

---

# Synthetic Shifts to Initial Seed Vector Exposes the Brittle Nature of Latent-Based Diffusion Models

---

Mao Po-Yuan\* Shashank Kotyan\* Tham Yik Foong Danilo Vasconcellos Vargas  
Laboratory of Intelligent Systems  
Kyushu University, Fukuoka, Japan

## Abstract

Recent advances in Conditional Diffusion Models have led to substantial capabilities in various domains. However, understanding the impact of variations in the initial seed vector remains an underexplored area of concern. Particularly, latent-based diffusion models display inconsistencies in image generation under standard conditions when initialized with suboptimal initial seed vectors. To understand the impact of the initial seed vector on generated samples, we propose a reliability evaluation framework that evaluates the generated samples of a diffusion model when the initial seed vector is subjected to various synthetic shifts. Our results indicate that slight manipulations to the initial seed vector of the state-of-the-art Stable Diffusion (Rombach et al., 2022) can lead to significant disturbances in the generated samples, consequently creating images without the effect of conditioning variables. In contrast, GLIDE (Nichol et al., 2022) stands out in generating reliable samples even when the initial seed vector is transformed. Thus, our study sheds light on the importance of the selection and the impact of the initial seed vector in the latent-based diffusion model.

## 1 Introduction

In recent years, diffusion models have risen to the forefront as state-of-the-art instruments for content creation and the precision generation of high-quality synthetic data empowered by deep neural networks and extensive datasets (Bao et al., 2022). Their influence spans across multiple domains, including images (Ho et al., 2020; Dhariwal & Nichol, 2021; Ho et al., 2022; Ho & Salimans, 2022), audio (Kong et al., 2021; Huang et al., 2022a,b; Kim et al., 2022), texts (Li et al., 2022), molecules (Xu et al., 2022), solidifying their status as leading technologies in data synthesis. Notably, the release of Stable Diffusion (Rombach et al., 2022), an advanced open-source text-based image generation model, has sparked diverse applications and workloads (Lugmayr et al., 2022; Jeong et al., 2023).

Samuel et al. (2023b) links the quality of generating a rare concept to the initial seed vector, suggesting the importance of choosing the initial seed vector. Using this understanding, researchers want to drive creativity and facilitate high-quality synthetic data generation with the help of manipulation the initial seed vector.

In the Figure 1, we show that the latent-based diffusion model is more susceptible to generating non-reliable output based on shifts to the initial seed vector, highlighting the importance of selection of the initial seed vector in the latent-based diffusion model. Since a comprehensive understanding of the diffusion process to generate samples is still lacking (Li et al., 2023b,a), we also focus on understanding the intricate landscape of the generated samples by the diffusion models. To comprehensively evaluate the impact of the initial seed vector in generating samples, we meticulously investigate the ability of a diffusion model to adapt from shifts to the initial seed vector. These

---

\*These authors contributed equally to this work

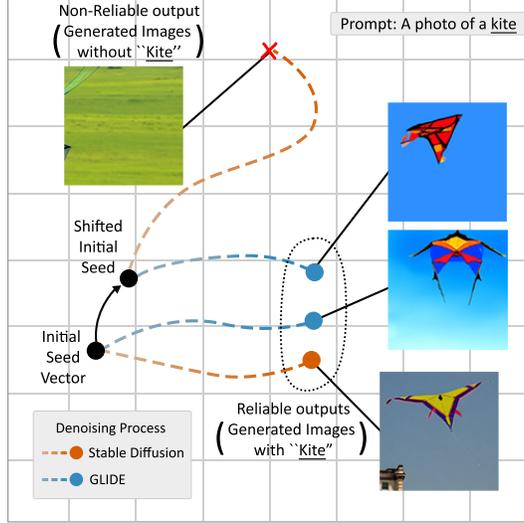


Figure 1: **Illustration of performance of Stable Diffusion (Rombach et al., 2022) compared to GLIDE (Nichol et al., 2022) in the Presence of Synthetic Shifts to the Initial Seed Vector.** This figure illustrates the trajectories of the diffusion process in pixel space for both Stable Diffusion v2.1 and GLIDE models when subjected to seed vector shifts. We notice that with a slight shift to the initial seed vector, the generated image by Stable Diffusion diverges towards a non-reliable output. On the other hand, GLIDE consistently generates reliable outputs, demonstrating robustness to shift in the initial seed vector.

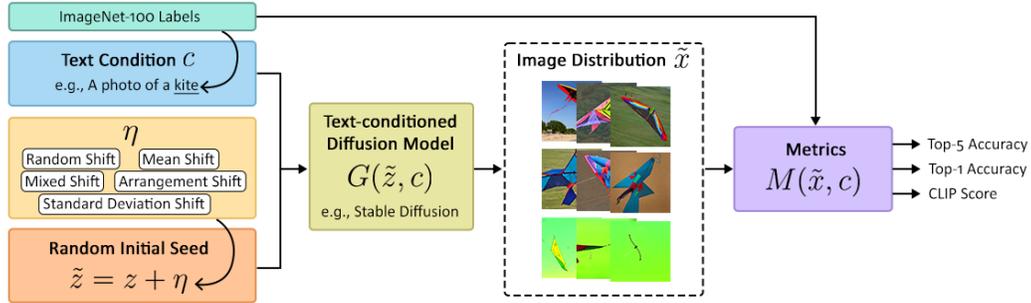


Figure 2: **Illustration of the Reliability Evaluation Framework.** Prompts are constructed by filling the sentence "A photo of a \_\_\_\_" with ImageNet100 (SHEKHAR, 2021) labels, and the random initial seed vector  $z$  is transformed by  $\eta$  to form a shifted initial seed vector  $\tilde{z}$ . Subsequently, the text-conditioned diffusion model generates an image distribution  $\tilde{x}$  based on these inputs. The metrics  $M$  then evaluate the generated image distribution with the known label used in the conditioning variable.

synthetic shifts encompass a spectrum, including Random Shift, Mean Shift, Standard Deviation Shift, Mixed Shift, and Arrangement Shift. Thus, to assess the generated samples and evaluate their correctness, we propose a model-agnostic reliability evaluation framework, as illustrated in Figure 2.

### Contributions:

**Reliability Evaluation Framework:** We propose a simple framework (Figure 2) for systematically evaluating how diffusion-based models can handle the shift to the initial seed vector.

**Brittle nature of Stable Diffusion (Rombach et al., 2022):** Our results show that slight variations to the initial random vector break the Stable Diffusion as it creates undesired samples.

**Robustness of GLIDE (Nichol et al., 2022):** Through our experiments, we provide empirical evidence substantiating the reliable nature of samples generated by GLIDE compared to Stable Diffusion.

## 2 Background and Related Works

The diffusion model (Sohl-Dickstein et al., 2015) is a latent variable model that can be described as a Markov chain with learned Gaussian transitions. It consists of two main components: the diffusion process and the reverse process. The reverse process is a trainable model that is trained to reduce the Gaussian noise introduced by the diffusion process systematically.

To illustrate, consider input data represented as  $x \in \mathbb{R}$ , the approximate posterior ( $q$ ) is expressed by:

$$q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1 - \beta_t} \cdot x_{t-1}, \beta_t I) \quad (1)$$

which is defined as a fixed Markov chain. This Markov chain progressively introduces Gaussian noise to the data in accordance with a predefined schedule of variances, denoted as  $\beta_1, \beta_2, \dots, \beta_T$ :

$$q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t-1}). \quad (2)$$

Subsequently, the reverse process with trainable parameters  $p_\theta(x_{0:T})$  revert the diffusion process returning the data distribution:

$$p_\theta(x_{0:t}) := p(x_T) \cdot \prod_{t=1}^T p_\theta(x_{t-1}|x_t), \quad (3)$$

$$p_\theta(x_{t-1}|x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)). \quad (4)$$

where  $p_\theta$  contains the mean  $\mu_\theta(x_t, t)$  and the variance  $\Sigma_\theta(x_t, t)$ , both of them are trainable models predict the value by using the current time step and the current noise.

