

# DREAMGUIDER: IMPROVED TRAINING FREE DIFFUSION-BASED CONDITIONAL GENERATION

Anonymous authors

Paper under double-blind review

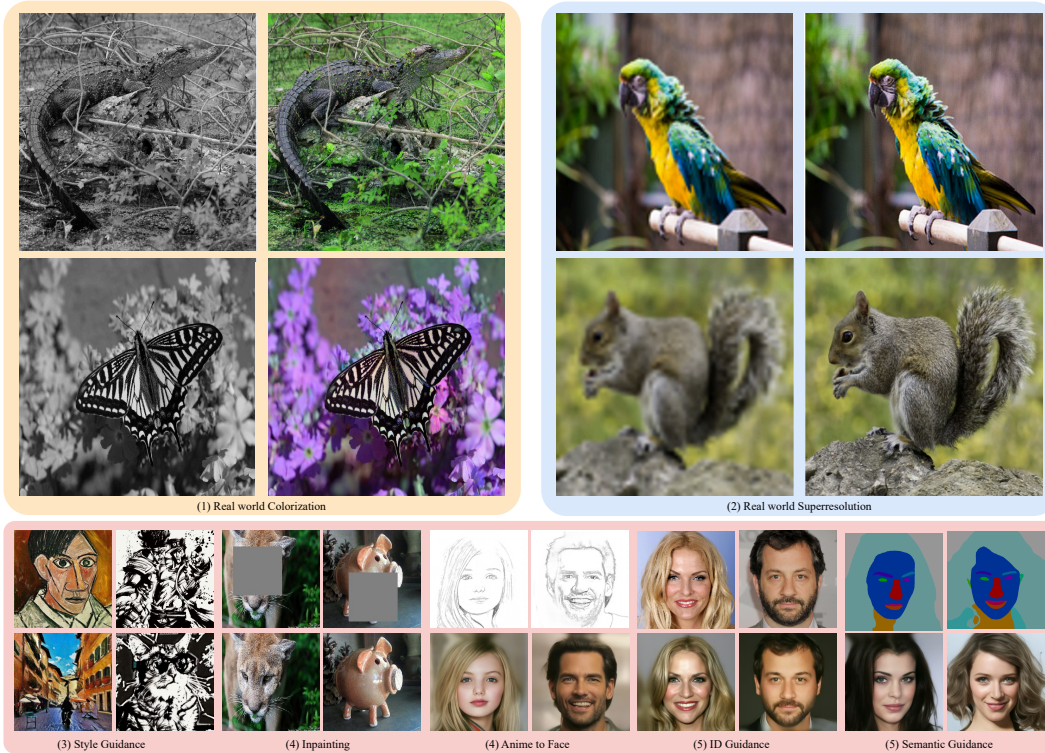


Figure 1: An illustration of the different applications of our method. We utilize a pretrained diffusion model to generate images satisfying a predefined condition without backpropagation through the diffusion UNet or any hand-crafted parameter tuning. We present results on (1) Real-world colorization, (2) Real-world super-resolution, (3) Style-guided Text-to-Image Generation, (4) Inpainting, (5) Sketch-to-Face, (6) Face ID Guidance, and (7) Face Semantics-to-Face synthesis.

## ABSTRACT

Diffusion models have emerged as a formidable tool for training-free conditional generation. However, a key hurdle in inference-time guidance techniques is the need for compute-heavy backpropagation through the diffusion network for estimating the guidance direction. Moreover, these techniques often require handcrafted parameter tuning on a case-by-case basis. Although some recent works have introduced minimal compute methods for linear inverse problems, a generic lightweight guidance solution to both linear and non-linear guidance problems is still missing. To this end, we propose Dreamguider, a method that enables inference-time guidance without compute-heavy backpropagation through the diffusion network. The key idea is to regulate the gradient flow through a time-varying factor. Moreover, we propose an empirical guidance scale that works for a wide variety of tasks, hence removing the need for handcrafted parameter tuning. We further introduce an effective lightweight augmentation strategy that significantly boosts the performance during inference-time guidance. We present experiments using Dreamguider on multiple tasks across multiple datasets and models to show the effectiveness of

054 the proposed modules. To facilitate further research, we will make the code public  
055 after the review process.  
056  
057

## 058 1 INTRODUCTION 059

060 Generative modeling utilizing Denoising Diffusion Probabilistic Models (DDPMs) [Sohl-Dickstein](#)  
061 [et al. \(2015\)](#); [Ho et al. \(2020\)](#); [Dhariwal & Nichol \(2021\)](#); [Song et al. \(2021b\)](#) has massively improved  
062 over the past few years. Multiple works have extended the use of diffusion models for text-to-image  
063 synthesis [Balaji et al. \(2022\)](#); [Rombach et al. \(2021\)](#); [Saharia et al. \(2022b\)](#), 3D synthesis [Poole](#)  
064 [et al. \(2022\)](#); [Jun & Nichol \(2023\)](#), video generation [Ho et al. \(2022\)](#); [Blattmann et al. \(2023\)](#); [Wu](#)  
065 [et al. \(2023a\)](#), as well as for conditioning to solve inverse problems. Moreover, like conditional  
066 generative adversarial networks (GANs) [Goodfellow et al. \(2020\)](#); [Arjovsky et al. \(2017\)](#), DDPMs  
067 can be adapted to tasks based on a labels [Rombach et al. \(2021\)](#); [Dhariwal & Nichol \(2021\)](#) or visual  
068 prior-based conditioning [Saharia et al. \(2022a\)](#). However, like conditional GANs [Wang et al. \(2018\)](#);  
069 [Radford et al. \(2015\)](#), DDPMs also need to be trained with annotated pairs of labels and instructions  
070 to obtain satisfactory results. This poses a limitation in many cases where there is a lack of paired  
071 data to train large diffusion models. Due to this reason, there has been recent interest in models that  
072 can perform conditional generation without the need for explicit training [Yu et al. \(2023\)](#); [Chan et al.](#)  
[\(2016\)](#); [Nguyen et al. \(2017\)](#); [Graikos et al. \(2022\)](#).

073 Progressing towards this direction is prior research in plug-and-play models. First introduced in  
074 [Nguyen et al. \(2017\)](#), the initial research on plug-and-play models [Nguyen et al. \(2017\)](#); [Graikos](#)  
075 [et al. \(2022\)](#) enabled conditional sampling from GANs trained with unlabeled data. For this, a  
076 pre-trained classifier [Simonyan & Zisserman \(2014\)](#); [Hossain et al. \(2019\)](#) or a captioning model was  
077 used to estimate the deviation between the GAN-generated image and a given label, and based on  
078 this deviation, the GAN input noise was modulated until the generated sample satisfied the given  
079 text or class label. A similar approach that has been attempted for diffusion models to facilitate  
080 conditional sampling from unconditional diffusion models is classifier guidance [Dhariwal & Nichol](#)  
081 [\(2021\)](#); [Graikos et al. \(2022\)](#), where a noise-robust classifier is trained along with the diffusion  
082 model to guide the sampling towards a particular direction. However, classifier guidance brings  
083 in the computational costs of training a classifier, which is often undesirable. Some recent works  
084 have performed conditional generation without explicit training for the condition by utilizing the  
085 implicit guidance capabilities of the diffusion model [Chung et al. \(2023b\)](#); [Yu et al. \(2023\)](#); [Nair](#)  
086 [et al. \(2023\)](#); [Bansal et al. \(2023\)](#); [Chung et al. \(2023a\)](#). Diffusion posterior sampling (DPS) [Chung](#)  
087 [et al. \(2023b\)](#) proposed a technique of using an  $L_2$  norm-based loss function to solve linear inverse  
088 problems using unconditional diffusion models. However, DPS often requires a large number of  
089 sampling steps for photorealistic results. Freedom [Yu et al. \(2023\)](#), yet another work, proposed the  
090 use of general loss functions during sampling to achieve training-free conditional sampling. Some  
091 variants of DPS have also been proposed in the literature [Song et al. \(2023\)](#). All the aforementioned  
092 loss-guided posterior sampling techniques involve a guidance function at each timestep that requires  
093 backpropagation through the diffusion UNet. Recently, [He et al. \(2023\)](#) proposed Manifold Preserving  
094 Guided Diffusion Models (MGD) that remove the need for backpropagating through the diffusion  
095 U-Net by performing a gradient descent with respect to the Minimum Mean Square Error (MMSE).  
096 Although MGD [He et al. \(2023\)](#) works remarkably well for linear tasks that require more guidance  
097 towards the start of the guidance process, it may fail in some tasks where guidance happens earlier,  
098 for example, face semantics-to-image and sketch-to-image, where stronger guidance is required  
099 from a much earlier stage. Moreover, like [Yu et al. \(2023\)](#); [Nair et al. \(2023\)](#), MGD also requires a  
100 case-by-case handcrafted parameter. Hence, a generic lightweight method that works well for both  
linear and non-linear guidance functions is still missing. Moreover, the need to find a handcrafted  
guidance parameter on a case-by-case basis still remains an open challenge.

101 In this paper, we introduce a new framework that can adaptively perform zero-shot generation using  
102 diffusion models without the need for any manual intervention by the user. We found a rather  
103 simple fix to the problem during the initial timesteps of diffusion, i.e., by utilizing the gradient  
104 with respect to the diffusion output noise in initial steps of inference. Combined with the guidance  
105 with respect to the MMSE estimate, we found that the combination generalizes well to tasks that  
106 require guidance at very early stages of guidance. Figure 2 presents the visualization of our approach  
107 over existing works present in the literature. Utilizing the correction term along with the correction  
with respect to the MMSE estimate significantly boosts the performance in non-linear tasks. We

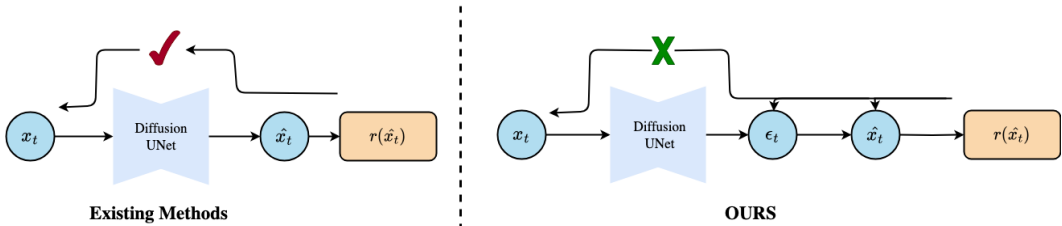


Figure 2: An illustration of the difference between the existing method and our method. Existing works backpropagate through the diffusion network to perform guidance at each timestep, whereas we find the gradients with respect to the MMSE estimate and the predicted noise based on the timesteps, thereby bypassing the expensive backpropagation operation.

