
VAEs meet Diffusion Models: Efficient and High-Fidelity Generation

Kushagra Pandey
IIT Kanpur
kushagra@cse.iitk.ac.in

Avideep Mukherjee
IIT Kanpur
avideep@cse.iitk.ac.in

Piyush Rai
IIT Kanpur and Google Research
piyush@cse.iitk.ac.in

Abhishek Kumar
Google Research
abhishk@google.com

Abstract

Diffusion Probabilistic models have been shown to generate state-of-the-art results on several competitive image synthesis benchmarks but lack a low-dimensional, interpretable latent space, and are slow at generation. On the other hand, Variational Autoencoders (VAEs) have access to a low-dimensional latent space but, despite recent advances, exhibit poor sample quality. We present VAEDM, a novel generative framework for *refining* VAE generated samples using diffusion models while also presenting a novel conditional forward process parameterization for diffusion models. We show that the resulting parameterization can improve upon the unconditional diffusion model in terms of sampling efficiency during inference while also equipping diffusion models with the low-dimensional VAE inferred latent code. Furthermore, we show that the proposed model exhibits out-of-the-box capabilities for downstream tasks like image superresolution and denoising.

1 Introduction

Generative modeling is the task of capturing the underlying data distribution and learning to generate novel samples from a posited explicit/implicit distribution of the data in an unsupervised manner. Variational Autoencoders (VAEs) [17, 29] are a type of explicit-likelihood based generative models which can also learn a low-dimensional latent representation for the data. The resulting framework is flexible and can be used for several downstream applications [3, 11, 18, 27]. However, in image synthesis applications, VAE generated samples are usually blurry and fail to incorporate high-frequency information [7]. Despite recent advances [4, 28, 37, 39] in improving VAE sample quality, there is still a significant gap in sample quality between VAEs and their implicit-likelihood counterparts like GANs [8, 14, 15, 16]. In contrast, Denoising diffusion probabilistic models (DDPM) [12, 33] have been recently shown to achieve impressive performance on several image synthesis benchmarks, even surpassing GANs on several such benchmarks [6, 13]. However, conventional diffusion models require an expensive iterative sampling procedure and lack a low-dimensional latent representation, limiting these models' practical applicability for downstream applications.

In this preliminary work, we present VAEDM, a novel framework which combines the best of both VAEs and DDPMs. VAEDM consists of a *generator-refiner* framework in which blurry samples generated from a VAE are *refined* using a conditional DDPM. Using qualitative and quantitative experiments, we show that the proposed method, unlike standard VAEs, can generate high-quality samples that can be controlled using the low-dimensional VAE latent space. Moreover, we show that the proposed method requires fewer reverse diffusion sampling steps during inference, and exhibits *out-of-the-box generalization* to downstream tasks like image super-resolution and denoising.

2 The VAEDM framework

Before introducing the VAEDM framework, we recommend the readers to refer to Appendix A for a background on VAEs and diffusion models. Given a high-resolution image x_0 , an auxiliary conditioning signal y to be modelled using a VAE, a latent representation z associated with y , and a sequence of T representations $x_{1:T}$ learned by a diffusion model, the VAEDM generative process can be formulated as follows:

$$p(x_{0:T}, y, z) = p(z)p_\theta(y|z)p_\phi(x_{0:T}|y, z) \quad (1)$$

where θ and ϕ are the parameters of the VAE decoder and the reverse process of the conditional diffusion model, respectively. Since computation of the likelihood for this generative process is intractable, we can approximate it by computing a lower bound (ELBO) with respect to the joint posterior over the unknowns $(x_{1:T}, y, z)$ which can be formulated as follows:

$$q(x_{1:T}, z|x_0, y) = q_\psi(z|y, x_0)q(x_{1:T}|y, z, x_0) \quad (2)$$

where ψ are the parameters of the VAE recognition network ($q_\psi(z|y, x_0)$). We keep the DDPM forward process ($q(x_{1:T}|y, z, x_0)$) fixed throughout training. It can be shown that the specified probabilistic framework yields the following ELBO objective (See Appendix C.1 for proof)

$$\log p(x_0, y) \geq \underbrace{\mathbb{E}_{q_\psi(z|y, x_0)}[p_\theta(y|z)] - \mathcal{D}_{KL}(q_\psi(z|y, x_0)||p(z))}_{\mathcal{L}_{VAE}} + \underbrace{\mathbb{E}_{z \sim q(z|x_0, y)} \left[\mathbb{E}_{q(x_{1:T}|y, z, x_0)} \left[\frac{p_\phi(x_{0:T}|y, z)}{q(x_{1:T}|y, z, x_0)} \right] \right]}_{\mathcal{L}_{DDPM}} \quad (3)$$

Therefore, the overall VAEDM training objective decomposes into a sum of VAE and a conditional DDPM objectives. In addition to the above objective, we make the following simplifying assumptions:

1. We assume the conditioning signal y to be x_0 itself. Given this choice, we do not condition the reverse diffusion process on y and take it as $p_\phi(x_{0:T}|z)$.
2. For ease of optimization, we train Eq. (3) in a sequential two-stage manner, *i.e.*, first optimizing \mathcal{L}_{VAE} and then optimize for \mathcal{L}_{DDPM} while fixing θ and ψ .
3. Lastly, instead of conditioning the reverse diffusion directly on the latent code z , we condition the second stage DDPM model on the VAE reconstruction \hat{x}_0 which is a deterministic function of z .