By fixing the forward process variances, Denoising Diffusion Probabilistic Models (DDPM) (Ho et al., 2020) modify the Equation 4 to :

$$p_\theta(x_{t-1}|x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma^2 I). \quad (5)$$

This smart design achieved higher-quality image synthesis than Generative Adversarial Networks (GANs) (Goodfellow et al., 2014).

**Diffusion Models and its variants:** Similar to other types of generative models (Mirza & Osindero, 2014; Sohn et al., 2015), the generation process can also be conditioned. For instance, GLIDE (Nichol et al., 2022) learns to generate images according to an input textual sentence on the image space, DALL-E-2 (Ramesh et al., 2022) uses a DDPM to learn a prior distribution on the CLIP (Radford et al., 2021). Text-to-image generation is also explored in Stable Diffusion (Rombach et al., 2022) and Imagen (Saharia et al., 2022). Furthermore, the release of Stable Diffusion has catalyzed a surge in diverse applications and workloads. However, the recursive sampling process makes the diffusion model a time-consuming model. To address this problem, Song et al. (2020) proposed Denoising Diffusion Implicit Models (DDIM) (Song et al., 2020), a non-Markovian inference process that faster the sampling process. Salimans & Ho (2022) propose to distill the prediction network into new networks, which progressively reduce the number of sampling steps. Rombach et al. (2022) speed up sampling by splitting the process into a compression stage and a generation stage and applying the DDPM on the compressed (latent) space.

**Diffusion Model lack reliable explanations:** Other than improving the application site, a sort of research notice that Diffusion Model lacks reliable explanations (Ning et al., 2023; Li et al., 2023a; Daras et al., 2023). This problem is called exposure bias or sampling drift. The researchers claim that error propagation happens to diffusion models because the models are of a cascade structure. Li et al. (2023b) further develop a theoretical framework for analyzing the error propagation of diffusion models.

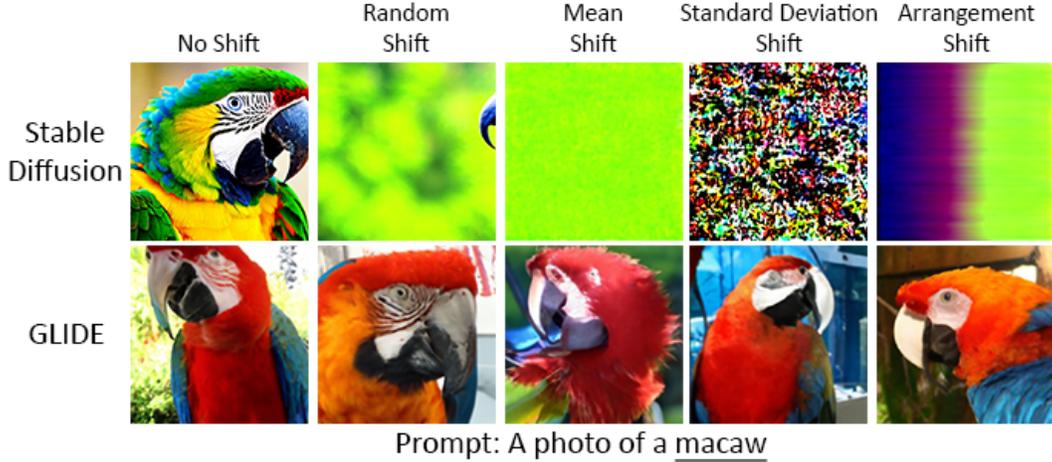


Figure 3: **Examples showcasing the impact of various shifts to initial seed vector on image generation of Stable Diffusion (Top) and GLIDE (Bottom).** We do the following shifts to the initial seed vector (Left to Right): a) No Shift, b) Random Shift ( $\eta_r = 0.2$ ), c) Mean Shift ( $\eta_m = 0.2$ ), d) Standard Deviation Shift ( $\eta_s = 0.3$ ), and e) Arrangement Shift ( $\eta_a = 64$ ).

**Impact of Seed Vector to Diffusion Models:** Recently, latent space and seed vectors have been shown to be highly correlated to the final result (Samuel et al., 2023a; Wu & De la Torre, 2022; Ge et al., 2023). Better seed vectors can consistently generate high-quality images, and even facilitate conditional control to achieve desired results (Mao et al., 2023; Singh et al., 2022b). This impact is compounded when generating the rare distribution such as rare fine-grained concepts or rare combinations (Liu et al., 2022; Zhao et al., 2019; Feng et al., 2022; Chefer et al., 2023). Samuel et al. (2023b) shows that finely selected seed vectors can generate rare distribution, which raises our interest in understanding how seed vectors behave differently.

### 3 Evaluating Reliability of Diffusion Models

#### 3.1 Problem Definition

Consider a condition-based diffusion model  $G(\cdot)$  in which we generate samples  $x$  using the equation  $x = G(z, c)$ . Here,  $z$  is the initial seed vector following a simple and tractable normal distribution,  $z \sim \mathcal{N}(\mu, \alpha^2)$ , as described in prior work (Sohl-Dickstein et al., 2015) and  $c$  is the conditioning variable to which  $x$  is closely related. This conditioning variable can represent various data types, such as images, sentences, or sounds. The strength of the correlation between  $x$  and  $c$  can be quantified as  $p(x|c) \propto M(x, c)$ . Here  $M$  is a function that evaluates the variable  $x$  to the conditioning variable  $c$ .

Further, the susceptibility of the diffusion model  $G(\cdot)$  to the initial seed vector  $z$  can be evaluated by manipulating  $z$ . To introduce shifts to the initial seed vector  $z$ , we transform the initial seed vector  $z$  with  $\eta$ , resulting in  $\tilde{z} = z + \eta$ . Here,  $\tilde{z}$  is the modified seed vector. Further, the sample generated by the modified seed vector  $\tilde{z}$  can be represented as  $\tilde{x} = G(\tilde{z}, c)$ . Particularly, we are interested in identifying instances where there exists a  $\tilde{z}$  such that  $M(\tilde{x}, c)$  is significantly lesser than  $M(x, c)$ . This case implies that the modified seed vector  $\tilde{z}$  breaks the diffusion model’s  $G(\cdot)$  relationship between generated sample  $\tilde{x}$  and conditioning variable  $c$ , which was evident prior to the shift to initial seed vector  $z$  with the generated sample  $x$ .

#### 3.2 Synthetic Shifts to the Initial Seed Vector

In our investigation, we delve into the influence of the parameter  $\eta$  on the generated samples through distinct transformations, comprehensively evaluating the impact of modifications to the initial seed vector. In order to align with the original definition from the diffusion model (Sohl-Dickstein et al., 2015), we ensure that our modified initial vector maintains the normal distribution-like characteristics

Table 1: **Performance of Stable Diffusion v2.1 for all shifts proposed in Section 3.2.** We evaluate top-1 accuracy and top-5 accuracy calculated using pre-trained ViT and CLIP Score calculated using the OpenCLIP model. We observe that shifts to the initial seed vector degrade the reliability of generated samples by diffusion models as the shifts increase.