Table 1: Table illustrating the difference over existing methods performing inference-time guidance.

Method	Zeroth order	Linear Tasks	Non-Linear Tasks	Automatic scaling
DPS <a href="#">Chung et al. (2023a)</a>	✗	✓	✗	✗
$\pi$ GDM <a href="#">Song et al. (2022)</a>	✗	✓	✗	✗
Freedom <a href="#">Yu et al. (2023)</a>	✗	✗	✓	✗
MGD <a href="#">He et al. (2023)</a>	✓	✓	✗	✗
OURS	✓	✓	✓	✓

present the corresponding results in Section 5. Moreover, we treat the energy-based inference-time guidance [Chung et al. \(2023b\)](#); [Yu et al. \(2023\)](#) as a stochastic gradient optimization of the MMSE estimate and the noise present in the image. This formulation enabled us to leverage recent research in parameter-free learning [Defazio & Mishchenko \(2023\)](#); [Ivgi et al. \(2023\)](#) to develop a dynamic step size schedule. This step size adjusts itself adaptively based on the initial noise seed input of the diffusion model and guidance functions, hence removing the need for manual parameter tuning for inference-time guidance. Moreover, motivated by the effectiveness of differentiable augmentations while training GANs [Zhao et al. \(2020\)](#), we found that utilizing multiple levels of matching differentiable augmentations to the MMSE estimate and guidance reference significantly improves the sampling quality, enabling very high-quality sampling with a low number of guidance steps. We present an overview of the different applications of our method in Figure 1 and an illustration of the difference of dreamguider with existing methods in Table 1. Namely, we present results using Stable Diffusion [Rombach et al. \(2021\)](#), unconditional diffusion models released by [Nichol & Dhariwal \(2021\)](#) for  $256 \times 256$  guidance, and class-conditional diffusion models for high-resolution  $512 \times 512$  conditional synthesis. The different functionalities of Dreamguider are tabulated in Figure 2.

We present experiments on publicly released models on generic images, face images, and stable diffusion to show the relevance of our method. We focus on the tasks of (1) Inpainting, (2) Super-resolution, (3) Colorization, (4) Gaussian Deblurring, (5) Semantic label-to-image generation, (6) Face sketch-to-image, (7) ID guidance and identity generation, and beat existing benchmarks that utilize diffusion models for these tasks, obtaining a significant boost in performance over existing methods leveraging loss-guided models. To summarize, our contributions are:

- We propose a zeroth-order loss-guided diffusion guidance that is applicable to both linear inverse problems and non-linear inverse problems.
- We remove the need for a manually tuned guidance scale for classifier guidance by proposing a scaling function that works for a wide variety of tasks.
- We propose a time-varying guidance scale for improving sampling quality.
- We propose a differentiable augmentation strategy to improve sampling quality.

## 2 BACKGROUND

### 2.1 TRAINING-FREE CONDITIONAL SAMPLING USING DIFFUSION MODELS

Recently, there has been a rise in multiple works that propose utilizing unconditional diffusion models for conditional sampling [Nair et al. \(2023\)](#); [Bansal et al. \(2023\)](#); [Chung et al. \(2023c\)](#); [Kawar et al. \(2022\)](#). The earlier works proposed solving linear inverse problems using diffusion

models with the help of priors dependent on the inverse transform of degradation. Recently, diffusion posterior sampling [Chung et al. \(2023b\)](#) considered the degradation to be conditioned on a Gaussian distribution given any intermediate timestep and derived an  $L_2$  norm-based regularization at each intermediate timestep to solve for linear inverse problems. Recent works such as Freedom [Yu et al. \(2023\)](#) explored an energy-based perspective and extended guidance to non-linear functions using general loss functions. Universal diffusion guidance [Aggarwal et al. \(2018\)](#) extended this guidance process to stable diffusion and improved the performance by using forward-backward guidance. More recent works, such as manifold-guided diffusion [He et al. \(2023\)](#), further proposed to constrain the manifold space by projecting for the latent space alone.

## 2.2 PERTURBED MARKOVIAN KERNEL FOR DIFFUSION TRANSITION

Let us assume that  $r(x_t, y)$  gives a measure of the distance between an intermediate  $x_t$  and the condition  $y$  and is a positive bounded function. Hence, in the reverse process, the diffusion trajectory should proceed through distributions with a higher probability of being closer to the desired cases. We model these trajectory intermediate distributions with

$$\hat{p}(x_t) = p(x_t)r(x_t, y). \quad (1)$$

Dickenson et al. [Sohl-Dickstein et al. \(2015\)](#) first proposed the use of Markovian kernels to estimate the distribution of diffusion intermediates. Specifically, given the state  $x_t$  at the equilibrium of the training process for a diffusion model, the intermediate of a diffusion model at a time instant, the distribution at a timestep  $t - 1$  can be estimated as

$$p(x_{t-1}) = \int p(x_t)p_\theta(x_{t-1}|x_t)dx_t. \quad (2)$$

As we know, the kernel  $p(x_{t-1}|x_t)$  is a Gaussian distribution whose mean can be estimated using the diffusion UNet and  $x_t$ . To estimate a perturbed kernel  $\hat{p}_\theta(x_{t-1}|x_t)$ , the perturbed distribution is

$$p(x_{t-1})r(x_{t-1}, y) = \int r(x_t, y)p(x_t)\hat{p}_\theta(x_{t-1}|x_t)dx_t. \quad (3)$$

By merging the constant terms in the transition into the normalization factor, the transition step is

$$\hat{p}_\theta(x_{t-1}|x_t) = p_\theta(x_{t-1}|x_t)r(x_{t-1}, y). \quad (4)$$

The proof is given in the supplementary material. Hence, we can see that rather than considering a Gaussian posterior, as in DPS [Chung et al. \(2023b\)](#), any distance or loss function can be used. Similarly, one other valid transition step of the perturbed process is

$$\hat{p}_\theta(x_{t-1}|x_t) = p_\theta(x_{t-1}|x_t)\frac{r(x_{t-1}, y)}{r(x_t, y)}, \quad (5)$$

which adopts the notion of reciprocal distance from the previous timestep.

## 2.3 INFERENCE-TIME GUIDANCE OF DIFFUSION MODELS

For conditional generation tasks using an unconditional diffusion model, ideally, the model would predict intermediates closer to the condition. The formulation can be seen in terms of transition probabilities. Consider a pretrained unconditional diffusion model on a specific domain. The problem at hand needs to guide the diffusion model during inference time conditioned with a condition  $y$ . Dhariwal et al. [Dhariwal & Nichol \(2021\)](#) proposed a general strategy to perform this by conditioning on the condition  $y$  and finding the resultant marginal distribution

$$p(x_t|x_{t+1}, y) = p(x_t|x_{t+1})p(y|x_t). \quad (6)$$

By assuming the distribution  $p(y|x_t)$  has much lower curvature compared to  $p(x_t|x_{t+1})$ , considering the marginal distribution close to  $x_t$ ,

$$\begin{aligned} \log p(y|x_t) &= (x_t - \mu)\nabla_{x_t}\log p(y|x_t), \\ g &= \nabla_{x_t}\log p(y|x_t). \end{aligned} \quad (7)$$

Plugging back to  $\log(p(x_t|x_{t+1}, y))$ ,

$$\begin{aligned} \log(p(x_t|x_{t+1}, y)) &= (x - \mu - \Sigma g)^T \Sigma^{-1} (x - \mu - \Sigma g) + C, \\ p(x_t|x_{t+1}, y) &\sim N(\mu + \Sigma g, \Sigma). \end{aligned} \quad (8)$$

Hence, the reverse sampling equation becomes,

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t) \right) + \sigma_t \epsilon + \Sigma \frac{dr(x_{t-1}, y)}{dx_{t-1}}, \epsilon \sim \mathcal{N}(0, I). \quad (9)$$

## 2.4 SHORTCOMINGS OF THE EXISTING METHODS

Although the energy-based guidance theory supports guidance as a function of the current latent estimate, almost all loss-based guidance techniques derive the distance function as a function of  $x_t$  rather than  $x_{t-1}$  and derive the gradient based on the previous sample. Although this approach works for many tasks, it requires backpropagating through the neural network and modeling the score function for the guidance correction term. This limits the use of classifier guidance since existing diffusion architectures that produce photorealistic results are often very bulky. One can see why the existing framework utilizes the derivative with respect to the previous sample works by taking a better look at Equation (5). As we can see, a reciprocal distance over the previous timestep diffusion latent  $x_t$  is a perfectly valid distance guidance function. In the next section, we elaborate on Dreamguider.

## 3 PROPOSED METHOD

Suppose  $x_{t-1}$  denotes the current step and  $x_t$  denotes the previous step in the inference process of the diffusion module. As mentioned in the previous section, existing works utilize the derivative with respect to the previous step for guidance; one reason for this is to use an off-the-shelf auxiliary distance function on the MMSE estimate at each step  $\hat{x}_t$ , which enables the use of general functions defined on image space for guidance. Here, the MMSE estimate is defined as

$$\hat{x}_t = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t)}{\sqrt{\bar{\alpha}_t}}, \quad (10)$$

where  $\bar{\alpha}$  denotes the variance schedule of the diffusion process and  $\epsilon_\theta(x_t)$  is the noise estimated by the network. One other observation to note is that finding the derivative with respect to the current step requires finding  $\hat{x}_{t-1}$ , which again requires an additional propagation through the diffusion network. Hence, the dilemma of backpropagating through the UNet for guidance still remains unresolved.