Moreover, in this work, we only consider a single-stage VAE (with a single stochastic layer) for the VAE training stage. However, due to the flexibility of the VAEDM two-stage training, multi-stage VAE approaches as proposed in [4, 28, 37] can also be utilized. For the second stage DDPM training, we consider the following two DDPM variants in this work.

Formulation-1: In this DDPM formulation, we assume that the forward process is independent of the VAE reconstructions \hat{x} (and hence the latent code information z) *i.e.* $q(x_{1:T}|z, x_0) \approx q(x_{1:T}|x_0)$. Moreover, given the simplifying assumptions discussed previously, the reverse process transitions take the following form *i.e.* $p(x_{0:T}|z) \approx p(x_{0:T}|\hat{x}_0)$.

Formulation-2: In this DDPM formulation, we assume that the forward process is also dependent on the VAE reconstructions \hat{x} *i.e.* $q(x_{1:T}|z, x_0) \approx q(x_{1:T}|x_0, \hat{x}_0)$. The form of the reverse process remains the same as in Formulation-1. Since, the forward process is also dependent on \hat{x}_0 in addition to the input data x_0 , the form of the conditional marginal $q(x_t|x_0)$ in the forward process changes, which leads to a modified DDPM training and inference procedure. We refer interested readers to Appendix B for a detailed discussion of this formulation.

3 Experiments

We evaluate the effectiveness of our method by making qualitative and quantitative comparisons between the proposed VAEDM framework and the unconditional DDPM model [12] on the

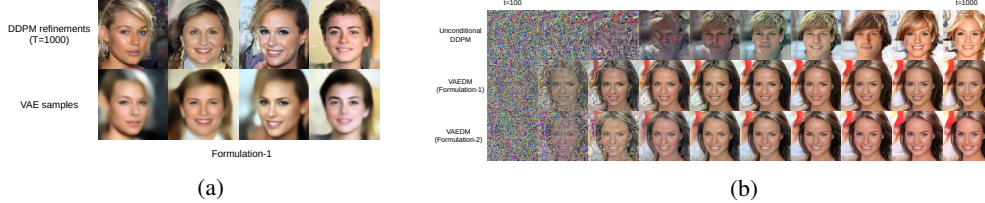


Figure 1: (Left) Illustration of the VAEDM generator-refiner framework. (Right) Illustration of sampling speed improvements in VAEDM (Best viewed with zoom-in)

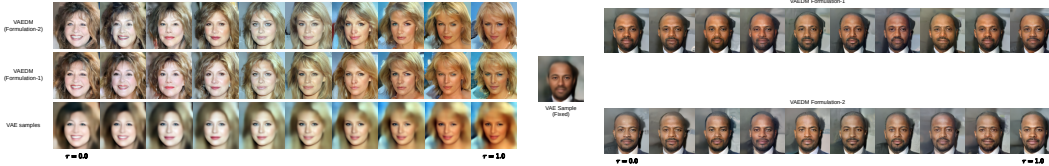


Figure 2: (Left) Illustration of the VAEDM generator-refiner framework. (Right) Illustration of sampling speed improvements in VAEDM (Best viewed with zoom-in)

CelebAMask-HQ [22] dataset. All the images in the CelebAMask-HQ dataset were downsampled to 128×128 resolution for efficiency. For all the experiments, we set number of DDPM timesteps, $T = 1000$, the noise schedule, $\beta_1 = 10^{-4}$ and $\beta_2 = 0.02$ and use our unconditional DDPM implementation to compare with the proposed approach. More details regarding the network architecture and the training and evaluation procedures can be found in Appendix E.

3.1 Sample quality

Fig. 1a shows some samples generated from the VAEDM model (using Formulation-1). As can be observed from the visualization, the final generated samples (Fig. 1a (Top row)) are a refinement of the *blurry* samples generated by our single-stage VAE model (Fig. 1a (Bottom row)). Some more unconditional samples from the VAEDM model are provided in Appendix F.

3.2 Number of Reverse process steps during sampling

Next, we qualitatively compared the unconditional DDPM model proposed in [12] and VAEDM in terms of the number of reverse process sampling steps required to produce high-fidelity samples. The visualization in Fig. 1b suggests that our method requires around half the number of reverse process steps to generate high-fidelity samples than the unconditional DDPM model. These qualitative results are further supported by Fig. 3a where VAEDM largely outperforms the unconditional DDPM model in terms of FID in the low-step regime. Intuitively, these results might not be surprising as the DDPM model in VAEDM only needs to refine a pre-generated blurry template while the unconditional DDPM in [12] requires a longer sampling schedule to model all image details, which is a harder task than refining a pre-generated template.