	← (Negative Shift)				(No Shift)	→ (Positive Shift)					
	<b>Random Shift (<math>\eta_r</math>)</b>										
	-0.30	-0.20	-0.15	-0.10	-0.05	0.00	0.05	0.10	0.15	0.20	0.30
Top-1 Accuracy (↑)	41.8%	65.6%	69.4%	71.1%	71.2%	71.6%	<u>72.9%</u>	72.1%	70.1%	65.5%	33.1%
Top-5 Accuracy (↑)	65.1%	87.5%	89.5%	89.8%	89.8%	90.0%	<u>90.7%</u>	90.6%	89.6%	86.6%	52.4%
CLIP Score (↑)	27.7	31.3	32.2	<u>32.6</u>	<u>32.6</u>	32.5	32.4	32.2	31.8	31.1	26.7
	<b>Mean Shift (<math>\eta_m</math>)</b>										
		-0.20	-0.15	-0.10	-0.05	0.00	0.05	0.10	0.15	0.20	
Top-1 Accuracy (↑)		12.9%	42.0%	65.4%	71.0%	71.6%	<u>72.2%</u>	65.5%	33.2%	6.4%	
Top-5 Accuracy (↑)		24.6%	65.8%	87.8%	89.8%	90.0%	<u>90.5%</u>	86.6%	51.8%	12.3%	
CLIP Score (↑)		22.5	27.6	31.2	<u>32.6</u>	32.5	32.1	31.1	26.6	21.7	
	<b>Standard Deviation Shift (<math>\eta_s</math>)</b>										
			-0.30	-0.20	-0.10	0.00	0.10	0.20	0.30		
Top-1 Accuracy (↑)			5.0%	41.0%	67.5%	<u>71.6%</u>	62.9%	7.8%	0.0%		
Top-5 Accuracy (↑)			7.7%	60.7%	87.9%	<u>90.0%</u>	85.6%	16.7%	0.0%		
CLIP Score (↑)			21.5	28.0	31.8	<u>32.5</u>	31.6	23.6	20.4		
	<b>Mixed Shift (<math>\eta_s, \eta_m</math>)</b>										
			(-0.30, -0.15)	(-0.20, -0.10)	(-0.10, -0.05)	(0.00, 0.00)	(0.10, 0.05)	(0.20, 0.10)	(0.30, 0.15)		
Top-1 Accuracy (↑)			2.1%	23.6%	66.3%	<u>71.6%</u>	65.5%	7.7%	0.0%		
Top-5 Accuracy (↑)			3.6%	38.2%	87.6%	<u>90.0%</u>	87.1%	17.1%	0.0%		
CLIP Score (↑)			17.9	23.2	31.3	<u>32.5</u>	31.6	24.1	22.2		
	<b>Arrangement Shift (<math>\eta_a</math>)</b>										
						0	8	16	32	64	
Top-1 Accuracy (↑)						<u>71.6%</u>	71.4%	56.3%	1.4%	0.0%	
Top-5 Accuracy (↑)						<u>90.0%</u>	90.0%	78.3%	2.7%	0.0%	
CLIP Score (↑)						<u>32.5</u>	32.2	30.8	22.1	18.2	

after our transformations. We evaluate the following five transformations: a) Random Shift, b) Mean Shift, c) Standard Deviation Shift, d) Mixed Shift, e) Arrangement Shift described in brief below. Figure 3 shows the generated image by Stable Diffusion (Rombach et al., 2022) and GLIDE (Nichol et al., 2022), when subjected to these transformations.

**Random Shift** offers a modification method that introduces randomness in the latent values. It incorporates randomness by sampling from a uniform distribution within the range of 0 to 1 and multiplying it by a scale factor  $\eta_r$ .

$$\tilde{z}_{\text{Random}} = z + \eta_r \cdot \mathcal{U}[0, 1] \quad (6)$$

**Mean Shift** represents the straightforward mean adjustment by adding a constant value  $\eta_m$  to all pixels.

$$\tilde{z}_{\text{Mean}} = z + \eta_m \quad (7)$$

**Standard Deviation Shift** modifies the variance of the distribution  $z$  with a scale factor  $\eta_s$ .

$$\tilde{z}_{\text{StandardDeviation}} = (1 + \eta_s) \cdot z \quad (8)$$

**Mixed Shift** combines both mean and standard deviation shifts, providing insights into their combined influence.

$$\tilde{z}_{\text{Mix}} = (1 + \eta_s) \cdot z + \eta_m \quad (9)$$

**Arrangement Shift** differs from the above methods by preserving the ideal normal distribution of mean  $\mu$  and standard deviation of  $\sigma^2$  while locally disrupting the normality. It achieves this by rearranging latent values without directly altering their values.

$$\tilde{z}_{\text{Arrangement}} = T(z, \eta_a) \quad (10)$$

Here,  $T$  is a sorting function that organizes the upper-left  $\eta_a$  submatrix. For instance, for Stable Diffusion the latent vector  $z \in \mathbb{R}^{64 \times 64 \times 4}$  and when we apply Arrangement Shift with  $\eta_a = 2$ , the latent values of upper-left  $2 \times 2$  submatrix along all 4 dimensions, i.e., a total of 16 latent values are sorted.

Note that in this study, we specifically selected shifted distributions with more than 86% overlap with the ideal normal distribution. This criterion ensures that we avoid investigating trivial cases, such as those resulting from scale issues causing the model to fail. Collectively, these transformations offer a comprehensive understanding of the impact of a slight shift in the initial seed vector in the resulting generated samples and the adaptivity of different diffusion models to handle such cases.

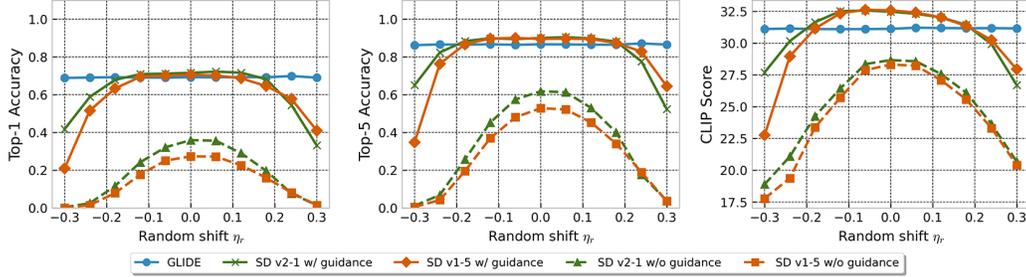


Figure 4: **Performance of different variants of Stable Diffusion and GLIDE for Random Shift.** We evaluate Stable Diffusion v1.5 with and without guidance, Stable Diffusion v2.1 with and without guidance and GLIDE. Note that GLIDE shows remarkable consistency in generating samples, whereas all variants Stable Diffusion shows degradation, showing that the shifts to the initial seed vector break the correlation of the generated sample and the conditioning prompt.

### 3.3 Reliability Evaluation Framework

To assess the reliability and correctness of generated with shifts to the initial seed vector, we propose a reliability evaluation framework as illustrated in Figure 2. We employ a simple strategy to include the label from a dataset as a conditioning variable to the diffusion model defined as Conditioning Variable  $c = \text{“A photo of a \_\_\_\_”}$  to evaluate different varieties of objects generated by diffusion models. In this study, the  $\_\_\_\_$  is dynamically filled from the ImageNet-100 dataset (SHEKHAR, 2021). For example, “A photo of a macaw”, “A photo of a sea lion”, “A photo of a crane” etc. Before initiating the reverse diffusion process to generate the samples with the conditioning, the initial seed vector undergoes transformation through the shifts described in Section 3.2. Subsequently, various metrics are employed to evaluate the performance of the generated samples of diffusion models under each shift scale.

## 4 Experimental Results

**Common Experimental Setup.** In all our experiments, we generate 100 images for each class of ImageNet-100; thus, a total of 10,000 images are evaluated across 100 classes (SHEKHAR, 2021). We compute top-1 accuracy and top-5 accuracy of the generated images using the state-of-the-art ViT-H/14 model pre-trained by SWAG (Singh et al., 2022a) for image classification task. This model has an impressive top-1 accuracy of 88.55% and top-5 accuracy of 98.69% on the ImageNet-1k dataset. We also compute the CLIP score using the OpenAI-CLIP model (Radford et al., 2021).