### 3.1 TIME VARIANT CLASSIFIER GUIDANCE

We found a simple yet effective solution for this dilemma; if we take a look at the ODE estimate at each step proposed by Song et al. [Song et al. \(2021a\)](#). Hence, rather than perturbing the Gaussian kernel at each timestep, we perturb the components  $\hat{x}_t$  and  $\epsilon_\theta(x_t)$  by a small amount. Specifically, we perform the following operations:

$$\begin{aligned} \hat{x}_t &= \hat{x}_t - c \Sigma \frac{dr(\hat{x}_t, y)}{d\hat{x}_t}, t > t_0 \\ \epsilon_\theta(x_t) &= \epsilon_\theta(x_t) - d \Sigma \frac{dr(\hat{x}_t, y)}{d\epsilon_\theta(x_t)}, t < t_0 \\ x_{t-1} &= \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t) \right) + \sigma_t \epsilon - c_t \Sigma \frac{dr(\hat{x}_t, y)}{d\hat{x}_t} - d_t \Sigma \frac{dr(\hat{x}_t, y)}{d\epsilon_\theta(x_t)} \end{aligned} \quad (11)$$

where  $r(\hat{x}_t, y)$  is a non negative distance function that measures the distance between the MMSE estimate and condition,  $\Sigma$  is the variance of the latent estimate at each timestep as in Equation (8). Please note that we perform a double descent here. The intuition behind the double descent is that performing descent on one of the components, say  $\hat{x}_t$ , guides effectively at the end of the diffusion process where  $\alpha_{t-1}$  is one and vice versa. Hence, during the guidance with the gradient w.r.t.  $\hat{x}_t$ , the maximum component of shift that happens to the sample is when we consider the flow of this correction through  $\hat{x}_t$ . Hence, we define the value as the maximum component of  $x_{t-1}$  present in  $\hat{x}_t$ .

$$c_t = c \sqrt{\alpha_{t-1}}. \quad (12)$$

Similarly, we define  $d_t$  as the maximal component of  $\epsilon_\theta(x_t)$  in  $x_{t-1}$ . Hence,

$$d_t = -d \cdot \frac{1 - \alpha_t}{\sqrt{\alpha_t} \sqrt{1 - \bar{\alpha}_t}}. \quad (13)$$

Hence, this term gives efficient guidance at all timesteps, bypassing the guidance at the later timesteps alone as in MGD He et al. (2023). In the following section, we proceed to propose an effective empirical estimate for  $c$  and  $d$  that works for a wide range of tasks.

### 3.2 A GRADIENT-DEPENDENT SCALING FACTOR ESTIMATE

Recently, Distance over Gradients (DOG) Ivgi et al. (2023) was proposed as an effective parameter-free dynamic step size schedule for SGD problems. Given any Stochastic Gradient Descent (SGD) optimization problem, the Distance over Gradient works as an effective learning rate. Recent works Wu et al. (2023b) have found the diffusion process as a stochastic optimization problem and have derived an SGD-based interpretation of the diffusion sampling process. Hence, inspired by both of these works, we attempted an empirical guidance estimate of the form:

$$\gamma_t = \begin{cases} \frac{1\epsilon^{-5}}{\sqrt{g_t^2}}, & \text{if } t = T \\ \frac{\max_{i>t} |f_i - f_T|}{\sqrt{\sum_{i=i}^T g_i^2}}, & \text{otherwise} \end{cases} \quad (14)$$

where  $g_t$  is the gradient of the loss function as defined in the equation,  $f_t$  can be any of  $\hat{x}_t, x_t, \epsilon_\theta(t)$  at timestep  $t$  and  $f_0$  is the initial estimate of  $f_t$ . We noticed that this empirical estimate works well for the first-order sampling involving DPS Chung et al. (2023b) as well. We illustrate more results on the effect of this plug-in value for different cases in the appendix. Hence, utilizing Equation (14), we estimate  $c$  and  $d$  accordingly by substituting  $f_i$  as  $\hat{x}_t$  and  $\epsilon_\theta(x_t)$

### 3.3 DIFFERENTIAL AUGMENTATION CLASSIFIER GUIDANCE

A common practice while performing classifier guidance to augment diffusion models with specific regularization for guidance is to use the noisy estimate at timestep  $t$  and utilize it to compute the loss function to regularize the current prediction. However, in many cases, such guidance can give results with artifacts and color shifts, as portrayed in Figure 3 and Figure 5, due to excessive guidance or insufficient guidance at intermediate timesteps that shift the results off manifold or cause color shifts. One effective solution for this is to imitate different versions of artifacts or color shifts on both the source image and the target image and utilize these augmented versions for a boost in performance. Hence, to perform guidance with a much more robust guidance loss, we introduce DiffuseAugment, an augmentation strategy for diffusion guidance during inference time. Specifically, given an intermediate sample  $x_t$  and condition  $y$ , we augment  $\hat{x}_t$  and  $y$  with differentiable augmentations denoted by

$$\hat{x}_t^{aug}, y^{aug} = T(\hat{x}_t^{aug}, y^{aug}). \quad (15)$$

We choose three different types of augmentations for  $T$  comprising random cutouts, random translations, and color saturations. Please note that the augmentation of  $y$  is dependent on the input signal. For label-based conditioning such as identity or text, we do not perform augmentation for  $y$ . For image space augmentations, we augment  $y$  with the same random augmentation as that of  $x$ . While computing the effective loss, we find the average across all augmentations. We find that DiffuseAugment significantly boosts the sampling fidelity and quality of the reconstructed image. We present these results in Section 5.

## 4 EXPERIMENTS

Since our method comprises both linear and non-linear inverse tasks, for linear inverse tasks, we follow DPS and evaluate our method utilizing two different benchmarks: (1) ImageNet Deng et al. (2009) and (2) CelebA Liu et al. (2015). For non-linear tasks, we follow Freedom and evaluate using the CelebA dataset. For linear tasks, we evaluate our method quantitatively for Super-resolution ( $\times 4$ ), Colorization, Inpainting (Box), and Gaussian deblurring tasks. For non-linear tasks, we evaluate

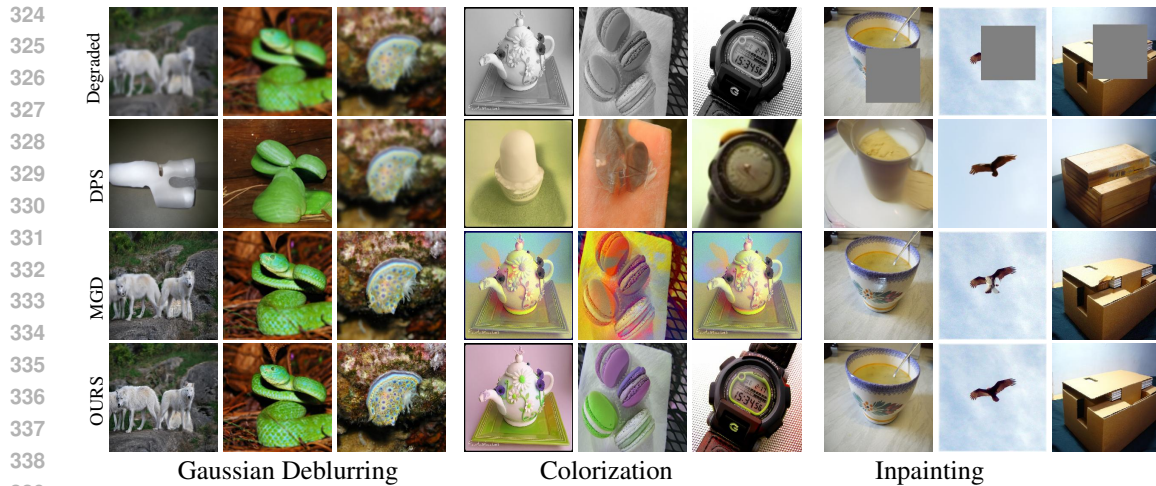


Figure 3: Qualitative comparisons for Linear Tasks on ImageNet for 100 inference steps

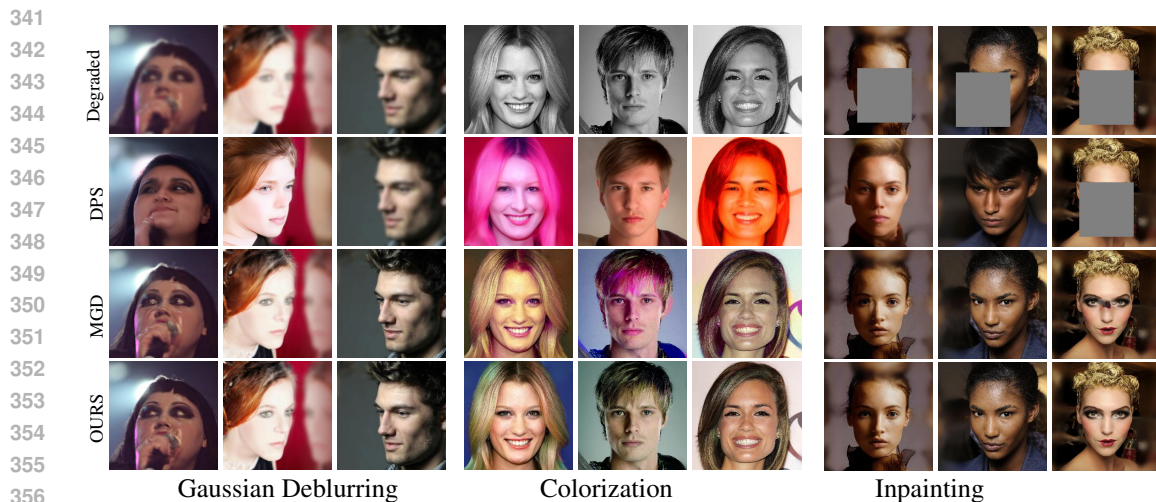


Figure 4: Qualitative comparisons for Linear Tasks on CelebA dataset for 100 inference steps

357  
358  
359  
360  
361  
362  
363  
364  
365  
366

for Face Sketch guidance, Face Parse maps guidance, and Face ID guidance. Since our method falls into the category of loss-guided diffusion models, we perform all quantitative evaluations using existing methods that follow this kind of sampling. Please note that although we acknowledge the parallel field of research in tackling inverse problems without backpropagation Wang et al. (2023); Kawar et al. (2021), we excluded these methods for comparison as they tackle solely Linear inverse problems. In contrast, loss-guided models are generic and applicable to a wider range of problems.

