
Interpretability of Diffusion Models

Sirui, Li

223040027

School of Data Science

The Chinese University of Hong Kong, Shenzhen

223040027@link.cuhk.edu.cn

Yangkuan, Li

223040028

School of Data Science

The Chinese University of Hong Kong, Shenzhen

223040028@link.cuhk.edu.cn

Linnmeng, Li

223040026

School of Data Science

The Chinese University of Hong Kong, Shenzhen

223040026@link.cuhk.edu.cn

Abstract

1 Generative models are essential tools in machine learning, with Variational Au-
2 toencoders (VAEs) and diffusion models serving as primary instruments for syn-
3 thesizing complex data distributions. Despite their strengths, current models face
4 challenges in controllability and disentanglement, limiting their utility in refined
5 application scenarios. Addressing these challenges, this study introduces a novel
6 generative framework that amalgamates the high-quality generative capabilities of
7 diffusion models with the expressive and interpretable latent space of VAEs. This
8 integration aims to leverage the structured nature of VAEs' latent space to direct the
9 generative prowess of diffusion models, thereby achieving refined control and aug-
10 mented sample diversity. To empirically validate the model's efficacy, we conducted
11 extensive experiments focusing on image generation tasks, employing the CelebA-
12 HQ Mask dataset to benchmark against standalone VAE and diffusion model base-
13 lines. Our findings indicate that the DiffusionVAE not only excels in generating im-
14 ages with superior fidelity, as evidenced by lower Frechet Inception Distance (FID)
15 scores, but also demonstrates enhanced capability in controlling and disentangling
16 generative factors. This research contributes a novel perspective to the generative
17 model domain, offering theoretical and empirical insights into achieving controlled
18 and interpretable generation processes and opening avenues for future explorations
19 in highly controllable generative applications. The code for replicating the method
20 is available at <https://github.com/Naukode/ML-course-Project>.

21 **1 Introduction**

22 With the rapid development of machine learning technology, generative models have become an
23 important branch in the field of deep learning, especially demonstrating significant capabilities in the
24 synthesis of images, audio, and text. Among them, diffusion models and VAEs serve as primary tools
25 for generating complex data distributions, each showcasing unique advantages. Diffusion models are

26 particularly noted for their ability to produce high-quality images, while VAEs are widely applied in
 27 various tasks due to their powerful latent space expressiveness and higher interpretability. Nonetheless,
 28 current generative models still face challenges in terms of controllability and disentanglement, which
 29 limit their utility in more refined application scenarios.
 30 Controllability refers to the ability of a model to manipulate specific features in the generation process
 31 in a predetermined manner, while disentanglement involves the model's ability to independently
 32 control and represent different generative factors of the data. High controllability and good disentan-
 33 glement are crucial for enhancing the practicality and interpretability of generative models, especially
 34 when tasks require precise adjustments in generated content.
 35 Based on reading a lot of papers, Pandey et al. (2022) gives us direction and ideas. To address these
 36 issues, this study introduces the DiffusionVAE model, a new model that attempts to integrate the
 37 advantages of VAEs and diffusion models, aiming to improve both the quality of generation and
 38 the model's controllability and disentanglement. By combining the latent space of VAEs with the
 39 generative process of diffusion models, DiffusionVAE aims to utilize the structured latent space of
 40 VAEs to control the output of diffusion models, thereby achieving finer control and greater sample
 41 diversity.
 42 Furthermore, this research delves into controlling the generative process through operations in
 43 the latent space. Particularly in the DiffusionVAE model, we experimentally verify the impact of
 44 disentangling and controlling operations on the latent variables on the quality and diversity of the
 45 generated samples. We find that precise control over latent variables allows for adjustments to
 46 specific image features during the generation process, broadening the possibilities for applications of
 47 generative models.
 48 Through these studies, we not only demonstrate the effectiveness of DiffusionVAE in image generation
 49 tasks but also provide new theoretical and empirical foundations for research into the controllability
 50 and disentanglement of generative models. These results foreshadow potential future directions
 51 for generative models, particularly in applications requiring highly controllable and interpretable
 52 generation processes.

53 2 Background

54 2.1 VAEs

55 Developed by Kingma and Welling (2013), Variational Autoencoders (VAEs) represent a significant
 56 advancement in generative models, employing a probabilistic framework to encode and reconstruct
 57 data efficiently. The core architecture of VAE comprises two primary components: an encoder
 58 that maps input data into a latent space distribution characterized typically by Gaussian parame-
 59 ters—means and variances—and a decoder that reconstructs the input data from sampled latent
 60 variables. VAEs aim to optimize the Evidence Lower Bound (ELBO), where the ELBO can be
 61 satisfied as

$$ELBO = E_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x)||p(z)) \quad (1)$$

62 where $q_\phi(z|x)$ is the data's marginal likelihood, effectively balancing data likelihood maximization
 63 with latent distribution regularization towards a standard Gaussian. This optimization not only helps
 64 in accurate data reconstruction but also ensures that the latent space maintains meaningful and
 65 generalizable properties, making VAEs particularly useful for tasks requiring robust and interpretable
 66 generative models. The specific structure can be shown as Figure 1 below.