### 4.1 Evaluating Reliability Scores across Different Synthetic Shifts

**Experimental setup.** In this experiment, we study the robustness of the Stable Diffusion v2.1 model against shifts to the initial seed vector by applying the five shifting techniques described in Section 3.2. The diffusion model setting follows the original paper (Rombach et al., 2022) with 50 sampling time steps and a 7.5 classifier-free guidance scale. We vary the shift factor for Random Shift  $\eta_r$  and Standard Deviation Shift  $\eta_s$  within the range  $[-0.3, 0.3]$ . For Mean Shift we vary the shift factor  $\eta_m$  within the range  $[-0.2, 0.2]$ . In Mixed Shift, the effects of  $\eta_m$  were examined over  $[-0.15, 0.15]$  with 0.05 intervals, and  $\eta_s$  over  $[-0.3, 0.3]$  with 0.1 intervals. For the Arrangement Shift, we choose  $\eta_a$  as 8, 16, 32, and 64.

**Result.** Table 1 shows the top-1 accuracy, top-5 accuracy, and CLIP Score of generated images across all types of proposed shifts and multiple scales for Stable Diffusion v2.1 model. Despite the pre-trained model with an 88% top-1 accuracy and a 98% top-5 accuracy, the situation is quite different for the Stable Diffusion v2.1 model. It struggles to achieve only a 71% top-1 accuracy and a 90% top-5 accuracy. At the same time, we notice a drop in performance for both positive and negative shifts, and interestingly, the rate of decline is twice as fast for Mean Shift compared to a Random Shift. Further, we observe that introducing a positive standard deviation shift tends to deteriorate performance even more rapidly than the negative counterpart.

In both Mean Shift and Random Shift, the Stable Diffusion v2.1 demonstrates optimal performance in top-1 and top-5 accuracy (improvement of around 2%) with a slight positive shift ( $\eta = 0.05$ ). This observation is further supported by comparing mixed shift results with standard deviation shift outcomes. However, there is a contrasting trend in CLIP scores where the model achieves the best CLIP score with a slight negative shift ( $\eta = 0.05$ ).

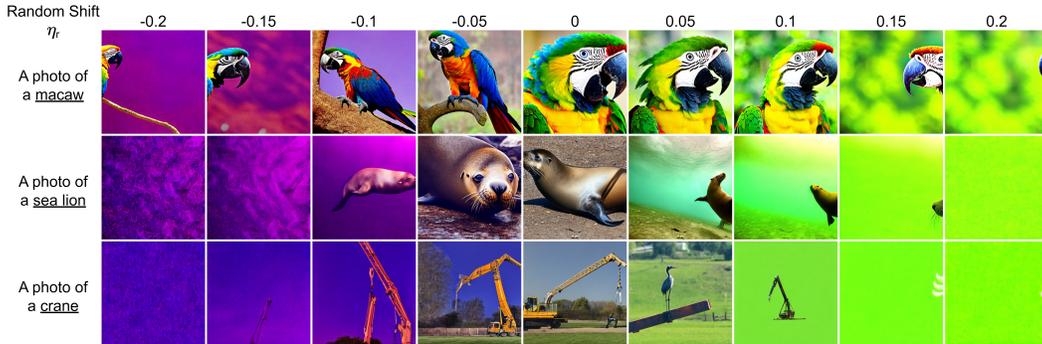


Figure 5: **Visual Inspection of image generated with Random Shift ( $\eta_r$ ) by Stable Diffusion v2.1.** We give the text prompt, “A photo of a \_\_\_” (row), and manipulate the initial random noise across varying levels of  $\eta_r$  (column). Note that as  $\eta_r$  deviates from zero, the images transition from accurate object representations to progressively loss of detail and color shift. Correspondingly, negative shifts cause purple hues, and positive shifts result in green hues.

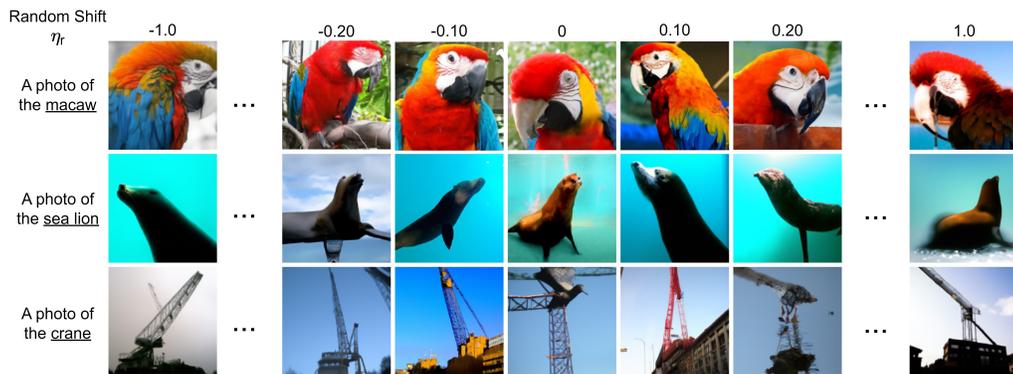


Figure 6: **Visual Inspection of image generated with Random Shift ( $\eta_r$ ) by GLIDE.** We give the text prompt, “A photo of a \_\_\_” (row), and manipulate the initial random noise across varying levels of  $\eta_r$  (column). Unlike Stable Diffusion, GLIDE exhibits consistency in generating images regardless of the degree of Random Shift applied, highlighting GLIDE’s superior reliability in generating images against shifts in the initial vector.

## 4.2 Evaluating Reliability Scores for Different Diffusion Models

**Experimental setup.** In this experiment, we compare variants of Stable Diffusion (Stable Diffusion v1.5 with and without guidance, Stable Diffusion v2.1 with and without guidance) and GLIDE using Random Shift( $\eta_r$ ). Although extremely rare, such a shift can occur due to sampling variability, which is possible in real-world scenarios. We compare by varying the perturbation  $\eta_r$  within the range  $[-0.3, 0.3]$  with an interval of 0.6. For GLIDE, we follow the originally proposed setting with 150 diffusion steps for the low-resolution image and 27 diffusion steps for upscaling it to the high-resolution image. For all Stable Diffusion models, we follow the original paper (Rombach et al., 2022) with 50 sampling time steps and a 7.5 classifier-free guidance scale.

**Result.** Figure 4 shows the performance of different diffusion models. The results reveal that, across all metrics evaluated, the performance of the latent-based diffusion models notably decreased as the

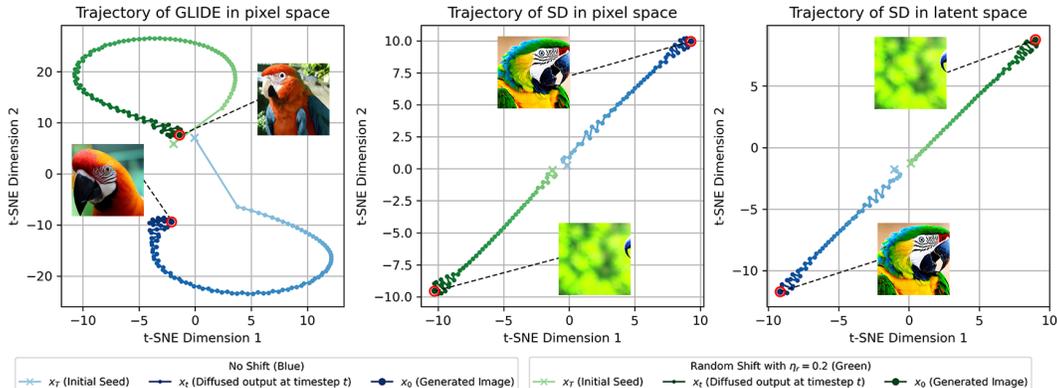


Figure 7: **The trajectories of the reverse diffusion process of GLIDE (Left) and Stable Diffusion (Center and Right) Models in Pixel Space (Left and Center) and Latent Space (Right) using t-SNE.** We plot the trajectories of the reverse diffusion process from the initial seed vector  $z$  and shifted seed vector  $\tilde{z}$ . Notably, we observe a phenomenon where GLIDE’s trajectory contains a large discontinuity in the early time step, visible between  $x_T$  (seed vector) and  $x_{T-1}$  (first noise correction). In contrast, the Stable Diffusion does not exhibit the same behavior, and the diffusion process diverges for modified seed vector  $\tilde{z}$  compared to the initial seed vector  $z$ . This comparison highlights the relative stability of GLIDE in generating consistent images regardless of shifts to initial seed vectors.