#### 367 4.1 IMPLEMENTATION DETAILS

368  
369  
370  
371  
372  
373  
374  
375  
376  
377

We perform all experiments on NVIDIA A6000 GPUs. For ImageNet Deng et al. (2009) based tasks, we utilize the unconditional model released by Guided Diffusion. For Linear Tasks involving faces, we use the model trained on the FFHQ dataset Karras et al. (2017) and perform experiments on the CelebA dataset Liu et al. (2015) similar to DPS. For non-linear tasks, we follow Freedom and utilize the model trained unconditionally on the CelebA dataset. We evaluate using conditions derived from existing networks. For the high-resolution results presented in Figure 2, we utilized the class-conditional model of resolution  $512 \times 512$  released by Guided Diffusion. For all experiments, we used 100 sampling steps. For style transfer, we utilized Stable Diffusion Rombach et al. (2021) v1.5. Please note that our sampling method is generic, and any sampler can be used. We fix the number of augmentations in DiffuseAugment for all the experiments to 8. For linear inverse problems we set the value of  $t_0$  to 5 in Equation (11) to 30 and for linear inverse problems we set  $t_0$  to 5

Method	Inpaint (Box)				Colorization				SR ( $\times 4$ )				Gaussian Deblur			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	Cons $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Score-SDE Song et al. (2021b)	9.57	0.329	0.634	94.33	0.1627	0.3996	0.6609	118.86	20.75	0.5844	0.3851	53.22	23.39	0.632	0.361	66.81
ILVR Song et al. (2021b)	-	-	-	-	-	-	-	-	26.14	0.7403	0.2776	52.82	-	-	-	-
DPS Chung et al. (2023a)	19.39	0.610	0.3766	58.89	0.0069	0.5404	0.5594	55.61	17.36	0.4969	0.4613	56.08	20.52	0.5824	0.3756	52.64
MGD Chung et al. (2023a)	27.21	0.7460	0.2197	11.83	0.0018	0.6865	0.4549	38.22	27.51	0.7852	0.2464	60.21	27.23	<b>0.7695</b>	0.2327	51.59
Ours	<b>28.84</b>	<b>0.8491</b>	<b>0.1432</b>	<b>5.96</b>	<b>0.0014</b>	<b>0.7775</b>	<b>0.3036</b>	<b>20.89</b>	<b>29.47</b>	<b>0.8429</b>	<b>0.1757</b>	<b>46.95</b>	<b>27.30</b>	0.7672	<b>0.2202</b>	<b>42.70</b>

Table 2: Quantitative evaluation of image restoration tasks on CelebA 256 $\times$ 256-1k with  $\sigma_y = 0.05$ , We utilize 100 inference steps for all methods

Method	Inpaint (Box)				Colorization				SR ( $\times 4$ )				Gaussian Deblur			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	Cons $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
Score-SDE Song et al. (2021b)	9.66	0.2087	0.7375	133.54	0.1723	0.3105	0.8197	194.87	14.07	0.2468	0.6766	129.91	15.39	0.3158	0.620	134.67
ILVR Song et al. (2021b)	-	-	-	-	-	-	-	-	15.51	0.4033	0.5253	64.13	-	-	-	-
DPS Chung et al. (2023a)	15.23	0.4261	0.6087	97.90	0.021	0.3774	0.8011	106.25	14.94	0.3258	0.6594	87.26	17.19	0.3980	0.5817	84.74
MGD Chung et al. (2023a)	21.94	0.6920	0.2410	40.30	0.0057	0.5809	0.5427	73.75	23.12	0.6025	0.3936	70.83	23.13	0.6092	0.3695	61.49
Ours	<b>23.49</b>	<b>0.7271</b>	<b>0.2001</b>	<b>30.72</b>	<b>0.0055</b>	<b>0.6804</b>	<b>0.3362</b>	<b>52.76</b>	<b>24.23</b>	<b>0.6818</b>	<b>0.2884</b>	<b>43.00</b>	<b>23.31</b>	<b>0.6157</b>	<b>0.3566</b>	<b>58.38</b>

Table 3: Quantitative evaluation of image restoration tasks on ImageNet 256 $\times$ 256-1k with  $\sigma_y = 0.05$ . **Bold**: best, We utilize 100 inference steps for all methods

## 4.2 QUALITATIVE ANALYSIS

We present results on Gaussian Deblurring, super-resolution, and colorization. As we can see, DPS fails since 100 steps of diffusion are used, and the DPS scaling factor is not strong enough to perform proper guidance within 100 steps of diffusion. We set the amount of posterior noise for the measurement as 0.05 in all experiments. MGD works remarkably well for the deblurring and inpainting tasks; however, it fails for colorization since early guidance is required for the flow of natural colors.

For ImageNet tasks, the performance of DPS falls more because the problem is more ill-posed. This can be seen in the eagle diagram, where the method is unable to reconstruct the eagle properly. In contrast, our method performs relatively better, producing much more realistic images. We highlight the performance improvement on colorization since we argue that these results are obtained because of the early flow of gradients. For non linear inverse problems, as we can see, Freedom is able to produce realistic-looking results for even the difficult task of Parse Maps to Faces. We argue that this is because backpropagation through the UNet purifies the gradient flow; hence, the generated images look much more naturalistic.

## 4.3 QUANTITATIVE ANALYSIS

We utilize Dreamguider and quantitatively evaluate CelebA and ImageNet datasets. The results for face restoration tasks are shown in Table 2 and 3. We evaluate these tasks utilizing four different metrics. SDEdit Meng et al. (2021) fails for the task of face inpainting and colorization as a single perturbation in the noisy domain throws the image off the manifold. DPS requires more inference steps for proper guidance. ILVR is originally designed for super-resolution. Hence, we quantitatively evaluate ILVR Choi et al. (2021) only for the task of super-resolution. Since DPS and MGD are applicable to all cases, we evaluate with these methods. As we can see, our approach obtains better results than the baselines because of the flow of gradients, which allows for better reconstruction quality. For faces, the difference is much more highlighted in the task of colorization, where we get a significant boost of 18 FID score above the baseline. General linear inverse problems in ImageNet are much more complex than in faces; hence, there is an overall drop in metrics for the natural domain images in ImageNet. In our case, DiffAugment purifies the gradient; hence, we look for much better realistic-looking images. However, MGD does not produce realistic results for sketch-to-image and anime-to-face synthesis.

## 5 ABLATION STUDIES

We perform extensive ablation studies with respect to the effect of DiffuseAugment as well as the effect of each guidance term. For the ablation experiments, rather than utilizing the whole testing dataset of 1000 images, we utilize 100 images and report the average LPIPS value.



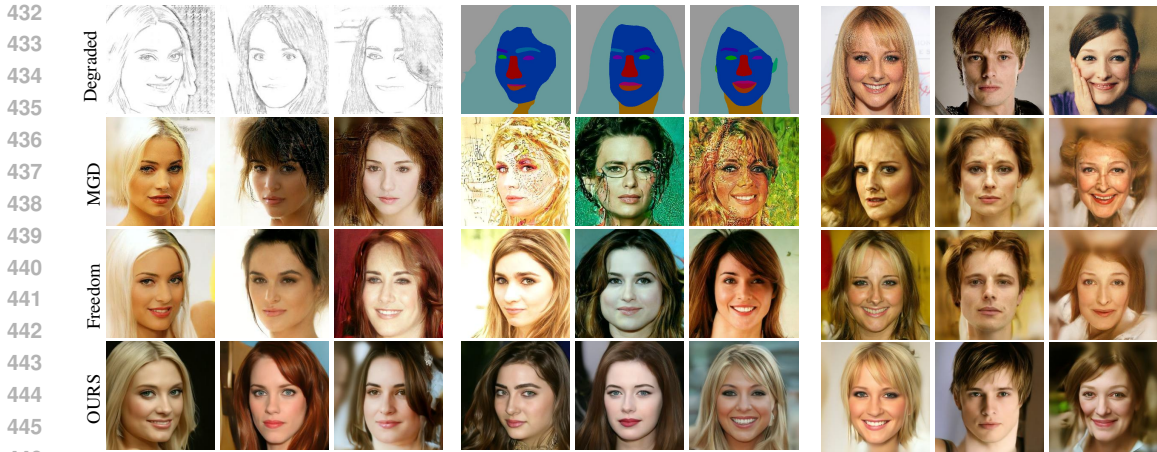


Figure 5: Qualitative comparisons for Non-linear Tasks on CelebA dataset for 100 inference steps

Method	Semantic Parsing			ID Guidance			Face Sketch		
	Distance↓	LPIPS↓	FID↓	Distance↓	LPIPS↓	FID↓	Distance↓	LPIPS↓	FID↓
Freedom <a href="#">Yu et al. (2023)</a>	1864.51	0.6030	66.89	0.3767	0.7058	81.40	39.05	0.6583	86.51
MGD <a href="#">He et al. (2023)</a>	<b>2698.27</b>	0.6995	104.32	0.4291	0.7178	92.61	39.34	0.6576	70.42
Ours	2722.51	<b>0.6199</b>	<b>79.42</b>	<b>0.3780</b>	<b>0.5932</b>	<b>82.70</b>	<b>39.03</b>	<b>0.5509</b>	<b>69.51</b>

Table 4: Non-linear tasks. Best results out of zeroth-order optimization algorithms are highlighted.