3.3 Interpolations in the latent space

The proposed VAEDM model consists of two latent codes: the VAE latent code z and the latent representations x_T associated with the reverse process base measure $p(x_T)$ (which is of the same size of the input image x_0 and thus might not be beneficial for downstream tasks). We next discuss the effects of manipulating both z and x_T . We consider the following two interpolation settings:

Varying z with fixed x_T : We first sample two VAE latent codes z_1 and z_2 using the standard Gaussian distribution. We then perform linear interpolation between z_1 and z_2 to obtain intermediate VAE latent codes $\tilde{z} = \lambda z_1 + (1 - \lambda)z_2$, which are then used to generate the final VAEDM samples with a shared x_T latent.

Fixed z with varying x_T : Next, we sampled the VAE latent code z using the standard Gaussian distribution. With a fixed z , we then sampled two latent DDPM representations x_{T_1} and x_{T_2} from the

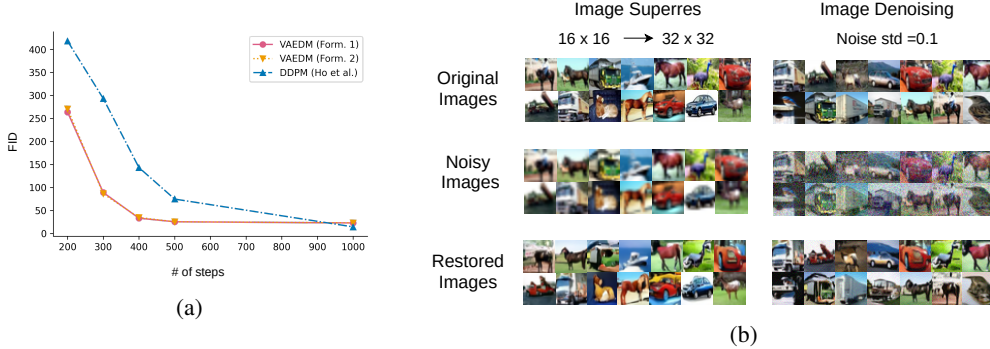


Figure 3: (a) Plot of FID vs. number of reverse process steps in VAEDM and DDPM [12] (b) Illustration of VAEDM generalization to image-superresolution and denoising (Best viewed with zoom-in)

reverse process base measure $p(x_T)$. We then performed linear interpolation between x_{T_1} and x_{T_2} with a fixed z to generate the final VAEDM samples.

As can be observed from Fig. 2, interpolating in the low-dimensional VAE latent space leads to changes in major features of the generated samples while changes in the DDPM latent conditioned on the same VAE reconstructions leads to minor differences between the generated samples which can be attributed to the stochastic nature of the DDPM sampling procedure. This observation implies that we can control the DDPM generated samples primarily by manipulating the low-dimensional VAE latent code z_{vae} and that the DDPM latents in VAEDM have low-entropy. To the best of our knowledge, our work is the first to equip diffusion models with such a low-dimensional latent representation.

3.4 Out-of-the-box super-resolution and denoising

Recently, [31] showed impressive image super-resolution results using diffusion models. To test if our generator-refiner framework can generalize over different types of noisy inputs, we conditioned the DDPM reverse process on the noisy input (instead of the VAE reconstruction). The samples obtained from such conditioning are visualized in Fig. 3. On a 16 x 16 to 32 x 32 image super-resolution task and an image denoising task (using gaussian noise corruption) for the CIFAR-10 dataset, our model is able to generate plausible reconstructions (Fig. 3b(Middle row)) when compared to the original samples. Intuitively, these results can be expected since the task of refining (blurry) VAE reconstructions might be more challenging than learning to upsample a downsampled version of an image. However, it is worth noting that certain artifacts in the generated reconstructions are evident, leaving scope for improvements.

4 Conclusion

In this work, we presented a novel unifying framework for training VAEs and diffusion models. We presented the effectiveness of the proposed approach in generating high-quality samples, requiring fewer reverse process steps during inference when compared with the unconditional DDPM formulation, equipping DDPM with a low dimensional latent code, and generalizing to additional tasks like image super-resolution and denoising. We also provide a detailed comparison of our proposed approach with other state-of-the-art generative model families in Appendix D. However, we look forward to combining our proposed approach with recent advances in diffusion models and make more quantitative comparisons with other state-of-the-art methods on image-synthesis benchmarks.

5 Impact statement

We note that the synthetic image generation techniques have the potential to mitigate bias and privacy issues for related ML models that require data collection and annotation. However, such techniques could be misused to produce fake or misleading information, and researchers should be aware of these risks and explore the techniques responsibly.