shift increased. In contrast, GLIDE maintains a consistent performance by remaining unaffected by the shift irrespective of the intensity. Initially, with no or subtle shifts, variations of the Stable Diffusion models with guidance exhibit a slightly better performance over GLIDE. However, as the shift intensity increases, a clear divergence in performance emerges. The latent diffusion models demonstrate a 50% drop in both top-1 and top-5 accuracy and a 30% drop in CLIP Score, while GLIDE’s performance remains comparatively unaffected. For visual comparison, we provide the images generated using multiple prompts by Stable Diffusion in Figure 5 and GLIDE in Figure 6. We note that in the case of no shift to the initial random vector, the image quality of Stable Diffusion is better compared to GLIDE. However, our aim is to evaluate the reliability of generated images in the context that the object given in the conditioning prompt is generated in the image rather than measuring the image quality.

Furthermore, a comparison within the latent-based diffusion models reveals distinct behavior based on the presence of guidance. Models with guidance not only outperform their counterparts without guidance, which aligns with findings from previous studies (Nichol et al., 2022), but also show greater reliability to increasing shift. This is reflected by a slower rate of performance drop in models with guidance under subtle shifts, as opposed to the steeper decline observed in models without guidance, reaching 0% in top-1 and top-5 accuracy, suggesting that the models fail to incorporate the conditioning variable in generating the outputs. Overall, the empirical result suggests that GLIDE’s reliable generations are more pronounced than that of the latent-based diffusion models.

## 5 Discussion

### 5.1 Extended Analysis of Stable Diffusion Model

**Disparities in Class-wise Reliable Generations:** The visualization in the Figure 5 distinctly showcases the varying levels of class-wise reliable generations exhibited by the model. Notably, even when the initial random vector is subjected to a Random Shift with  $\eta_r = 0.15$ , the model adeptly generates an image of a “macaw” with remarkable precision, isolating the subject from its background. In stark contrast, attempts to generate an image of a “crane” under similar conditions result in the object not being generated by the model. This discrepancy aligns with the observations of Samuel et al. (2023b), who underscore the notion that common concepts are reliably generated across a wider range of initial seed vectors. In contrast, generating images representing rare concepts

The Trajectories of Diffusion Models Generating Different Objects from the Same Seed Vectors

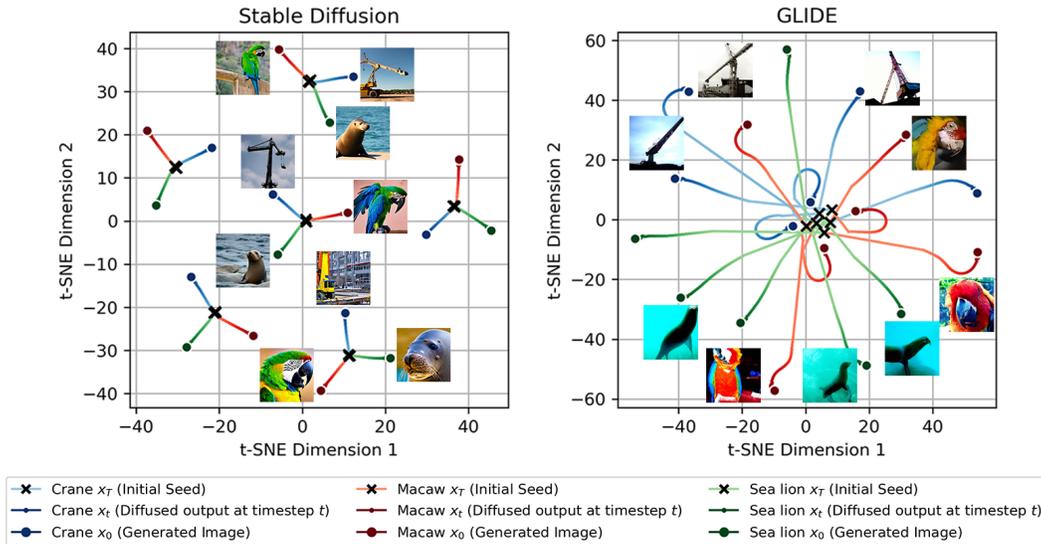


Figure 8: **The trajectories of reverse diffusion process of GLIDE and Stable Diffusion to generate different objects from the same initial seed vector.** The model is tasked to generate images originating from six sets of seed vectors, each set is conditioned by three prompts and each with the labels of "macaw", "crane", or "sea lion". We use t-SNE to reduce the visualisation of space to 2 dimensions. Originating from the same seed vectors  $x_T$ , Stable Diffusion’s sampling trajectory is inclined to traverse toward a local desired output. GLIDE, by comparison, tends to diverge, with its trajectory extending relatively farther from the seed vector toward the images  $x_0$ .

demands a meticulous selection of initial seed vectors, and a slight disturbance may break the diffusion model to generate the said object in the conditioning variable.

**Unintended Positional Shifts:** The complexity of the impact of shifts on the initial seed vector is further revealed when we notice the generated images of “macaw” in Figure 5. It is intriguing to note that the position of the generated object is strongly tied to the direction of the Random Shift. Positive Random Shift gradually displaces the “macaw” to the right, eventually causing it to vanish from the frame on the right side. Conversely, negative shifts induce a slow leftward movement of the “macaw”. This observation invites further exploration into the realm of position control (Mao et al., 2023), suggesting that manipulation of the initial seed could yield nuanced control over the final position of generated objects.

**Slight Random Shift boost performance:** Surprisingly, as evident in Table 1, Stable Diffusion outperforms its counterparts in cases arising from the slight shift from the standard normal distribution. In the context of Random Shift setting, Stable Diffusion attains the best top-1 and top-5 accuracy, particularly excelling in the presence of a slight positive shift  $\eta_r = \eta_m = 0.05$ . This noteworthy trend is also observed in Mean Shift experiments. Furthermore, when comparing Standard Deviation Shift and Mixed Shift scenarios, employing a mere Standard Deviation Shift of  $\eta_s = 0.1$  yields a top-1 Accuracy of 62.0% and a top-5 accuracy of 86.6%. Intriguingly, introducing a slight positive shift of  $\eta_m = 0.05$ , creating a Mixed Shift, significantly enhances performance to 65.5% top-1 and 87.1% top-5 accuracy. This observation suggests that shifts, when strategically applied, can propel the initial seed vector towards a more favorable starting position, thereby boosting overall performance as noted by (Samuel et al., 2023b,a; Mao et al., 2023).

## 5.2 Robustness of GLIDE to Synthetic Shifts

To discuss the GLIDE’s highly reliable generations compared to the Stable Diffusion model, we display the difference in sampling trajectories of the reverse diffusion process for both models with Figures 7 and 8.

Figure 7 shows the trajectories of the reverse diffusion process of GLIDE and Stable Diffusion using t-SNE. We plot the difference in trajectories of the reverse diffusion process for initial seed vector  $z$  and shifted seed vector  $\tilde{z}$ . Notably, GLIDE’s sampling trajectory exhibits a significant discontinuity in the early time step, a phenomenon we coin as Early Steps Discontinuity. Such discontinuity hinted that GLIDE "pulls" the shifted seed vectors back towards a standard normal distribution, maintaining its output quality. Consequently, this could account for GLIDE’s relative stability in consistently generating images despite transformations to initial shift vectors. On the contrary, the Stable Diffusion model does not display the same behavior. Its sampling trajectory with the shifted seed  $\tilde{z}$  diverges when compared to the trajectory of  $z$ .