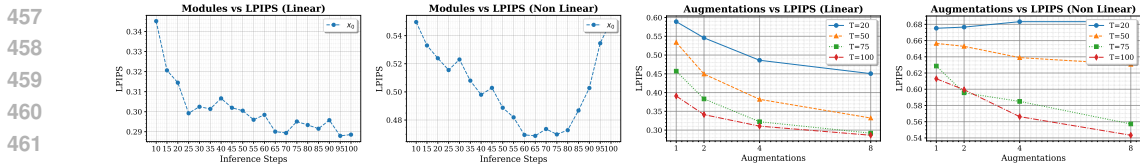


Figure 6: Ablation analysis on linear and non-linear tasks. FaceID guidance & ImageNet superresolution

### 5.1 EFFECT OF DIFFUSEAUGMENT

We notice that for linear tasks, even for low values of  $T$  such as  $T = 20$ , just by increasing the number of augmentations at the output to 8, the perceptual quality drastically improves, matching that of diffusion inference with  $T = 50$  with just 2 augmentations. Further, we notice that although the effect of augmentations is very significant for linear tasks, the performance is not that significant or rather drops in some cases for low  $T$  such as  $T = 20$ ; this is because with 20 diffusion steps, most intermediate MMSE estimates remain noisy, and hence the guidance network ArcFace [Deng et al. \(2019\)](#) cannot handle such input and hence returns irregular gradients affecting the quality. However, we can see that as  $T$  increases and when there are enough gradient steps, DiffuseAugment plays a significant role in boosting the performance.

### 5.2 EFFECT OF DIFFERENT COMPONENTS OF GUIDANCE

We present the ablation analysis of the effect of different terms of guidance in Figure 6. Please note that for this experiment, we set the number of augmentations from DiffuseAugment as 1. We also turn off time travel sampling for this experiment. For this experiment, we perform guidance with respect to  $\epsilon_\theta(t)$  until  $t_0$  and perform guidance with respect to  $\hat{x}_t$  for  $t > t_0$ . Here  $t = 100$  represents pure gaussian noise and  $t = 0$  represents the image. As we can see, guidance with  $\hat{x}_t$  alone faces a drop in performance initially for a low number of inference steps for non linear cases. We argue that this is because the guidance flow through the MMSE estimate is weak during the earlier steps of diffusion. Although time travel sampling helps to alleviate this issue, careful parameter tuning is

required to obtain satisfactory results. We also notice that guiding utilizing the gradients of the output noise of the network closer to the start of the generation process produces better results.

## 6 LIMITATIONS AND FUTURE WORKS

Although we illustrated the working across various tasks for pixel space diffusion models, the direct approach cannot be used for latent diffusion models for the task of linear inverse problems, and one might have to apply multiple steps of time travel sampling to fix this issue, making a large computational overhead of the overall sampling time. We emphasize that this problem arises due to the reconstruction error in the VAE that encodes the image to the latent space. In the future, we will attempt to improve upon this with better optimization techniques. Moreover, although the proposed empirical estimate based on distance over gradients works for most tasks and shows the existence of an optimal parameter estimate, a thorough mathematical evaluation and the most optimal parameters are still missing. We leave this problem up to future works to estimate the optimal guidance parameter.

## 7 CONCLUSION

In this paper, we proposed an improvement to existing loss-guided techniques for zero-shot conditional generation with an unconditional diffusion model. Specifically, we proposed a sampling technique that removes the need to backpropagate through the diffusion U-Net in order to tackle sampling for general inverse problems. We also present an empirical function for automatic scaling parameters that removes the need for manual scaling parameter tuning, which was previously a huge hurdle in using classifier-free guidance. The newly proposed scaling parameter also removes the need for model-specific tuning of start and end guidance steps. We also introduced a differentiable data augmentation method that significantly improves the sampling fidelity. We illustrated the working of our method across 4 linear and 3 non-linear tasks across faces and real image domains. Our sampling technique produces photorealistic samples with much lower sampling time and higher fidelity than existing methods.

## REFERENCES

- Hemant K Aggarwal, Merry P Mani, and Mathews Jacob. MoDL: Model-based deep learning architecture for inverse problems. *IEEE transactions on medical imaging*, 38(2):394–405, 2018.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. *arXiv preprint arXiv:2302.07121*, 2023.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22563–22575, 2023.
- Stanley H Chan, Xiran Wang, and Omar A Elgendy. Plug-and-play admm for image restoration: Fixed-point convergence and applications. *IEEE Transactions on Computational Imaging*, 3(1): 84–98, 2016.
- Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. *arXiv preprint arXiv:2108.02938*, 2021.

- 540 Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul  
541 Ye. Diffusion posterior sampling for general noisy inverse problems. In *International Confer-*  
542 *ence on Learning Representations*, 2023a. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=OnD9zGAGT0k)  
543 [OnD9zGAGT0k](https://openreview.net/forum?id=OnD9zGAGT0k).
- 544 Hyungjin Chung, Dohoon Ryu, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Solving  
545 3d inverse problems using pre-trained 2d diffusion models. *IEEE/CVF Conference on Computer*  
546 *Vision and Pattern Recognition*, 2023b.
- 547 Hyungjin Chung, Jong Chul Ye, Peyman Milanfar, and Mauricio Delbracio. Prompt-tuning latent  
548 diffusion models for inverse problems. *ArXiv*, abs/2310.01110, 2023c.
- 549 Aaron Defazio and Konstantin Mishchenko. Learning-rate-free learning by d-adaptation. *arXiv*  
550 *preprint arXiv:2301.07733*, 2023.
- 551 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
552 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
553 pp. 248–255. Ieee, 2009.
- 554 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin  
555 loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision*  
556 *and pattern recognition*, pp. 4690–4699, 2019.
- 557 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
558 *in Neural Information Processing Systems*, 34, 2021.
- 559 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
560 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the*  
561 *ACM*, 63(11):139–144, 2020.
- 562 Alexandros Graikos, Nikolay Malkin, Nebojsa Jojic, and Dimitris Samaras. Diffusion models as  
563 plug-and-play priors. *arXiv preprint arXiv:2206.09012*, 2022.
- 564 Yutong He, Naoki Murata, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Dongjun Kim, Wei-  
565 Hsiang Liao, Yuki Mitsufuji, J Zico Kolter, Ruslan Salakhutdinov, et al. Manifold preserving  
566 guided diffusion. *arXiv preprint arXiv:2311.16424*, 2023.
- 567 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
568 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 569 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P  
570 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition  
571 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- 572 MD Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. A comprehensive  
573 survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):1–36, 2019.
- 574 Maor Ivgi, Oliver Hinder, and Yair Carmon. Dog is sgd’s best friend: A parameter-free dynamic step  
575 size schedule. *arXiv preprint arXiv:2302.12022*, 2023.
- 576 Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint*  
577 *arXiv:2305.02463*, 2023.
- 578 Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for  
579 improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- 580 Bahjat Kawar, Gregory Vaksman, and Michael Elad. Stochastic image denoising by sampling from  
581 the posterior distribution. In *Proceedings of the IEEE/CVF International Conference on Computer*  
582 *Vision (ICCV) Workshops*, pp. 1866–1875, October 2021.
- 583 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration  
584 models. *arXiv preprint arXiv:2201.11793*, 2022.
- 585 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In  
586 *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

- 594 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool.  
595 Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the*  
596 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11461–11471, 2022.  
597
- 598 Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit:  
599 Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*,  
600 2021.
- 601 Nithin Gopalakrishnan Nair, Anoop Cherian, Suhas Lohit, Ye Wang, Toshiaki Koike-Akino, Vishal M  
602 Patel, and Tim K Marks. Steered diffusion: A generalized framework for plug-and-play conditional  
603 image synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
604 pp. 20850–20860, 2023.
- 605
- 606 Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. Plug & play  
607 generative networks: Conditional iterative generation of images in latent space. In *Proceedings of*  
608 *the IEEE conference on computer vision and pattern recognition*, pp. 4467–4477, 2017.
- 609 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
610 In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.  
611
- 612 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d  
613 diffusion. *arXiv*, 2022.
- 614 Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep  
615 convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.  
616
- 617 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
618 resolution image synthesis with latent diffusion models, 2021.
- 619
- 620 Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet,  
621 and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022*  
622 *Conference Proceedings*, pp. 1–10, 2022a.
- 623 Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi.  
624 Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and*  
625 *Machine Intelligence*, 2022b.
- 626
- 627 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image  
628 recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 629 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
630 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*,  
631 pp. 2256–2265. PMLR, 2015.  
632
- 633 Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion  
634 models for inverse problems. In *International Conference on Learning Representations*, 2022.
- 635 Jiaming Song, Qinsheng Zhang, Hongxu Yin, Morteza Mardani, Ming-Yu Liu, Jan Kautz, Yongxin  
636 Chen, and Arash Vahdat. Loss-guided diffusion models for plug-and-play controllable generation.  
637 In *International Conference on Machine Learning*, pp. 32483–32498. PMLR, 2023.  
638
- 639 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of  
640 score-based diffusion models. *Advances in Neural Information Processing Systems*, 34, 2021a.
- 641 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
642 Poole. Score-based generative modeling through stochastic differential equations. In *International*  
643 *Conference on Learning Representations*, 2021b. URL [https://openreview.net/forum?](https://openreview.net/forum?id=PXTIG12RRHS)  
644 [id=PXTIG12RRHS](https://openreview.net/forum?id=PXTIG12RRHS).
- 645
- 646 Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-  
647 resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of*  
*the IEEE conference on computer vision and pattern recognition*, pp. 8798–8807, 2018.

648 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion  
 649 null-space model. In *The Eleventh International Conference on Learning Representations, 2023*.  
 650 URL <https://openreview.net/forum?id=mRieQgMtNTQ>.  
 651

652 Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu,  
 653 Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion  
 654 models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on*  
 655 *Computer Vision*, pp. 7623–7633, 2023a.