67 2.2 Diffusion Models

68 Diffusion models (Croitoru et al. (2023)) are a class of generative models that transform data into
 69 pure noise through a forward process and then reverse this transformation to regenerate the data
 70 from noise, effectively simulating a Markov chain where Gaussian noise is incrementally added
 71 to corrupt the data progressively. These models are particularly noted for producing high-quality,
 72 detailed outputs in image and audio generation. The forward process progressively shifts the data
 73 towards a standard Gaussian distribution, while the reverse process, which involves training a neural
 74 network to estimate and correct the noise at each step, aims to reconstruct the original data from

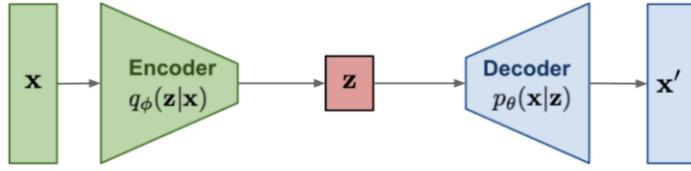


Figure 1: The Structure of VAEs

75 its noisy version, optimizing the model parameters to achieve accurate data recovery. The specific
76 structure can be shown as Figure 2 below.

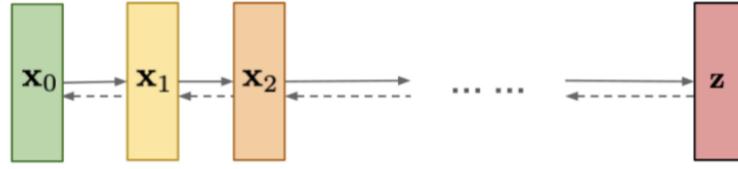


Figure 2: The Structure of Diffusion Models

77 3 Approach

78 3.1 Model Architecture and Training Objective

79 We present an innovative generative model architecture that combines the strengths of VAE and
80 DDPM. Our model's architecture operates in two stages: the first focuses on optimizing the VAE,
81 and the second optimizes the DDPM, integrating VAE's reconstructions into the diffusion process for
82 enhanced interpretability and generative performance. The detailed structure is shown in the Figure 3
83 below.

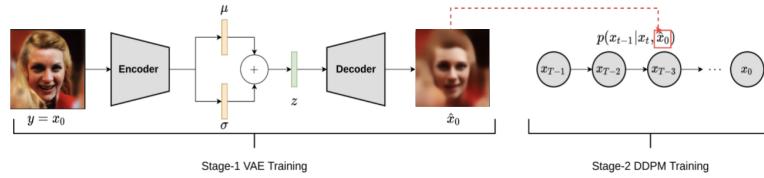


Figure 3: The Structure of our model

84 The joint distribution of our model is given by:

$$p(x_{0:T}, y, z) = p(z)p_\theta(y|z)p_\phi(x_{0:T}|y, z) \quad (2)$$

85 where x_0 is the original image, y is an auxiliary conditioning signal modeled using VAE, z represents
86 the latent space associated with y , and $x_{0:T}$ is a sequence of T images learned by the diffusion model.
87 Parameters θ and ϕ are associated with the VAE and DDPM components, respectively.

88 Due to the intractability of the joint posterior $p(x_{1:T}, z|y, x_0)$, we approximate it with a surrogate
 89 posterior:

$$q(x_{1:T}, z|y, x_0) = q_\psi(z|y, x_0)q(x_{1:T}|y, z, x_0) \quad (3)$$

90 The log-likelihood of the generative process is approximated by the Evidence Lower Bound (ELBO):

$$\log p(x_0, y) = \log \int p(x_{0:T}, y, z) dx_{1:T} dz \quad (4)$$

$$\geq \mathbb{E}_{q(x_{1:T}, z|x_0, y)} \left[\log \frac{p(x_{0:T}, y, z)}{q(x_{1:T}, z|x_0, y)} \right] \quad (5)$$

$$\geq \underbrace{\mathbb{E}_{q_\psi(z|y, x_0)} [p_\theta(y|z)] - D_{KL}(q_\psi(z|y, x_0)||p(z))}_{L_{VAE}} + \mathbb{E}_{z \sim q(z|y, x_0)} \underbrace{\left[\mathbb{E}_{q(x_{1:T}|y, z, x_0)} \left[\frac{p_\phi(x_{0:T}|y, z)}{q(x_{1:T}|y, z, x_0)} \right] \right]}_{L_{DDPM}} \quad (6)$$

93 To facilitate training, we assume:

- 94 • The conditioning signal y to be x_0 itself, ensuring a deterministic mapping.
- 95 • The second stage DDPM model is conditioned on the VAE reconstruction \hat{x}_0 , a deterministic
 96 function of z .
- 97 • A two-stage training process: first optimizing L_{VAE} and then L_{DDPM} , with θ and ψ fixed
 98 during the second stage.