We further discuss the contrasting trajectories of the reverse diffusion process of GLIDE and Stable Diffusion for generating different objects using conditioning variable  $c$  in response to the same seed vectors with Figure 8. Originating from the same seed vectors, Stable Diffusion’s sampling trajectory tends to traverse towards a local desired output, indicating that it is more inclined to search locally from the random seed vector for the desired output. On the other hand, GLIDE sampling process exhibits a divergent pattern, with its trajectory extending relatively further from the seed vector toward the generated images, reflecting a broader exploration of the solution space.

We speculate that these factors contribute to the generation of non-reliable images by Stable Diffusion when initial seed vector  $z$  is transformed using  $\eta$  compared to the generation of reliable images by GLIDE. In summary, this analysis sheds light on the differences in Stable Diffusion and GLIDE the reasoning for the more reliable generation of images by GLIDE compared to Stable Diffusion, emphasizing the importance of the chosen diffusion process and training strategy in influencing the overall performance and resilience of diffusion models.

## 6 Conclusion

This paper conducts a comprehensive analysis of the generated images by the Stable Diffusion and GLIDE models when the initial seed vector of diffusion models is subjected to transformations. Our findings indicate that the state-of-the-art latent-based diffusion model Stable Diffusion struggles to effectively manage diverse shifts generating objects irrespective of the initial seed vector as the scale of the shift increases. On the other hand, a relatively older diffusion model GLIDE demonstrates a higher resiliency to handle such shifts to the initial seed vector. Through a combination of experimental and theoretical approaches, we identify and elucidate the factors contributing to GLIDE’s superior reliable generations compared to Stable Diffusion. We anticipate that our work will serve as a foundational resource for researchers aiming to design diffusion models that are simultaneously stable and reliable in generating objects.

## References

- Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-DPM: an analytic estimate of the optimal reverse variance in diffusion probabilistic models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=0xiJLKH-ufZ>.
- Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–10, 2023.
- Giannis Daras, Yuval Dagan, Alexandros G Dimakis, and Constantinos Daskalakis. Consistent diffusion models: Mitigating sampling drift by learning to be consistent. *arXiv preprint arXiv:2302.09057*, 2023.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*, 2022.
- Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs, Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior for video diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22930–22941, 2023.

- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *The Journal of Machine Learning Research*, 23(1): 2249–2281, 2022.
- Rongjie Huang, Max W. Y. Lam, Jun Wang, Dan Su, Dong Yu, Yi Ren, and Zhou Zhao. Fastdiff: A fast conditional diffusion model for high-quality speech synthesis. In Lud De Raedt (ed.), *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pp. 4157–4163. International Joint Conferences on Artificial Intelligence Organization, 7 2022a. doi: 10.24963/ijcai.2022/577. URL <https://doi.org/10.24963/ijcai.2022/577>. Main Track.
- Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren. Prodiff: Progressive fast diffusion model for high-quality text-to-speech. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 2595–2605, 2022b.
- Yujin Jeong, Wonjeong Ryoo, Seunghyun Lee, Dabin Seo, Wonmin Byeon, Sangpil Kim, and Jinkyu Kim. The power of sound (tpos): Audio reactive video generation with stable diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7822–7832, 2023.
- Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-tts: A diffusion model for text-to-speech via classifier guidance. In *International Conference on Machine Learning*, pp. 11119–11133. PMLR, 2022.
- Zhifeng Kong, Wei Ping, Jiayi Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=a-xFK8Ymz5J>.
- Mingxiao Li, Tingyu Qu, Wei Sun, and Marie-Francine Moens. Alleviating exposure bias in diffusion models through sampling with shifted time steps. *arXiv preprint arXiv:2305.15583*, 2023a.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328–4343, 2022.
- Yangming Li, Zhaozhi Qian, and Mihaela van der Schaar. Do diffusion models suffer error propagation? theoretical analysis and consistency regularization. *arXiv preprint arXiv:2308.05021*, 2023b.
- Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVII*, pp. 423–439. Springer, 2022.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11461–11471, 2022.
- Jiafeng Mao, Xueting Wang, and Kiyoharu Aizawa. Guided image synthesis via initial image editing in diffusion model. *arXiv preprint arXiv:2305.03382*, 2023.
- Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and editing with text-guided diffusion models. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 16784–16804. PMLR, 2022. URL <https://proceedings.mlr.press/v162/nichol22a.html>.
- Mang Ning, Enver Sangineto, Angelo Porrello, Simone Calderara, and Rita Cucchiara. Input perturbation reduces exposure bias in diffusion models. *arXiv preprint arXiv:2301.11706*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=TTdIXIpzhoI>.
- Dvir Samuel, Rami Ben-Ari, Nir Darshan, Haggai Maron, and Gal Chechik. Norm-guided latent space exploration for text-to-image generation. *arXiv preprint arXiv:2306.08687*, 2023a.
- Dvir Samuel, Rami Ben-Ari, Simon Raviv, Nir Darshan, and Gal Chechik. It is all about where you start: Text-to-image generation with seed selection. *arXiv preprint arXiv:2304.14530*, 2023b.
- AMBESH SHEKHAR. ImageNet100. <https://www.kaggle.com/datasets/ambityga/imagenet100>, 2021.
- Mannat Singh, Laura Gustafson, Aaron Adcock, Vinicius de Freitas Reis, Bugra Gedik, Raj Prateek Kosaraju, Dhruv Mahajan, Ross Girshick, Piotr Dollár, and Laurens Van Der Maaten. Revisiting weakly supervised pre-training of visual perception models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 804–814, 2022a.
- Vedant Singh, Surgan Jandial, Ayush Chopra, Siddharth Ramesh, Balaji Krishnamurthy, and Vineeth N Balasubramanian. On conditioning the input noise for controlled image generation with diffusion models. *arXiv preprint arXiv:2205.03859*, 2022b.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- Chen Henry Wu and Fernando De la Torre. Making text-to-image diffusion models zero-shot image-to-image editors by inferring "random seeds". In *NeurIPS 2022 Workshop on Score-Based Methods*, 2022.
- Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric diffusion model for molecular conformation generation. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=PzcvxEMzvQC>.
- Bo Zhao, Lili Meng, Weidong Yin, and Leonid Sigal. Image generation from layout. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8584–8593, 2019.

## A Analysis of Overlap of Shifted Distributions with respect to Standard Normal Distribution

Sohl-Dickstein et al. (2015) assumes that the initial seed vector for diffusion models be sampled from the standard normal distribution for generating high-fidelity images from diffusion models. We propose to use the synthetic shifts, which follow a normal distribution but keep the difference to the standard normal distribution as small as possible. Moreover, we do not want to shift the transformations too drastically because it goes beyond the range of learned latent values by diffusion models. When designing the experiments, we carefully considered the distribution overlap to ensure that the transformations are not too drastic and have a high overlap with standard normal distribution. The overlap percentage of shifted distribution to standard normal distribution suggested by (Sohl-Dickstein et al., 2015) can be referred in Table 2. Most cases of scaled shifts have over 90% overlap with the standard normal distribution, while all the evaluated scaled shifts have an overlap of more than 80% with the standard normal distribution. This implies that the latent values in the shifted seed vectors have a relatively high likelihood of being sampled from the ideal normal distribution.

## B Extended Analysis for Robustness of GLIDE to Synthetic Shifts

Commencing with an initial seed vector  $z$  (depicted in the “blue” cross), distinct shifts are sequentially applied (depicted in various colored crosses). The initial seed vector  $z$ , and the shifted seed vector  $\tilde{z}$ , follow distinct reverse diffusion process trajectories, as illustrated in Figure 9. The temporal evolution of the reverse diffusion process is displayed in Figure 10.

In Figure 9, despite the visual similarity among trajectories, they yield entirely different outputs. Each time step closely approximates the preceding one, and the final output remains in proximity to the initial points, a phenomenon corroborated by Figure 10 (Top). This figure underscores that, early on, each path outlines the intended content to be generated. Subsequently, additional details are progressively incorporated to enhance image clarity. However, it is crucial to note that even a seemingly negligible perturbation can mislead Stable Diffusion, leading to the generation of unintended content.