Acknowledgements

We would like to thank Ben Poole for his insightful comments during the course of this project.

References

- [1] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders, 2016.
- [2] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation, 2020.
- [3] Ricky T. Q. Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders, 2019.
- [4] Rewon Child. Very deep vaes generalize autoregressive models and can outperform them on images, 2021.
- [5] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models, 2021.
- [6] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021.
- [7] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks, 2016.
- [8] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [9] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord: Free-form continuous dynamics for scalable reversible generative models, 2018.
- [10] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning, 2020.
- [11] I. Higgins, L. Matthey, A. Pal, Christopher P. Burgess, Xavier Glorot, M. Botvinick, S. Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017.
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.
- [13] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *arXiv preprint arXiv:2106.15282*, 2021.
- [14] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation, 2018.
- [15] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks, 2019.
- [16] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan, 2020.
- [17] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- [18] Diederik P. Kingma, Danilo J. Rezende, Shakir Mohamed, and Max Welling. Semi-supervised learning with deep generative models, 2014.
- [19] Diederik P. Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improving variational inference with inverse autoregressive flow, 2017.

- [20] Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models, 2021.
- [21] Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- [22] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [23] Wonkwang Lee, Donggyun Kim, Seunghoon Hong, and Honglak Lee. High-fidelity synthesis with disentangled representation, 2020.
- [24] Vaden Masrani, Tuan Anh Le, and Frank Wood. The thermodynamic variational objective, 2021.
- [25] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks, 2018.
- [26] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021.
- [27] Adrian Alan Pol, Victor Berger, Gianluca Cerminara, Cecile Germain, and Maurizio Pierini. Anomaly detection with conditional variational autoencoders, 2020.
- [28] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2, 2019.
- [29] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows, 2016.
- [30] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015.
- [31] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *arXiv preprint arXiv:2104.07636*, 2021.
- [32] Samarth Sinha and Adji B. Dieng. Consistency regularization for variational auto-encoders, 2021.
- [33] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics, 2015.
- [34] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution, 2020.
- [35] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations, 2021.
- [36] Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder variational autoencoders, 2016.
- [37] Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder, 2021.
- [38] Rianne van den Berg, Leonard Hasenclever, Jakub M. Tomczak, and Max Welling. Sylvester normalizing flows for variational inference, 2019.
- [39] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning, 2018.
- [40] Zhisheng Xiao, Karsten Kreis, Jan Kautz, and Arash Vahdat. Vaebm: A symbiosis between variational autoencoders and energy-based models, 2021.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]** See Section 3
 - (b) Did you describe the limitations of your work? **[Yes]** See Section 3.4 and 4
 - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See Section 5
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** See Sections 2,??,??
 - (b) Did you include complete proofs of all theoretical results? **[Yes]** See Appendix B,C
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[No]** The code is currently proprietary since the work is on-going
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See Appendix E
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[No]** Most experiments are qualitative in nature
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** See Appendix E
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? **[Yes]** See [21, 22]
 - (b) Did you mention the license of the assets? **[No]** The datasets used are standard benchmarks in the field
 - (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** Generated samples from our model
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[No]** The datasets used are open-sourced under a license
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[No]** The datasets used are standard benchmarks in the field
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**

A Appendix: Background on VAE’s and Diffusion models

A.1 Variational Autoencoders

VAEs [17] are based on a simple but principled encoder-decoder based formulation which tries to maximize the evidence lower bound (ELBO) of the data log-likelihood, which is intractable to compute in general. The VAE optimization objective can be stated as follows:

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - \mathcal{D}_{KL}[q_\phi(z|x)||p(z)] \quad (4)$$

Under the amortized variational inference scheme, the approximate posterior ($q_\phi(z|x)$) and the likelihood ($p_\theta(x|z)$) distributions can be modeled using deep neural networks with parameters ϕ and θ , respectively using the reparameterization trick [17, 29]. The choice of the prior distribution $p(z)$ is flexible and can vary from a standard Gaussian [17] to more expressive priors like normalizing flows [9, 19, 38].

A.2 Denoising Diffusion Probabilistic Models

DDPMs [12, 33] are latent-variable models consisting of a forward noising process ($q(x_{1:T}|x_0)$) which gradually destroys the structure of the data x_0 and a reverse process ($p(x_{0:T})$) which learns to recover the original data x_0 from the noisy input. The forward noising process is modeled using a first-order Markov chain with Gaussian transitions and is fixed throughout training, and the noise schedules β_1 to β_T can be fixed or learned. The form of the forward process and some notable properties of the forward process conditional distributions are summarized in the equations below

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t I) \quad (5)$$

The forward process of DDPMs admits a closed form for x_t for any t , as follows

$$q(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) \quad \text{where } \alpha_t = (1 - \beta_t) \quad \text{and } \bar{\alpha}_t = \prod_t \alpha_t \quad (6)$$