99 3.2 Parameterization

100 3.2.1 VAE Parameterization

101 Standard VAE is employed as discussed in the background section.

102 3.2.2 DDPM Parameterization

103 **Forward Generative Process** We integrate VAE reconstruction \hat{x}_0 into the forward process as a
 104 condition, modeled as:

$$q(x_t|x_{t-1}, \hat{x}_0) = N(\sqrt{1-\beta_t}x_{t-1} + (1-\sqrt{1-\beta_t})\hat{x}_0, \beta_t \mathbf{I}) \quad (7)$$

$$q(x_t|x_0, \hat{x}_0) = N(\sqrt{1-\bar{\alpha}_t}x_0 + (1-\sqrt{1-\bar{\alpha}_t})\hat{x}_0, (1-\bar{\alpha}_t)\mathbf{I}) \quad (8)$$

106
 107 **Reverse Generative Process** We assume the transitions only depend on the VAE reconstruction \hat{x}_0 :

$$p(x_{0:T}|z) \approx p(x_{0:T}|\hat{x}_0) \quad (9)$$

108 4 Experiment

109 4.1 Generator-Refiner Framework

110 Our proposed "generator-refiner" framework employs a two-stage training approach that significantly
 111 enhances the visual quality of generated images. In the first stage, a VAE constructs a basic visual
 112 outline. The second stage involves a DDPM, which refines this initial image to achieve superior detail
 113 and realism. This framework is particularly advantageous for high-quality applications such as facial
 114 generation and artistic creation. We trained our models using the CelebA-HQ Mask dataset (Karras
 115 et al. (2017)), and the results (illustrated in Figure 4 and detailed in Table 1) demonstrated that our
 116 model significantly outperforms the VAE baseline in terms of image fidelity, as quantified by Frechet
 117 Inception Distance (FID) scores (Heusel et al. (2017)), maintaining competitive benchmarks against
 118 standalone DDPM models.

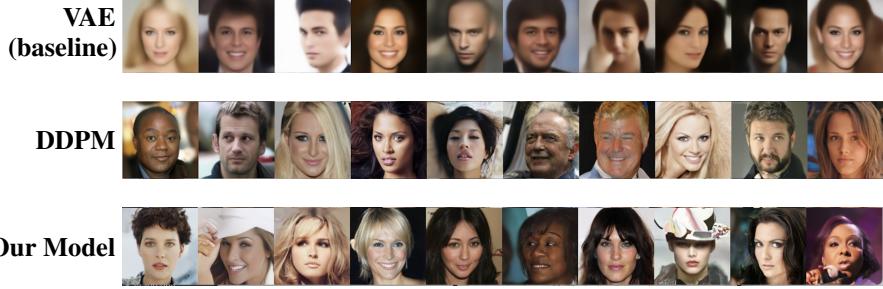


Figure 4: Sampling Variations Across Different Models

Models	FID on 5k samples
VAE (baseline)	115.82
DDPM	43.87
Our Model	51.44

Table 1: Comparison of FID Scores Across Different Models

119 4.2 Enhanced Controllable Generation through Latent Interpolation

120 4.2.1 Model Interpolation Techniques

121 Interpolation (Vahdat et al. (2022)) within our model is implemented through two distinct latent rep-
122 resentations: the low-dimensional VAE latent code z_{vae} and the DDPM intermediate representations
123 z_T .

- 124 • **Interpolation in VAE latent space z_{vae} :**

125 By sampling two VAE latent codes $z_{vae}^{(1)}$ and $z_{vae}^{(2)}$ from a standard Gaussian distribution and
126 interpolating between them, we generate an intermediate VAE latent code:

$$z_{vae} = \lambda z_{vae}^{(1)} + (1 - \lambda) z_{vae}^{(2)} \quad (10)$$

127 where λ is the interpolation factor ranging from 0 to 1. This method allows for the gen-
128 eration of Our model samples that, while maintaining overall image structure such as
129 facial expressions and hairstyles, exhibit enhanced details compared to their VAE-generated
130 counterparts.



Figure 5: Our model Samples Generated By Linearly Interpolating in the z_{vae} Latent Space

131 Figure 5 displays samples generated by our VAE through interpolation between two sampled
132 VAE codes, as outlined earlier. Correspondingly, Figure 4 (Top Row) presents samples
133 refined by our model, obtained by interpolating within the z_{vae} space. These refined samples,
134 derived from the initially blurry VAE outputs, notably retain the essential structure of the
135 image, demonstrating the efficacy of our generator-refiner framework in preserving key
136 image attributes while enhancing detail and clarity.