In the Figure 10 (Bottom), it becomes evident that a substantial discontinuity emerges in the early time steps regardless of the starting point. This discontinuity is crucial in steering the shifted seed vector to a more favorable starting point. This phenomenon is particularly pronounced in the Arrangement Shift case. Examining the last row of Figure 10 (Bottom), we observe an arrangement patch situated at the top left during the initial time step. However, after a few subsequent time steps, it vanishes. This observation underscores the significance of discontinuity in effectively managing shifted seed vectors.

Table 2: **Distribution overlap between initial and shifted seed vectors.** The table shows how the shifted seed vectors differ from the initial seed vectors. Random Shift and Mean Shift are more than 90% overlap with the initial seed vector. Standard Deviation Shift and Mixed Shift are all higher than 80%. Regarding Arrangement Shift, since we are only switching the values, it has 100% overlap with the initial seed vector, regardless of  $\eta_{a,*}$ .

	← (Negative Shift)			(No Shift)	→ (Positive Shift)						
	-0.30	-0.20	-0.15	-0.10	-0.05	0.00	0.05	0.10	0.15	0.20	0.30
Overlap of Distribution (†)	94.02%	96.01%	97.01%	98.00%	99.00%	100.0%	99.00%	98.00%	97.01%	96.01%	94.02%
				Mean Shift ( $\eta_m$ )							
	-0.20	-0.15	-0.10	-0.05	0.00	0.05	0.10	0.15	0.20		
Overlap of Distribution (†)	92.03%	94.02%	96.01%	98.00%	100.0%	98.00%	96.01%	94.02%	92.03%		
				Standard Deviation Shift ( $\eta_s$ )							
	-0.30	-0.20	-0.10	0.00	0.10	0.20	0.30				
Overlap of Distribution (†)	83.00%	89.24%	94.90%	100.0%	95.39%	91.20%	87.37%				
				Mixed Shift ( $\eta_s, \eta_m$ )							
	(-0.30, -0.15)	(-0.20, -0.10)	(-0.10, -0.05)	(0.00, 0.00)	(0.10, 0.05)	(0.20, 0.10)	(0.30, 0.15)				
Overlap of Distribution (†)	82.00%	88.58%	94.60%	100.0%	95.10%	90.66%	86.60%				
				Arrangement Shift ( $\eta_a$ )							
				0	8	16	32	64			
Overlap of Distribution (†)				100.0%	100.0%	100.0%	100.0%	100.0%			

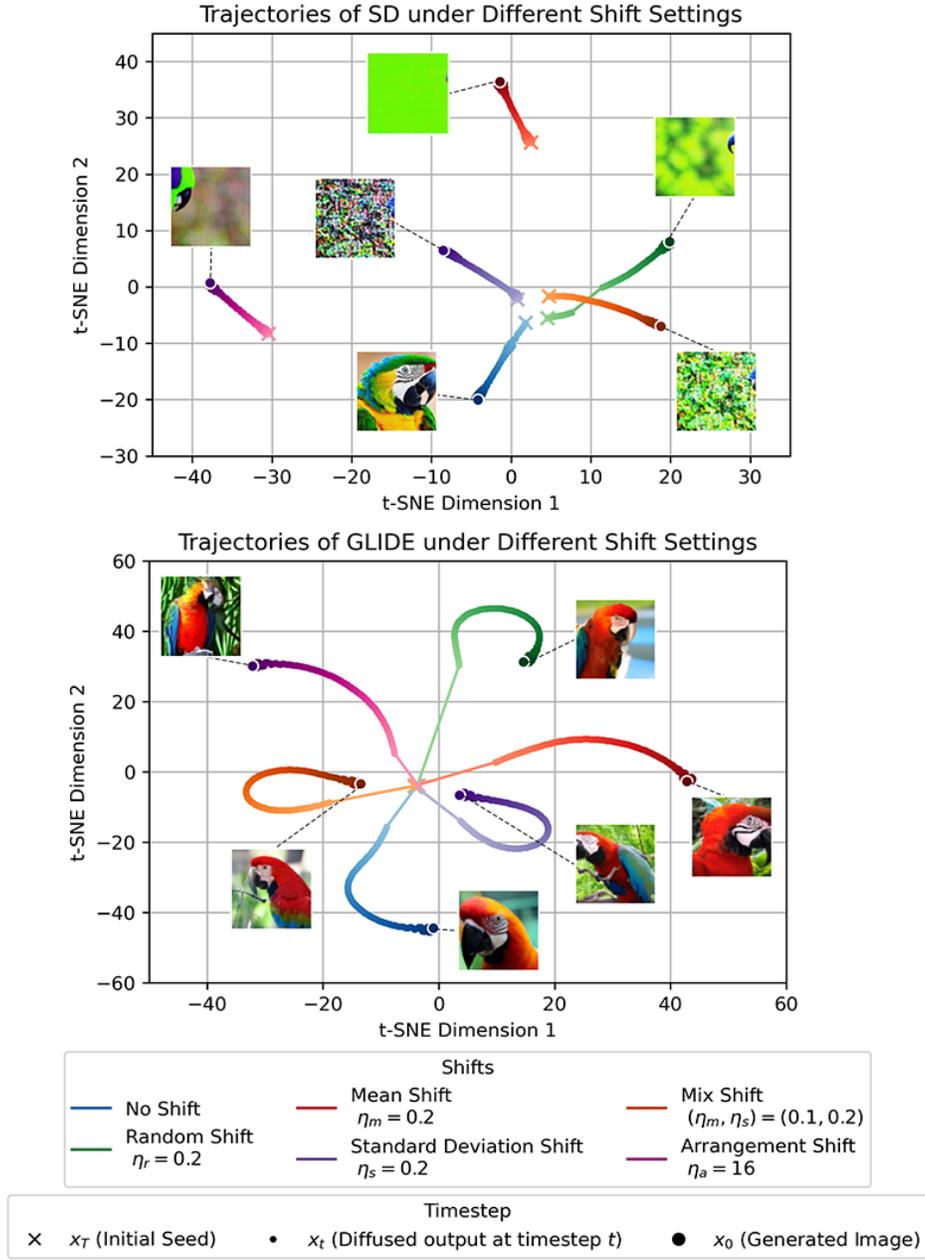


Figure 9: **Visual Inspection of reverse diffusion process trajectories of Stable Diffusion and GLIDE with shifts.**

### C Visual Examples of Image Generation by Different Diffusion Models for Synthetic Shifts

The outcomes of both Stable Diffusion and GLIDE shifts are illustrated in the following Figures 11 - 20. Notably, GLIDE consistently performs uniformly across various shifts. In both Random Shift (Figures 11 and 12) and Mean Shift (Figures 13 and 14) scenarios, Stable Diffusion demonstrates similar performance. As  $\eta_r$  or  $\eta_m$  deviates from zero, there is a perceptible shift in the images from accurate depictions of objects to a gradual loss of detail and a color shift. Negative shifts result in the presence of purple tones, while positive shifts bring about green tones. Regarding Standard Deviation Shift (Figures 15 and 16), positive shifts introduce noise and challenge item identification.