656 Zike Wu, Pan Zhou, Kenji Kawaguchi, and Hanwang Zhang. Fast diffusion model. *arXiv preprint*  
 657 *arXiv:2306.06991*, 2023b.

659 Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-free  
 660 energy-guided conditional diffusion model. *arXiv preprint arXiv:2303.09833*, 2023.

661 Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for  
 662 data-efficient gan training. *Advances in neural information processing systems*, 33:7559–7570,  
 663 2020.  
 664

## 665 A APPENDIX

### 666 B ALGORITHM OF DREAMGUIDER

667 We present the over algorithm of dreamguider without time travel sampling and the parameter  
 668 estimation algorithm in Algorithm 1  
 669

### 670 C PROOF FOR PERTURBED MARKOVIAN KERNEL EQUATION

671 In the main paper, we emphasized that any positive distance function can be utilized for performing  
 672 conditional generation using the perturbed Markovian kernel equation. Here we proceed to derive  
 673 the perturbed transition step. For the proof we closely follow the work from Dickenson et al [Sohl-  
 674 Dickstein et al. \(2015\)](#). Given a unconditional transition distribution  $p_\theta(x_{t-1}|x_t)$  and a distance  
 675 function  $r(\cdot, y)$ , where  $y$  is the condition provided Please note that we assume  $r(\cdot, y)$  has relatively  
 676 small variance compared to  $p_\theta(x_{t-1}|x_t)$ , We know that at equilibrium state, the distribution at any  
 677 timestep  $t$  ina diffusion model can be written as  
 678

$$679 p(x_{t-1}) = \int p(x_t)p_\theta(x_{t-1}|x_t)dx_t. \quad (16)$$

680 To estimate a perturbed transition kernel  $\hat{p}_\theta(x_{t-1}|x_t)$ , we start the perturbed distribution as  
 681

$$682 p(x_{t-1})r(x_{t-1}, y) = \int r(x_t, y)p(x_t)\hat{p}_\theta(x_{t-1}|x_t)dx_t. \quad (17)$$

683 By simple algebraic manipulations, taking  $r(x_{t-1}, y)$  to the other side, we get  
 684

$$685 p(x_{t-1}) = \int \frac{r(x_t, y)}{r(x_{t-1}, y)}p(x_t)\hat{p}_\theta(x_{t-1}|x_t)dx_t. \quad (18)$$

686 By comparing Equation (16) and Equation (18) we can see that one solution for the transitional  
 687 distribution is  
 688

$$689 \hat{p}_\theta(x_{t-1}|x_t) = p_\theta(x_{t-1}|x_t)\frac{r(x_{t-1}, y)}{r(x_t, y)}. \quad (19)$$

690 Also since normalization constants doesn't affect the score function or transition step, Absorbing  $x_t$   
 691 to the normalization factor of  $p_\theta(x_{t-1}|x_t)$ , another valid perturbed transition kernel is

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

---

**Algorithm 1** Dreamguider

---

**Input:** distance function  $r(\cdot, y)$ , condition  $y$ , Timesteps  $T$

```

1:  $x_T \sim \mathcal{N}(x_T; 0, I)$ 
2: for  $t = T - 1, \dots, 1$  do
3:    $\Sigma = \sqrt{I - \bar{\alpha}_t}$ 
4:    $\epsilon \sim \mathcal{N}(\epsilon; 0, I)$ 
5:    $\hat{x}_t = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t)}{\sqrt{\bar{\alpha}_t}}$ 
6:   Compute  $\frac{dr(\hat{x}_t, y)}{d\hat{x}_t}, \frac{dr(\hat{x}_t, y)}{d\epsilon_\theta(x_t)}$ 
7:   update  $c = ESTIMATE(t, \epsilon_\theta(x_t), \frac{dr(\hat{x}_t, y)}{d\epsilon_\theta(x_t)})$ 
8:   update  $d = ESTIMATE(t, \hat{x}_t, \frac{dr(\hat{x}_t, y)}{d\hat{x}_t})$ 
9:    $c_t = c\sqrt{\alpha_{t-1}}$ 
10:   $d_t = -d \cdot \frac{1 - \alpha_t}{\sqrt{\alpha_t}\sqrt{1 - \alpha_t}}$ 
11:  if  $t < t_0$  then
12:     $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t) \right) + \sigma_t \epsilon - d_t \Sigma \frac{dr(\hat{x}_t, y)}{d\epsilon_\theta(x_t)}$ 
13:  else
14:     $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t) \right) + \sigma_t \epsilon - c_t \Sigma \frac{dr(\hat{x}_t, y)}{d\hat{x}_t} -$ 
15:  end if
16: end for
17: function ESTIMATE( $t, f_i, g_t$ )
18:   if  $t = T$  then
19:      $\gamma_t = \frac{1e^{-5}}{\sqrt{g_T^2}}$ 
20:     Store  $f_T$ ,
21:   else
22:      $\gamma_t = \frac{\max_{i>t} |f_i - f_T|}{\sqrt{\sum_{i=i}^T g_i^2}}$ 
23:   end if
24:   Store  $\sqrt{\sum_{i=i}^T g_i^2}$ 
25:   return  $\gamma_t$ 
26: end function return  $x_0$ 

```

---

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

$$\hat{p}_\theta(x_{t-1}|x_t) = p_\theta(x_{t-1}|x_t) \frac{r(x_{t-1}, y)}{Z}. \quad (20)$$

Please note that the term  $Z$  does not affect the transition step in the reverse process when the variance of  $r(\cdot, y)$  is small.

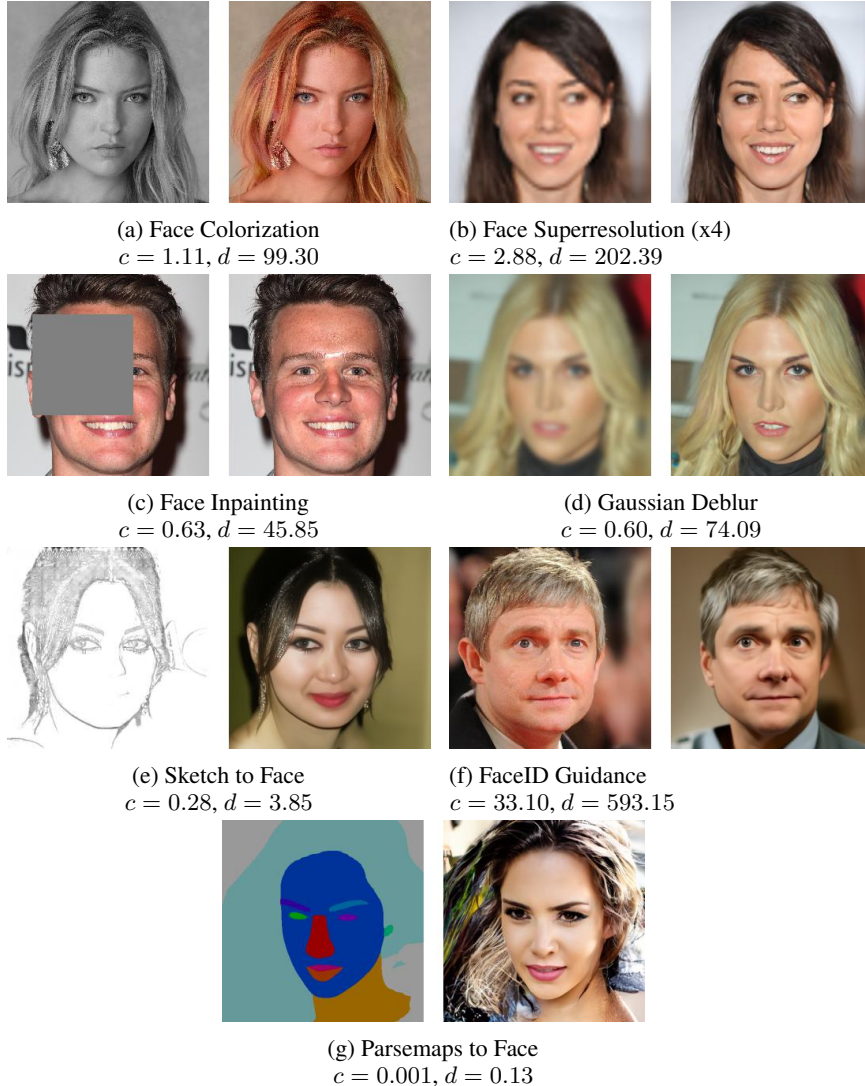


Figure 7: Figure illustrating the guidance scales for different tasks.

Method	Freedom	Dreamguider(1)	Dreamguider(2)	Dreamguider(3)
Sketch to Face	24.95	17.55	27.04	35.09
FaceID to Face	24.94	20.45	31.89	41.80
FaceParse to Face	56.25	48.35	75.43	107.02

Table 5: Non-linear tasks ablation analysis on time taken, the value is represented in seconds

## D TIME COMPARISON FOR DREAMGUIDER WITH TIMETRAVEL SAMPLING AND FREEDOM(FIRST ORDER) FOR NON LINEAR TASKS

We present the time taken by Freedom, a first order algorithm for one step of time travel sampling [Lugmayr et al. \(2022\)](#); [Yu et al. \(2023\)](#) in Table 5

810 E ESTIMATED PARAMETER VALUE FOR DIFFERENT TASKS

811

812

813

814

815

816

In this section, we present the result and the parameter estimated by our approach for different tasks. For this experiment, we use 100 timesteps of diffusion and present the value at the 100th timestep. Here we define  $d$  as the scaling factor of the scaling constant of the the loss derivative relative to  $\epsilon_\theta(x_t)$  and  $c$  as that of  $\hat{x}_t$  as in the main paper . The corresponding results are shown in Figure 7

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

860

861

862

863

F NON CHERRY PICKED RESULTS FOR DIFFERENT TASKS.