- 137 • **Interpolation in DDPM latent space X_T :**

138 With a fixed z_{vae} , we interpolate between two sampled DDPM representations $X_T^{(1)}$ and
139 $X_T^{(2)}$:

$$X_T = \lambda X_T^{(1)} + (1 - \lambda) X_T^{(2)} \quad (11)$$

140 This approach influences subtle features within the generated images, ensuring major
141 structure across samples.



Figure 6: Our model Samples Generated By Linearly Interpolating In The X_T Latent Space

142 Figure 6 illustrates the samples generated by our model with a fixed z_{vae} and interpolated x_T .
 143 It is evident that interpolating within the DDPM latent space results in variations in subtle
 144 features of the generated images. However, the overall structural integrity of the images
 145 remains consistent across different samples.

- **Consistency through Shared Stochasticity:**

147 By sharing the stochasticity in the DDPM reverse process across all generated samples,
 148 we achieve consistent latent interpolations and enable deterministic sampling. This shared
 149 approach implies the same style being applied to all refined samples, leading to smoother
 150 transitions and uniformity in style across different interpolations, as depicted in Figure 7.



Figure 7: Samples Generated With Sharing The Stochasticity

151 Figure 7 demonstrates that sharing stochasticity in the DDPM sampling process across
 152 samples results in a uniform stylization for all refined images. This uniformity promotes
 153 smooth transitions between interpolations.

154 4.3 Improved β -VAE for Disentanglement

155 To enhance latent space disentanglement, we employ an advanced β -VAE (Burgess et al. (2018))
 156 with a beta parameter greater than one, emphasizing KL divergence over reconstruction loss. This
 157 focus allows for significant disentanglement but could potentially degrade image quality. To mitigate
 158 this, we introduce a Dynamically Adjusted Parameter (C), initially prioritizing KL divergence to
 159 achieve disentanglement and subsequently adjusting to enhance reconstruction quality:

$$L(\theta, \phi; x, z, C) = E_{q_\phi(z|x)}[\log p_\theta(x|z)] - \gamma |D_{KL} q_\phi(z|x) || p(z) - C | \quad (12)$$

160 This dynamic adjustment ensures robust disentanglement while simultaneously maintaining high-
 161 quality image generation capabilities.

162 4.4 Experiment Details

163 To thoroughly evaluate the performance and robustness of our model framework, extensive exper-
 164 iments were conducted using carefully selected hyperparameters for both the β -VAE and DDPM
 165 stages. The following subsections detail the configuration settings and hyperparameters used, which
 166 have been optimized for the CelebA-HQ-Mask dataset to ensure that the results are both reliable and
 167 replicable.

168 4.4.1 Hyperparameters for improved β -VAE

169 The β -VAE stage of our model was rigorously trained with the following hyperparameters on Table 2:

- **Dataset:** High-resolution CelebA-HQ-Mask was chosen for its rich diversity in facial attributes, which is critical for assessing the performance of generative models in handling complex image details.
- **Resolution:** Each image was resized to 128x128 pixels to balance between computational efficiency and sufficient detail retention.
- **Batch Size:** Set at 128 to maximize hardware utilization without exceeding memory constraints.

- 177 • **Epochs:** The model was trained for 500 epochs to ensure thorough learning without
 178 overfitting.
- 179 • **Optimizer:** Adam optimizer was utilized for its adaptive learning rate capabilities, facilitating
 180 faster convergence.
- 181 • **Learning Rate:** Set at 1×10^{-4} , this rate was found optimal for steady progression in
 182 learning without causing instability in training dynamics.
- 183 • **Latent Code Size:** A size of 1024 was selected to capture a comprehensive amount of
 184 information about the data while maintaining manageable computational demands.
- 185 • **α Value:** Set at 100 to appropriately balance the trade-off between reconstruction fidelity
 186 and KL divergence.
- 187 • **Cmax and C_stop_iter:** These parameters were set at 25 and 100,000 iterations, re-
 188 spectively, to dynamically adjust the influence of KL divergence, promoting effective
 189 disentanglement over time.

Parameter	Value
Dataset	Celeb-A-HQ-Mask
Resolution	128x128 pixels
Batch Size	128
Epochs	500
Optimizer	Adam
Learning Rate	1×10^{-4}
Latent Code Size	1024
α value	100
Cmax	25
C_stop_iter	100,000

Table 2: Hyperparameters for the β -VAE Experiment

190 4.4.2 Hyperparameters for DDPM

The DDPM stage of our model was trained with the following hyperparameters on Table 3:

Parameter	Value
Dataset	CelebA-HQ-Mask
Resolution	128x128 pixels, 128 channels
Architecture Details	Scales of attention blocks: 16, Attention heads: 8, Residual blocks per scale: 2
Channel Multipliers	(1, 2, 2, 3, 4)
Total Parameters	95.2 million
Dropout	0.1
Noise Schedule	Linear from 1×10^{-4} to 0.02
Training Time Steps	100
EMA Decay Rate	0.9999
Batch Size	64
Optimizer	Adam
Learning Rate	2×10^{-5}
LR Annealing Steps	5000
Diffusion Loss Type	Noise prediction

Table 3: Hyperparameters and Configuration of the DDPM Experiment

- 191
- 192 • **Dataset and Resolution:** Consistent with the β -VAE, images from the CelebA-HQ-Mask
 193 dataset at 128x128 pixels were used.
- 194 • **Architecture Details:** The model features scales of attention blocks at 16, with 8 attention
 195 heads and 2 residual blocks per scale, ensuring detailed attention to features across various
 196 scales.

- 197 • **Channel Multipliers:** Configured as (1, 2, 2, 3, 4), these multipliers escalate the complexity
 198 and capacity of the model progressively through the network.
- 199 • **Total Parameters:** The model architecture comprises 95.2 million parameters, striking a
 200 balance between complexity and performance.
- 201 • **Dropout and Noise Schedule:** A dropout rate of 0.1 combined with a linear noise sched-
 202 ule from 1×10^{-4} to 0.02 was used to promote model robustness and effective noise
 203 management.
- 204 • **Training Time Steps, EMA Decay Rate, and Batch Size:** These were set at 100 steps, a
 205 decay rate of 0.9999, and a batch size of 64, respectively, to optimize the training process
 206 for quality and efficiency.
- 207 • **Optimizer and Learning Rate Adjustments:** The Adam optimizer with a learning rate
 208 of 2×10^{-5} and 5000 annealing steps ensured precise control over learning progression,
 209 minimizing potential for overfitting while accommodating for detailed feature capture.

210 **4.4.3 Compute Resources**

211 The experiments were conducted using high-performance computing resources to manage the substan-
 212 tial computational demands of training sophisticated generative models like β -VAE and DDPM. The
 213 specific hardware and execution times are detailed below, ensuring transparency and reproducibility
 214 of the results.

- 215 • **Hardware Configuration:** The training was performed on a robust setup involving high-end
 216 graphical processing units to effectively handle the training loads and memory requirements
 217 of the models. The specific data is shown in Table 4 below.

Resource Type	Specifications
Type of Compute Workers	Nvidia RTX 4090 GPU
Memory	24 GB

Table 4: Specifications of Compute Resources

- 218 • **Execution Time:** The execution time for each phase of the model training was carefully
 219 monitored to provide insights into the computational efficiency and practical deployment
 220 considerations of our framework. The data needed is shown below in Table 5.

Training Phase	Execution Time
Training for VAE	2 days
Training for DDPM	4 days

Table 5: Execution Time for Model Training Phases

221 The training of the β -VAE component was completed in 2 days, which involved extensive computa-
 222 tional processing to optimize the latent space and ensure robust feature capture. Following this, the
 223 DDPM component required an additional 4 days of training, reflecting the intricate calculations and
 224 extensive sampling processes intrinsic to diffusion models. The longer training time for the DDPM
 225 phase highlights its complexity and the computational effort required to refine the initial outputs from
 226 the β -VAE stage to achieve high fidelity and visually appealing results.

227 **5 Conclusion**

228 In summary, this study presents a novel generative framework that synergistically combines Varia-
 229 tional Autoencoders (VAEs) and diffusion models to address the challenges of controllability and
 230 disentanglement in generative models. Through extensive experimentation on the CelebA-HQ Mask
 231 dataset, our findings validate the model's superior capability to generate high-fidelity images while
 232 providing enhanced control over the generative process and improved disentanglement of generative

233 factors. Our model signifies a substantial step forward in the realm of generative models, offering a
234 nuanced approach that leverages the structured latent space of VAEs to guide the generative prowess
235 of diffusion models for controlled and diversified sample generation.

236 **Limitations** Despite the promising advancements introduced by our model, certain limitations warrant
237 further investigation. The model's computational demand is significant, necessitating considerable
238 resources for training and optimization, which could constrain its application in resource-limited
239 environments. Additionally, while the model exhibits commendable performance in image generation
240 tasks, its adaptability and efficacy across diverse data types and more complex generation tasks
241 remain to be fully explored.

242 **Future Work** Looking ahead, there are multiple avenues for future research to build upon the founda-
243 tions laid by this study. A primary area of interest involves exploring the potential of the framework
244 of our model in the realm of Natural Language Processing analogies, akin to the queen-king =
245 woman-man paradigm, for controllable image generation. This exploration could pave the way for
246 sophisticated applications where nuanced adjustments to generated images are necessary, further
247 bridging the gap between human cognitive processes and machine-generated content. Moreover,
248 efforts to optimize the computational efficiency of the model could significantly broaden its applica-
249 bility, making it more accessible for diverse applications. Additionally, extending the application
250 of our model to other data types and generation tasks promises to unlock new capabilities and
251 understandings in the field of generative models.