The forward process posteriors are also tractable and are given by

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(\tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t) \quad (7)$$

$$\text{where } \tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}\beta_t}}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \quad \text{and } \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t \quad (8)$$

The reverse process can also be parameterized using a first-order Markov chain with a learned Gaussian transition distribution as follows

$$p(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t), \quad p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \quad (9)$$

Given a large enough T and a well-behaved variance schedule of β_t , the distribution $q(x_T|x_0)$ will approximate an isotropic Gaussian. We can generate a new sample from the underlying data distribution $q(x_0)$ by sampling a latent from $p(x_T)$ (chosen to be an isotropic Gaussian distribution) and running the reverse process. As proposed in [12], the reverse process in DDPM is trained to minimize the following upper bound over the negative log-likelihood (See [33] for detailed proofs):

$$\mathbb{E}_q \left[\mathcal{D}_{KL}(q(x_T|x_0)||p(x_T)) + \sum_{t>1} \mathcal{D}_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) - \log p_\theta(x_0|x_1) \right] \quad (10)$$

A notable aspect of the above objective is that all the KL divergences involve Gaussians and, consequently, are available in closed form. Notably, [12] parameterize the reverse process conditional $p_\theta(x_{t-1}|x_t)$ using the forward process posterior $q(x_{t-1}|x_t, x_0)$. [12] show that such a parameterization simplifies the second term in Eq. 10 at any given time-step t to the following objective in Eq. 11.

$$\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|_2^2 \quad (11)$$

where $x_t = \sqrt{\bar{\alpha}_t}x_0 + \epsilon\sqrt{1 - \bar{\alpha}_t}$ and $\epsilon \sim \mathcal{N}(0, I)$. Intuitively, this means that the reverse process in DDPM is trained to predict the noise added to the input x_0 at any time-step t . We use this *simplified* training formulation throughout our work to train all proposed parameterizations of diffusion models as [12] show that this formulation yields superior sample quality than other forms of reverse process parameterizations. For further details on the exact training and inference processes, we encourage the readers to refer to [12].

B Appendix: DDPM training and inference under Formulation 2

Algorithm 1: Training (Form. 2)	Algorithm 2: Inference (Form. 2)
<ol style="list-style-type: none"> 1 repeat 2 $x_0 \sim q(x_0)$ 3 $t \sim \text{Uniform}(\{1 \dots T\})$ 4 $\epsilon \sim \mathcal{N}(0, I)$ 5 Take gradient descent step on 6 $\ \nabla_{\theta} \ \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon + \hat{x}_0, t)\ ^2$ 7 until convergence 	<ol style="list-style-type: none"> 1 $x_T \sim \mathcal{N}(y, I)$ 2 for $t = T$ to 1 do 3 $z = \mathcal{N}(0, I)$, if $t > 1$ else 0 4 $\hat{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \hat{x}_0 - \epsilon_{\theta}(x_t, \hat{x}_0, t)\sqrt{1 - \bar{\alpha}_t})$ 5 $\hat{x}_{t-1} = \gamma_0 \hat{x}_0 + \gamma_1 x_t + \gamma_2 \hat{x}_0$ 6 $x_{t-1} = \hat{x}_{t-1} + z\hat{\sigma}_t$ 7 return $x_0 - \hat{x}_0$

The DDPM training objective proposed in [12], has the following form:

$$\mathbb{E}_q \left[\underbrace{\mathcal{D}_{KL}(q(x_T|x_0)||p(x_T))}_{L_T} + \sum_{t>1} \underbrace{\mathcal{D}_{KL}(q(x_{t-1}|x_t, x_0)||p_{\theta}(x_{t-1}|x_t))}_{L_{t-1}} - \underbrace{\log p_{\theta}(x_0|x_1)}_{L_0} \right] \quad (12)$$

Reverse Process parameterization: Following [12], we parameterize the reverse process transition $p_{\theta}(x_{t-1}|x_t)$ using the functional form of the forward process posterior $q(x_{t-1}|x_t, x_0)$. For the VAEDM formulation-2 proposed in Section 2, we design the forward process conditional distributions as follows:

$$q(x_t|x_{t-1}, \hat{x}_0) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1} + (1 - \sqrt{1 - \beta_t})\hat{x}_0, \beta_t I) \quad \text{where } t > 1 \quad (13)$$

$$q(x_1|x_0, \hat{x}_0) = \mathcal{N}(\sqrt{1 - \beta_1}x_0 + \hat{x}_0, \beta_1 I) \quad (14)$$

Given this choice of $q(x_t|x_{t-1}, \hat{x}_0)$, it can be shown that the conditional marginal distribution $q(x_t|x_0, \hat{x}_0)$ can be specified as:

$$q(x_t|x_0, \hat{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0 + \hat{x}_0, (1 - \bar{\alpha}_t)I) \quad (15)$$