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917



Figure 8: Figure illustrating **Non cherry picked** results for ImageNet colorization

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

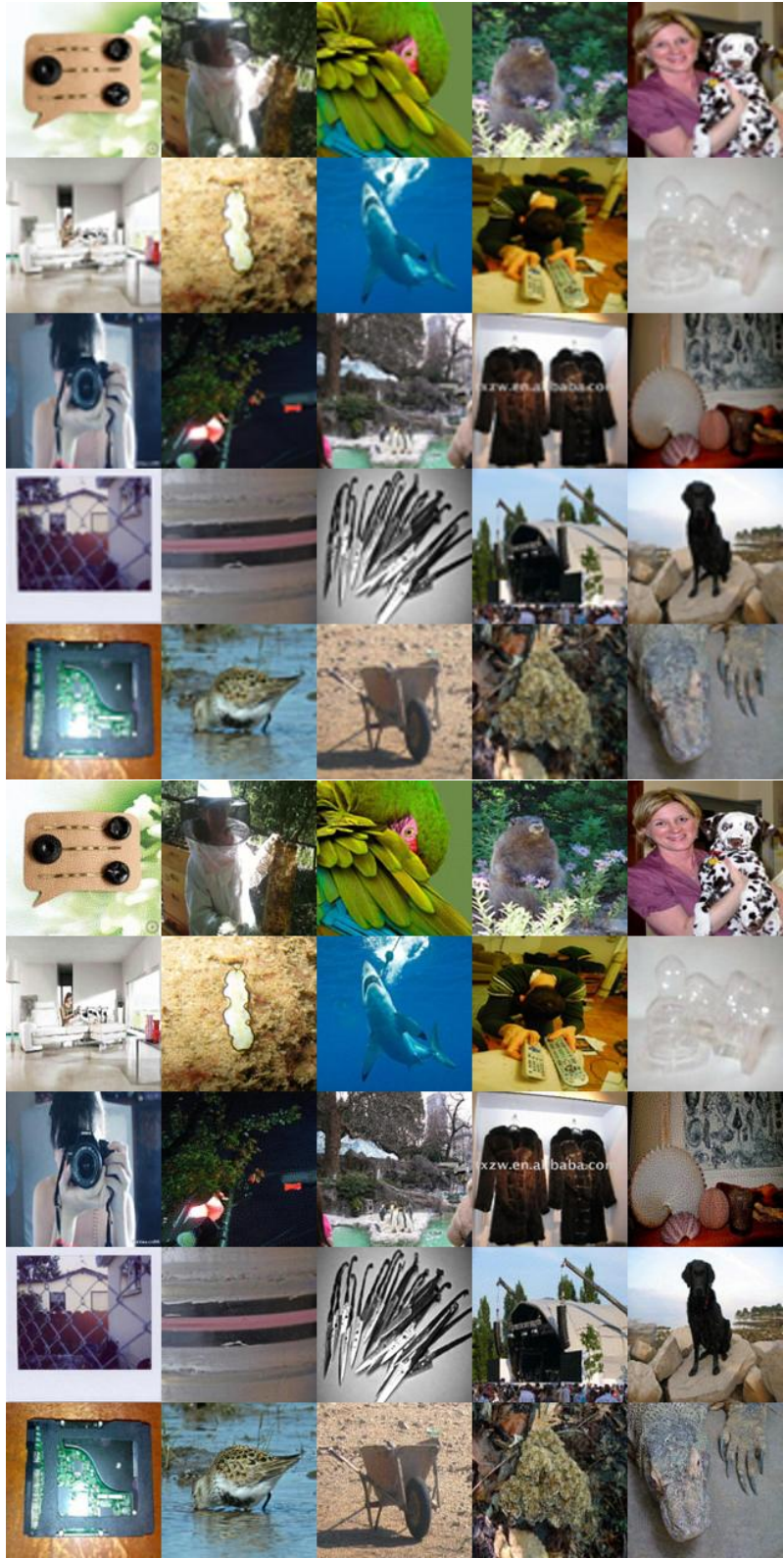


Figure 9: Figure illustrating **Non cherry picked** results for ImageNet superresolution

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

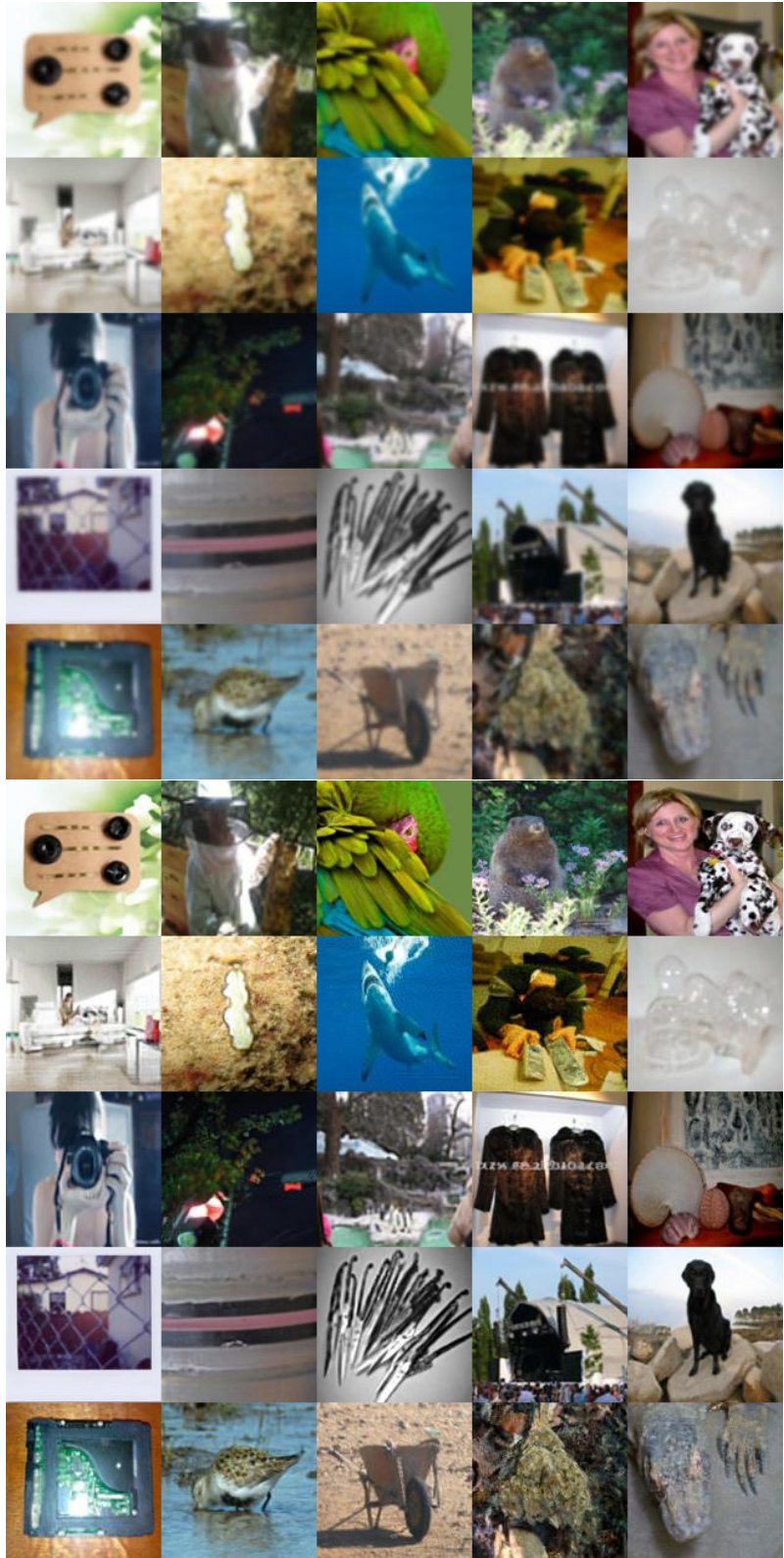


Figure 10: Figure illustrating **Non cherry picked** results for Gaussian deblurring on ImageNet

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079



Figure 11: Figure illustrating **Non cherry picked** results for face colorization

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133



Figure 12: Figure illustrating **Non cherry picked** results for face superresolution

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187



Figure 13: Figure illustrating **Non cherry picked** results for Gaussian Deblurring