252 In conclusion, our model framework introduces a compelling approach to enhancing the quality, con-
253 trollability, and interpretability of generative models. By addressing current limitations and pursuing
254 outlined future directions, this research can contribute to the advancement of highly controllable and
255 interpretable generative applications, fostering further innovation in the field of machine learning.

256 **NeurIPS Paper Checklist**

257 **1. Claims**

258 Question: Do the main claims made in the abstract and introduction accurately reflect the
259 paper's contributions and scope?

260 Answer: [Yes]

261 Justification: In the section of Abstract and Introduction, we accurately show the paper's
262 contributions and scope.

263 Guidelines:

- 264 • The answer NA means that the abstract and introduction do not include the claims
265 made in the paper.
- 266 • The abstract and/or introduction should clearly state the claims made, including the
267 contributions made in the paper and important assumptions and limitations. A No or
268 NA answer to this question will not be perceived well by the reviewers.
- 269 • The claims made should match theoretical and experimental results, and reflect how
270 much the results can be expected to generalize to other settings.
- 271 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
272 are not attained by the paper.

273 **2. Limitations**

274 Question: Does the paper discuss the limitations of the work performed by the authors?

275 Answer: [Yes]

276 Justification: The information can be found in the conclusion part, such as significant
277 computational cost.

278 Guidelines:

- 279 • The answer NA means that the paper has no limitation while the answer No means that
280 the paper has limitations, but those are not discussed in the paper.
- 281 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 282 • The paper should point out any strong assumptions and how robust the results are to
283 violations of these assumptions (e.g., independence assumptions, noiseless settings,
284 model well-specification, asymptotic approximations only holding locally). The authors
285 should reflect on how these assumptions might be violated in practice and what the
286 implications would be.
- 287 • The authors should reflect on the scope of the claims made, e.g., if the approach was
288 only tested on a few datasets or with a few runs. In general, empirical results often
289 depend on implicit assumptions, which should be articulated.
- 290 • The authors should reflect on the factors that influence the performance of the approach.
291 For example, a facial recognition algorithm may perform poorly when image resolution
292 is low or images are taken in low lighting. Or a speech-to-text system might not be
293 used reliably to provide closed captions for online lectures because it fails to handle
294 technical jargon.
- 295 • The authors should discuss the computational efficiency of the proposed algorithms
296 and how they scale with dataset size.
- 297 • If applicable, the authors should discuss possible limitations of their approach to
298 address problems of privacy and fairness.
- 299 • While the authors might fear that complete honesty about limitations might be used by
300 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
301 limitations that aren't acknowledged in the paper. The authors should use their best
302 judgment and recognize that individual actions in favor of transparency play an impor-
303 tant role in developing norms that preserve the integrity of the community. Reviewers
304 will be specifically instructed to not penalize honesty concerning limitations.

305 **3. Theory Assumptions and Proofs**

306 Question: For each theoretical result, does the paper provide the full set of assumptions and
307 a complete (and correct) proof?

308 Answer:[Yes]

309 Justification: In the Parameterization part, we assume the transitions only depend on the
310 VAE reconstruction.

311 Guidelines:

- 312 • The answer NA means that the paper does not include theoretical results.
- 313 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
314 referenced.
- 315 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 316 • The proofs can either appear in the main paper or the supplemental material, but if
317 they appear in the supplemental material, the authors are encouraged to provide a short
318 proof sketch to provide intuition.
- 319 • Inversely, any informal proof provided in the core of the paper should be complemented
320 by formal proofs provided in appendix or supplemental material.
- 321 • Theorems and Lemmas that the proof relies upon should be properly referenced.

322 4. Experimental Result Reproducibility

323 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
324 perimental results of the paper to the extent that it affects the main claims and/or conclusions
325 of the paper (regardless of whether the code and data are provided or not)?