The posterior distribution $q(x_{t-1}|x_t, x_0)$ will also be a Gaussian distribution with the following form:

$$q(x_{t-1}|x_t, x_0, \hat{x}_0) = \mathcal{N}(\hat{\mu}_t(x_t, x_0, \hat{x}_0), \hat{\beta}_t I) \quad (16)$$

$$\text{where } \hat{\mu}_t(x_t, x_0, \hat{x}_0) = \underbrace{\frac{\beta_t \sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} x_0 + \frac{(1 - \bar{\alpha}_{t-1}) \sqrt{\bar{\alpha}_t}}{1 - \bar{\alpha}_t} x_t}_{\hat{\mu}_t(x_t, x_0)} + \underbrace{\left(1 - \frac{(1 - \bar{\alpha}_{t-1}) \sqrt{\bar{\alpha}_t}}{1 - \bar{\alpha}_t}\right)}_{\kappa} \hat{x}_0 \quad (17)$$

$$\hat{\beta}_t = \frac{(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \beta_t \quad \text{and} \quad x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \hat{x}_0 - \epsilon \sqrt{1 - \bar{\alpha}_t}) \quad \text{where } \epsilon \sim \mathcal{N}(0, I) \quad (18)$$

Hence the forward process posterior in the VAEDM formulation is a shifted version of the forward process posterior proposed in [12]. Since the VAE reconstruction \hat{x}_0 for an image x_0 is constant during DDPM training, we can parameterize the reverse process posterior as $\hat{\mu}_{\theta}(x_t, x_0, \hat{x}_0, t) = \tilde{\mu}_{\theta}(x_t, x_0, t) + \kappa \hat{x}_0$. Additionally, we keep the variance of the reverse process conditional fixed and equal to $\hat{\beta}_t$ as proposed in [12]. Since $L_{t-1} \propto \|\hat{\mu}_t(x_t, x_0, \hat{x}_0) - \hat{\mu}_{\theta}(x_t, x_0, \hat{x}_0, t)\|^2$, the DDPM training objective in our formulation remains unchanged from the simplified denoising score matching objective proposed in [12].

Choice of the decoder, L_0 : One possible choice for the decoder is to set $p_{\theta}(x_0|x_1)$ to be a discrete independent decoder derived from the Gaussian $\mathcal{N}(\hat{\mu}_{\theta}(x_1, \hat{x}_0, 1), \hat{\beta}_1 I)$ [12]. However, at $t = 1$, we have $\hat{\mu}_{\theta}(x_1, \hat{x}_0, 1) = x_0(x_1, \hat{x}_0, \epsilon_{\theta}) + \hat{x}_0$. Therefore, to account for the VAE reconstruction bias in the final DDPM output, we set our decoder $p_{\theta}(x_0|x_1) = \mathcal{N}(\hat{\mu}_{\theta}(x_1, \hat{x}_0, 1) - \hat{x}_0, \hat{\beta}_1 I)$. Without using this adjustment, we found the final DDPM samples to be a bit blurry in our initial experiments. The final training and inference algorithms are summarized in Algorithms 1 and 2 respectively. In Algorithm 2, the coefficients γ_0, γ_1 and γ_2 denote the coefficients of the forward process posterior in Eqn. 17.

C Appendix: Detailed Proofs

C.1 Derivation of VAEDM ELBO objective

Given a high-resolution image x_0 , an auxiliary conditioning signal y to be modelled using a VAE, a latent representation z associated with y , and a sequence of T representations $x_{1:T}$ learned by a diffusion model, the DiffuseVAE generative process, $p(x_{0:T}, y, z)$ can be factorized as follows:

$$p(x_{0:T}, y, z) = p(z)p_\theta(y|z)p_\phi(x_{0:T}|y, z) \quad (19)$$

where θ and ϕ are the parameters of the VAE decoder and the reverse process of the conditional diffusion model, respectively. The log-likelihood of the training data can then be obtained as:

$$\log p(x_0, y) = \log \int p(x_{0:T}, y, z) dx_{1:T} dz \quad (20)$$

Furthermore, since the joint posterior $p(x_{1:T}, z|y, x_0)$ is intractable to compute, we approximate it using a surrogate posterior $q(x_{1:T}, z|y, x_0)$ which can also be factorized into the following conditional distributions:

$$q(x_{1:T}, z|y, x_0) = q_\psi(z|y, x_0)q(x_{1:T}|y, z, x_0) \quad (21)$$

where ψ are the parameters of the VAE recognition network ($q_\psi(z|y, x_0)$). Since computation of the likelihood in Eq. (20) is intractable, we can approximate it by computing a lower bound (ELBO) with respect to the joint posterior over the unknowns ($x_{1:T}, z$) as:

$$\log p(x_0, y) \geq \mathbb{E}_{q(x_{1:T}, z|x_0, y)} \left[\log \frac{p(x_{0:T}, y, z)}{q(x_{1:T}, z|x_0, y)} \right] \quad (22)$$