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241



Figure 14: Figure illustrating **Non cherry picked** results for face inpainting

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

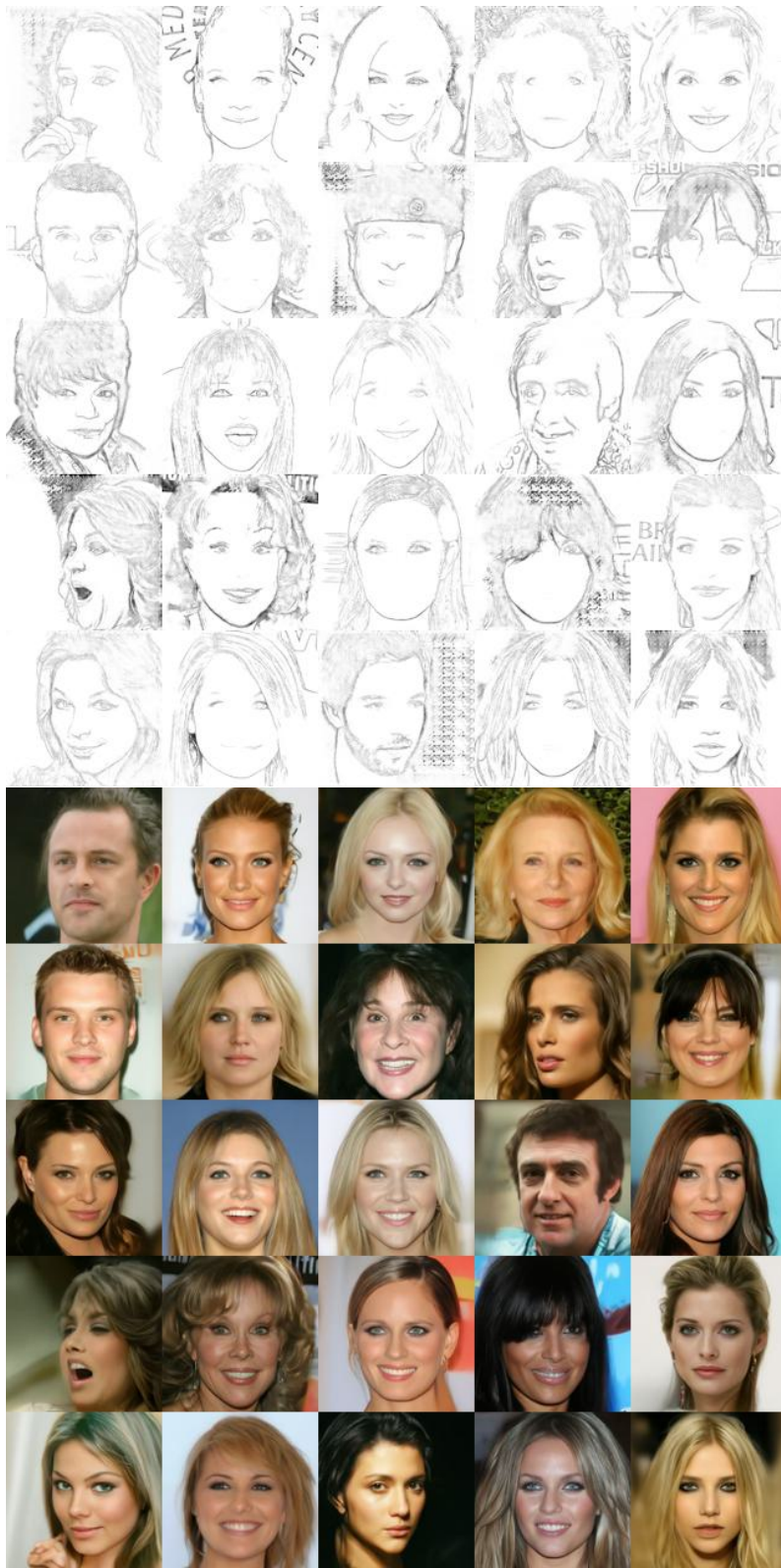


Figure 15: Figure illustrating **Non cherry picked** results for sketch to face synthesis



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

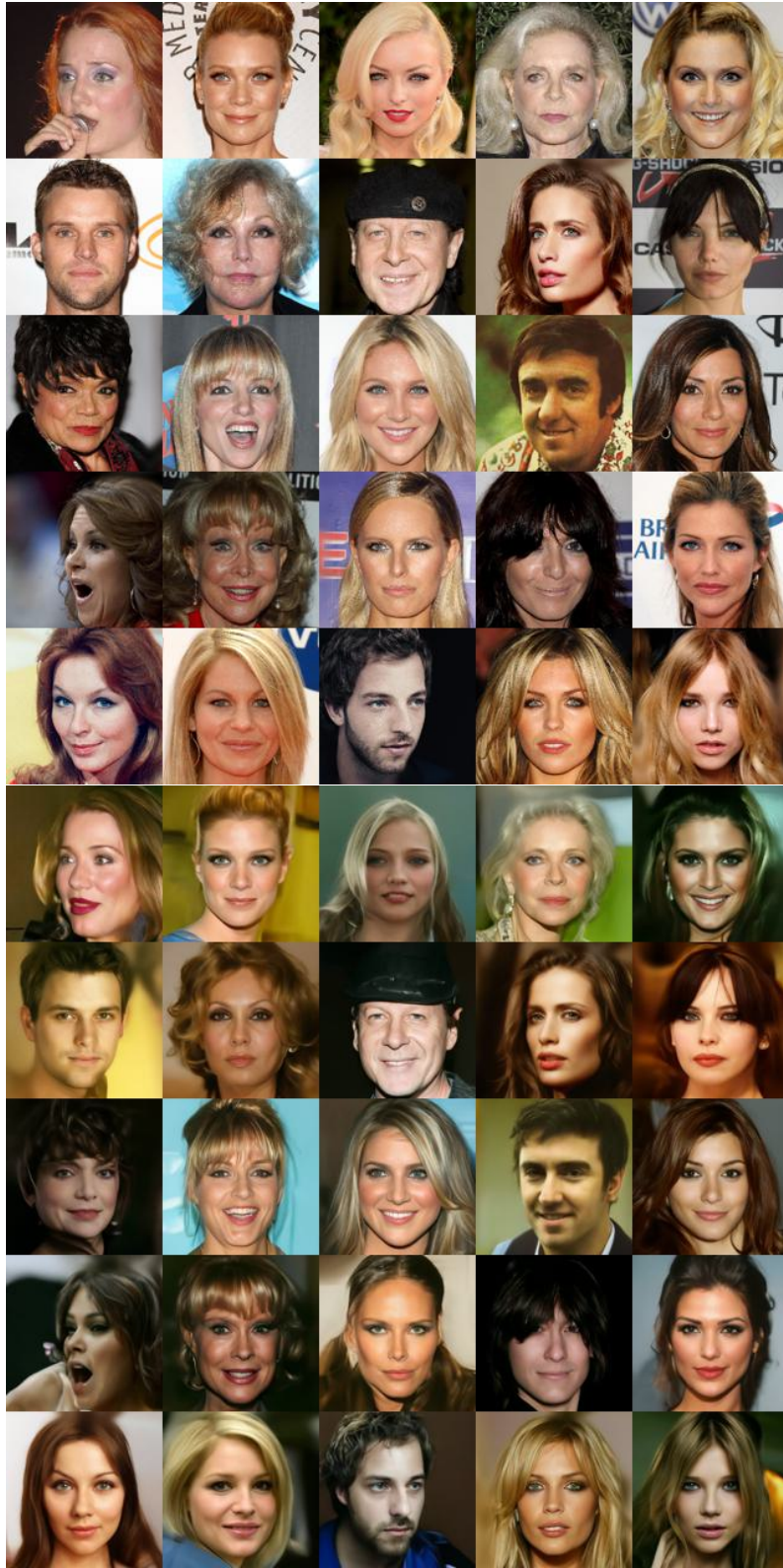


Figure 16: Figure illustrating **Non cherry picked** results for Face ID guidance

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

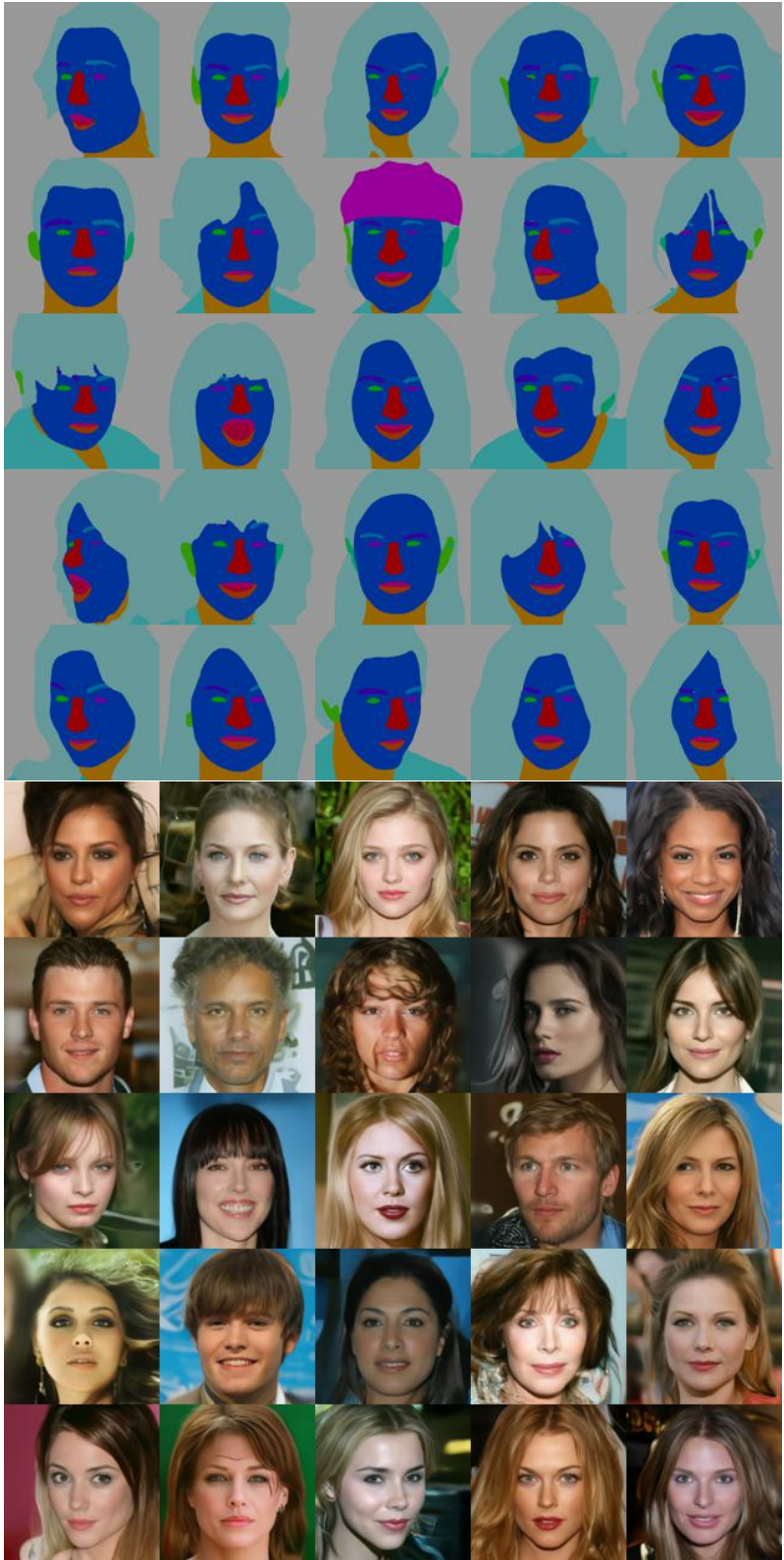


Figure 17: Figure illustrating **Non cherry picked** results for Face Parse Guidance