326 Answer: [Yes]

327 Justification: All the information can be found in Experiment part.

328 Guidelines:

- 329 • The answer NA means that the paper does not include experiments.
- 330 • If the paper includes experiments, a No answer to this question will not be perceived
331 well by the reviewers: Making the paper reproducible is important, regardless of
332 whether the code and data are provided or not.
- 333 • If the contribution is a dataset and/or model, the authors should describe the steps taken
334 to make their results reproducible or verifiable.
- 335 • Depending on the contribution, reproducibility can be accomplished in various ways.
336 For example, if the contribution is a novel architecture, describing the architecture fully
337 might suffice, or if the contribution is a specific model and empirical evaluation, it may
338 be necessary to either make it possible for others to replicate the model with the same
339 dataset, or provide access to the model. In general, releasing code and data is often
340 one good way to accomplish this, but reproducibility can also be provided via detailed
341 instructions for how to replicate the results, access to a hosted model (e.g., in the case
342 of a large language model), releasing of a model checkpoint, or other means that are
343 appropriate to the research performed.
- 344 • While NeurIPS does not require releasing code, the conference does require all submis-
345 sions to provide some reasonable avenue for reproducibility, which may depend on the
346 nature of the contribution. For example
 - 347 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
348 to reproduce that algorithm.
 - 349 (b) If the contribution is primarily a new model architecture, the paper should describe
350 the architecture clearly and fully.
 - 351 (c) If the contribution is a new model (e.g., a large language model), then there should
352 either be a way to access this model for reproducing the results or a way to reproduce
353 the model (e.g., with an open-source dataset or instructions for how to construct
354 the dataset).
 - 355 (d) We recognize that reproducibility may be tricky in some cases, in which case
356 authors are welcome to describe the particular way they provide for reproducibility.
357 In the case of closed-source models, it may be that access to the model is limited in
358 some way (e.g., to registered users), but it should be possible for other researchers
359 to have some path to reproducing or verifying the results.

360 5. Open access to data and code

361 Question: Does the paper provide open access to the data and code, with sufficient instruc-
362 tions to faithfully reproduce the main experimental results, as described in supplemental
363 material?

364 Answer: [Yes]

365 Justification: The accessible Github link can be found in Abstract.

366 Guidelines:

- 367 • The answer NA means that paper does not include experiments requiring code.
- 368 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 369 • While we encourage the release of code and data, we understand that this might not be
370 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
371 including code, unless this is central to the contribution (e.g., for a new open-source
372 benchmark).
- 373 • The instructions should contain the exact command and environment needed to run to
374 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 375 • The authors should provide instructions on data access and preparation, including how
376 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 377 • The authors should provide scripts to reproduce all experimental results for the new
378 proposed method and baselines. If only a subset of experiments are reproducible, they
379 should state which ones are omitted from the script and why.
- 380 • At submission time, to preserve anonymity, the authors should release anonymized
381 versions (if applicable).
- 382 • Providing as much information as possible in supplemental material (appended to the
383 paper) is recommended, but including URLs to data and code is permitted.

386 6. Experimental Setting/Details

387 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
388 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
389 results?

390 Answer: [Yes]

391 Justification: The detailed information can be found in Experiment part.

392 Guidelines:

- 393 • The answer NA means that the paper does not include experiments.
- 394 • The experimental setting should be presented in the core of the paper to a level of detail
395 that is necessary to appreciate the results and make sense of them.
- 396 • The full details can be provided either with the code, in appendix, or as supplemental
397 material.

398 7. Experiment Statistical Significance

399 Question: Does the paper report error bars suitably and correctly defined or other appropriate
400 information about the statistical significance of the experiments?

401 Answer: [Yes]

402 Justification: The information can be found in Experiment part. Through the comparison
403 between the VAE and DDPM, the results are significantly increasing.

404 Guidelines:

- 405 • The answer NA means that the paper does not include experiments.
- 406 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
407 dence intervals, or statistical significance tests, at least for the experiments that support
408 the main claims of the paper.
- 409 • The factors of variability that the error bars are capturing should be clearly stated (for
410 example, train/test split, initialization, random drawing of some parameter, or overall
411 run with given experimental conditions).

- 412 • The method for calculating the error bars should be explained (closed form formula,
 413 call to a library function, bootstrap, etc.)
 414 • The assumptions made should be given (e.g., Normally distributed errors).
 415 • It should be clear whether the error bar is the standard deviation or the standard error
 416 of the mean.
 417 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 418 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 419 of Normality of errors is not verified.
 420 • For asymmetric distributions, the authors should be careful not to show in tables or
 421 figures symmetric error bars that would yield results that are out of range (e.g. negative
 422 error rates).
 423 • If error bars are reported in tables or plots, The authors should explain in the text how
 424 they were calculated and reference the corresponding figures or tables in the text.

425 **8. Experiments Compute Resources**

426 Question: For each experiment, does the paper provide sufficient information on the com-
 427 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 428 the experiments?

429 Answer: [Yes]

430 Justification: The specific information can be found in the Experiment Details part.

431 Guidelines:

- 432 • The answer NA means that the paper does not include experiments.
- 433 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 434 or cloud provider, including relevant memory and storage.
- 435 • The paper should provide the amount of compute required for each of the individual
 436 experimental runs as well as estimate the total compute.
- 437 • The paper should disclose whether the full research project required more compute
 438 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 439 didn't make it into the paper).

440 **9. Code Of Ethics**

441 Question: Does the research conducted in the paper conform, in every respect, with the
 442 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

443 Answer: [Yes]

444 Justification: Our research strictly adheres to the NeurIPS Code of Ethics.

445 Guidelines:

- 446 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 447 • If the authors answer No, they should explain the special circumstances that require a
 448 deviation from the Code of Ethics.
- 449 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 450 eration due to laws or regulations in their jurisdiction).

