770			
1//1			
1772	6 6 6		6 6 6
1774			
1775			
1776			
1777			
1779			
1770	State 1 Here	E ILIN	
1780			
1781	1-A- (-A-	1-A-1 6	A A A A A
1/01			

Figure 18: Random Samples from Blonde to Black Hair Female