Plugging the factorial forms of the DiffuseVAE generative process and the joint posterior defined above in eqn. (22), we can simplify the ELBO as follows:

$$\log p(x_0, y) \geq \mathbb{E}_{q(x_{1:T}, z|y, x_0)} \left[\log \frac{p(x_{0:T}, y, z)}{q(x_{1:T}, z|y, x_0)} \right] \quad (23)$$

$$\geq \mathbb{E}_{q(x_{1:T}, z|x_0, y)} \left[\log \frac{p(z)p_\theta(y|z)p_\phi(x_{0:T}|y, z)}{q_\psi(z|y, x_0)q(x_{1:T}|y, z, x_0)} \right] \quad (24)$$

$$\geq \mathbb{E}_{q(x_{1:T}, z|x_0, y)} \left[\log \frac{p(z)}{q_\psi(z|y, x_0)} + \log p_\theta(y|z) + \log \frac{p_\phi(x_{0:T}|y, z)}{q(x_{1:T}|y, z, x_0)} \right] \quad (25)$$

$$\geq \mathbb{E}_{q(z|y, x_0)} \left[\log \frac{p(z)}{q_\psi(z|y, x_0)} + \log p_\theta(y|z) \right] + \mathbb{E}_{q(x_{1:T}, z|x_0, y)} \left[\log \frac{p_\phi(x_{0:T}|y, z)}{q(x_{1:T}|y, z, x_0)} \right] \quad (26)$$

$$\geq \mathbb{E}_{z \sim q(z|y, x_0)} \left[\underbrace{\mathbb{E}_{q(x_{1:T}|y, z, x_0)} \left[\log \frac{p_\phi(x_{0:T}|y, z)}{q(x_{1:T}|y, z, x_0)} \right]}_{\mathcal{L}_{\text{DDPM}}} \right] \quad (27)$$

$$+ \underbrace{\mathbb{E}_{q_\psi(z|y, x_0)} [p_\theta(y|z)] - \mathcal{D}_{KL}(q_\psi(z|y, x_0) || p(z))}_{\mathcal{L}_{\text{VAE}}} \quad (28)$$

C.2 Derivation of the forward conditional marginal in Formulation 2

Given the forward process transitions for VAEDM (Formulation-2)(See Appendix B):

$$q(x_1|x_0, \hat{x}_0) = \mathcal{N}(\sqrt{1 - \beta_1}x_0 + \hat{x}_0, \beta_1 I) \quad (29)$$

$$q(x_t|x_{t-1}, \hat{x}_0) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1} + (1 - \sqrt{1 - \beta_t})\hat{x}_0, \beta_t I) \quad (30)$$

From Eqn.(30), we can write,

$$x_t = \sqrt{1 - \beta_t}x_{t-1} + (1 - \sqrt{1 - \beta_t})\hat{x}_0 + \epsilon\sqrt{\beta_t}, \quad \text{where } \epsilon \sim \mathcal{N}(0, I) \quad (31)$$

Taking expectations both sides,

$$\mathbb{E}(x_t) = \sqrt{1 - \beta_t}\mathbb{E}(x_{t-1}) + (1 - \sqrt{1 - \beta_t})\hat{x}_0 \quad (32)$$

$$\mathbb{E}(x_t) = \sqrt{1 - \beta_t}[\sqrt{1 - \beta_{t-1}}\mathbb{E}(x_{t-2}) + (1 - \sqrt{1 - \beta_{t-1}})\hat{x}_0] + (1 - \sqrt{1 - \beta_t})\hat{x}_0 \quad (33)$$

$$\mathbb{E}(x_t) = \sqrt{(1 - \beta_t)(1 - \beta_{t-1})}\mathbb{E}(x_{t-2}) + (1 - \sqrt{(1 - \beta_t)(1 - \beta_{t-1})})\hat{x}_0 \quad (34)$$

$$\vdots \quad (35)$$

$$\mathbb{E}(x_t) = \sqrt{\prod_{t=2}^t (1 - \beta_t)}\mathbb{E}(x_1) + \hat{x}_0(1 - \sqrt{\prod_{t=2}^t (1 - \beta_t)}) \quad (36)$$

Substituting $\mathbb{E}(x_1) = \sqrt{1 - \beta_1}x_0 + \hat{x}_0$ from Eqn.(29) into the above formulation we get,

$$\mathbb{E}(x_t) = \sqrt{\prod_{t=1}^t (1 - \beta_t)}x_0 + \hat{x}_0 = \sqrt{\bar{\alpha}_t}x_0 + \hat{x}_0 \quad (37)$$

Similarly it can be shown that $Var(x_t) = (1 - \bar{\alpha}_t)I$. Therefore,

$$q(x_t|x_0, y) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0 + \hat{x}_0, (1 - \bar{\alpha}_t)I) \quad (38)$$

D Appendix: Detailed Comparison with existing approaches

Following the seminal work of [12, 33] in diffusion models, there has been a lot of recent progress in both unconditional [6, 20, 26] and conditional diffusion models [2, 5, 13, 31] (including score-based models [34, 35], based on a connection proposed in [12]) for a variety of downstream tasks including image synthesis, audio synthesis and likelihood estimation among others. Similarly there has also been progress in improving the ELBO estimates [1, 24, 32] and image synthesis [4, 23, 37, 40] quality using VAE’s [17, 29]. Next, we compare our proposed approach in detail with several of these related existing model families.