451 **10. Broader Impacts**

452 Question: Does the paper discuss both potential positive societal impacts and negative
 453 societal impacts of the work performed?

454 Answer: [Yes]

455 Justification: The paper gives the positive societal impacts because it discussed a new model
 456 which can produce the picture with both great quality and reconstructions. Also, improved
 457 the model to β -VAE. Enhancing the interpretability.

458 Guidelines:

- 459 • The answer NA means that there is no societal impact of the work performed.
- 460 • If the authors answer NA or No, they should explain why their work has no societal
 461 impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes, the creators or original owners of all the assets we use in our papers (such as code, data, models, etc.) have been properly recognized, and we have made explicit reference to the licenses and terms of use of these assets, ensuring that our use complies with these regulations.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 515 • If assets are released, the license, copyright information, and terms of use in the
516 package should be provided. For popular datasets, paperswithcode.com/datasets
517 has curated licenses for some datasets. Their licensing guide can help determine the
518 license of a dataset.
519 • For existing datasets that are re-packaged, both the original license and the license of
520 the derived asset (if it has changed) should be provided.
521 • If this information is not available online, the authors are encouraged to reach out to
522 the asset's creators.

523 **13. New Assets**

524 Question: Are new assets introduced in the paper well documented and is the documentation
525 provided alongside the assets?

526 Answer: [Yes]

527 Justification: The accessible link of code can be found in the abstract part.

528 Guidelines:

- 529 • The answer NA means that the paper does not release new assets.
530 • Researchers should communicate the details of the dataset/code/model as part of their
531 submissions via structured templates. This includes details about training, license,
532 limitations, etc.
533 • The paper should discuss whether and how consent was obtained from people whose
534 asset is used.
535 • At submission time, remember to anonymize your assets (if applicable). You can either
536 create an anonymized URL or include an anonymized zip file.

537 **14. Crowdsourcing and Research with Human Subjects**

538 Question: For crowdsourcing experiments and research with human subjects, does the paper
539 include the full text of instructions given to participants and screenshots, if applicable, as
540 well as details about compensation (if any)?

541 Answer: [NA]

542 Justification: The paper does not involve crowdsourcing nor research with human subjects.

543 Guidelines:

- 544 • The answer NA means that the paper does not involve crowdsourcing nor research with
545 human subjects.
546 • Including this information in the supplemental material is fine, but if the main contribu-
547 tion of the paper involves human subjects, then as much detail as possible should be
548 included in the main paper.
549 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
550 or other labor should be paid at least the minimum wage in the country of the data
551 collector.

552 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
553 Subjects**

554 Question: Does the paper describe potential risks incurred by study participants, whether
555 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
556 approvals (or an equivalent approval/review based on the requirements of your country or
557 institution) were obtained?

558 Answer: [NA]

559 Justification: The paper does not involve crowdsourcing nor research with human subjects.

560 Guidelines:

- 561 • The answer NA means that the paper does not involve crowdsourcing nor research with
562 human subjects.
563 • Depending on the country in which research is conducted, IRB approval (or equivalent)
564 may be required for any human subjects research. If you obtained IRB approval, you
565 should clearly state this in the paper.

- 566
- We recognize that the procedures for this may vary significantly between institutions
- 567 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 568 guidelines for their institution.
- 569
- For initial submissions, do not include any information that would break anonymity (if
- 570 applicable), such as the institution conducting the review.

571 **References**

- 572 Burgess, C. P., Higgins, I., Pal, A., Matthey, L., Watters, N., Desjardins, G., and Lerchner, A. (2018).
573 Understanding disentangling in β -vae. *arXiv preprint arXiv:1804.03599*.
- 574 Croitoru, F.-A., Hondu, V., Ionescu, R. T., and Shah, M. (2023). Diffusion models in vision: A
575 survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- 576 Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. (2017). Gans trained by a
577 two time-scale update rule converge to a local nash equilibrium. *Advances in neural information
578 processing systems*, 30.
- 579 Karras, T., Aila, T., Laine, S., and Lehtinen, J. (2017). Progressive growing of gans for improved
580 quality, stability, and variation. *arXiv preprint arXiv:1710.10196*.
- 581 Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint
582 arXiv:1312.6114*.
- 583 Pandey, K., Mukherjee, A., Rai, P., and Kumar, A. (2022). Diffusevae: Efficient, controllable and
584 high-fidelity generation from low-dimensional latents. *arXiv preprint arXiv:2201.00308*.
- 585 Vahdat, A., Williams, F., Gojcic, Z., Litany, O., Fidler, S., Kreis, K., et al. (2022). Lion: Latent point
586 diffusion models for 3d shape generation. *Advances in Neural Information Processing Systems*,
587 35:10021–10039.