Unconditional DDPM: DDPM as introduced in [12] generates images unconditionally which has limited application scope. On the other hand, the proposed VAEDM model can be used for tasks like image enhancement, super-resolution etc. Essentially, the DDPM formulations proposed in our method can also be trained in a stand-alone manner with the reverse process conditioned on an auxiliary input like a grayscale image or an image with missing pixels to generate colorized or inpainted results, respectively. Moreover, our approach requires lesser number of reverse process steps in general to produce plausible samples.

Conditional DDPM: Conditional DDPM as introduced in [13] and [31] use multiple diffusion models for generating high-resolution images in a cascaded fashion. However, for even a two-stage pipeline the sampling time of such models would be effectively much higher than VAEDM. Given the flexibility of our approach, we hypothesize that a single-stage VAE can also be replaced by a complex multi-stage VAE architecture as proposed in [4, 37] for comparable sample quality to cascaded diffusion models without affecting the sampling time significantly. Moreover, such cascades lack a low-dimensional latent code which might be a limiting factor for certain downstream applications. It is worth noting that, [13] use a conditioning augmentation scheme where the high resolution image is generated by conditioning on a blurred/noisy low resolution image. This augmentation scheme is empirically chosen and can be sub-optimal across datasets. On the contrary, our model is conditioned on a reconstruction generated by a VAE and no explicit augmentation is required.

Hierarchical VAEs Hierarchical VAEs [4, 28, 36, 37] can suffer from posterior collapse and heuristics like gradient skipping and spectral normalization [25] might be required to stabilize training. Moreover, these models require a large dimensionality of the latent codes to generate high-fidelity samples [28, 39]. In contrast, VAEDM training does not suffer from such instabilities and the proposed method requires a single latent code layer (with dimensionality comparable to GANs) to generate high-fidelity samples. Interestingly, in principle, our approach can be considered similar to Hierarchical VAEs in some sense since the first stage VAE is used to model the high-level attributes of the generated image while the DDPM model in the second stage can be considered as a refiner that adds low-level details. In Hierarchical VAEs, the subsequent latent hierarchies perform this function. [4].

E Model Architecture and training

Unless otherwise mentioned, we re-used the model architectures and training hyperparameters proposed in [12]. For all experiments, we used the EMA approach from [10] for updating the target model parameters with an initial decay rate of 0.9999. However, we did not experiment with the choice of the EMA scheme.

Data preprocessing: For experiments with the CelebAMask-HQ dataset [22], we downsampled the training data to 128 x 128 resolution since we believe this setting provides the correct balance between visual details and training efficiency. We did not use any form of data augmentation during training for both CelebAMask-HQ and CIFAR-10 datasets. For Stage-1 VAE training, the training data was scaled between [0.0, 1.0]. For Stage-2 DDPM training, the images were scaled between [-1.0, 1.0] for training both the unconditional DDPM and the VAEDM (both formulations).

Model architecture: The VAE architecture used for Stage-1 training on the CelebAMask-HQ dataset consists of around 21M parameters with residual block architectures inspired from [4]. The VAE latent code size was set to 1024 for CelebAMask-HQ and 512 for the CIFAR-10 dataset. The U-Net [30] decoder used in Stage-2 DDPM training is around 113M parameters. For CIFAR-10, the VAE model consists of around 9M parameters whereas the size of the U-Net decoder used is 34.8M parameters.

Training: The VAE model was trained for around 180 epochs and 500 epochs for the CelebAMask-HQ and the CIFAR-10 datasets with batch sizes of 40 and 128 respectively. All DDPM models were trained for around 250 epochs with a batch size of 64 for the CelebAMask-HQ dataset and for around 1000 epochs with a batch size of 128 for the CIFAR-10 dataset. We used a mix of 4 Nvidia 1080Ti GPUs (44GB memory) and a cloud TPUv2-8 (64GB memory) provided by Colab Pro for training the models. For DDPM training, it took around 0.74s and 0.2s per optimization step for the CelebAHQ and CIFAR-10 datasets respectively.

F Unconditional samples generated using VAEDM



Figure 4: Selected unconditional samples (128 x 128) generated from VAEDM using formulation 2. (t=500)



Figure 5: Selected unconditional samples (128 x 128) generated from VAEDM using formulation 1. (t=500)