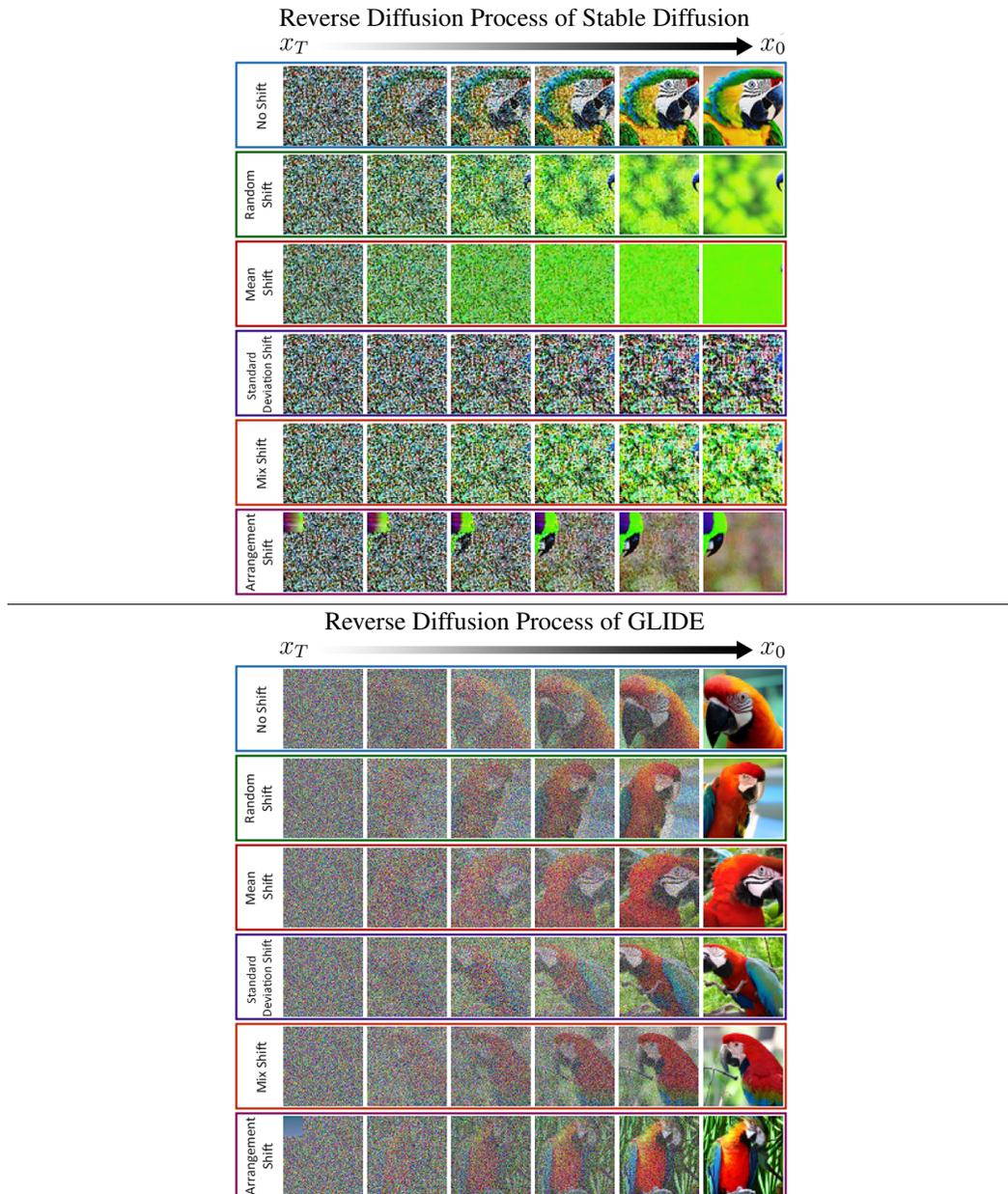


Figure 10: **Visual Inspection of the images generated in the reverse diffusion process of Stable Diffusion and GLIDE with shifts.**

Conversely, negative shifts reduce detail, as exemplified in the “A photo of a cock” row, where feather texture and cock face detail diminish with further negative shifts. Mixed Shift (Figures 17 and 18) combines both effects, altering color and details. Regarding Arrangement Shift (Figures 19 and 20), arranging 8 pixels from the top has no adverse impact on image quality. However, increasing it to 16 pixels leads to the inability to generate the left top position and further increases blur in the remaining portion. Starting from arranging 32 pixels, no image can be generated. These systematic experiments offer insights into how Stable Diffusion responds to perturbations across various scenarios.

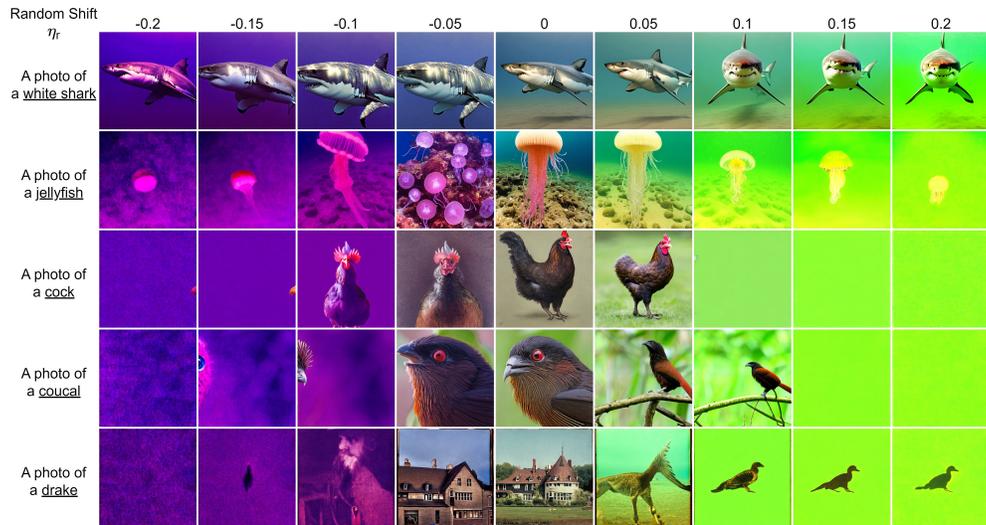


Figure 11: Visual Inspection of image generated with Random Shift ( $\eta_r$ ) by Stable Diffusion v2.1.

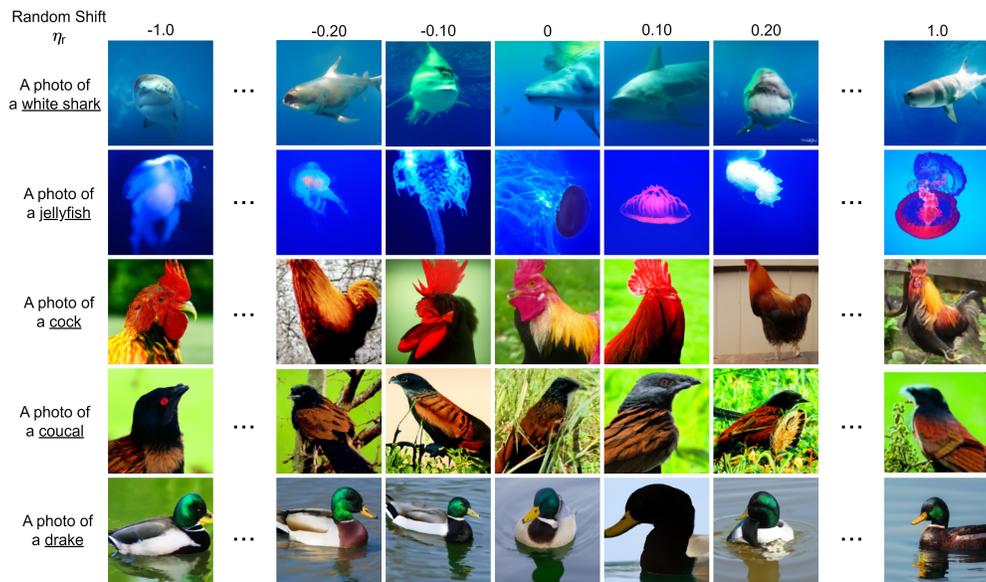


Figure 12: Visual Inspection of image generated with Random Shift ( $\eta_r$ ) by GLIDE.



Figure 13: Visual Inspection of image generated with Mean Shift ( $\eta_m$ ) by Stable Diffusion v2.1.

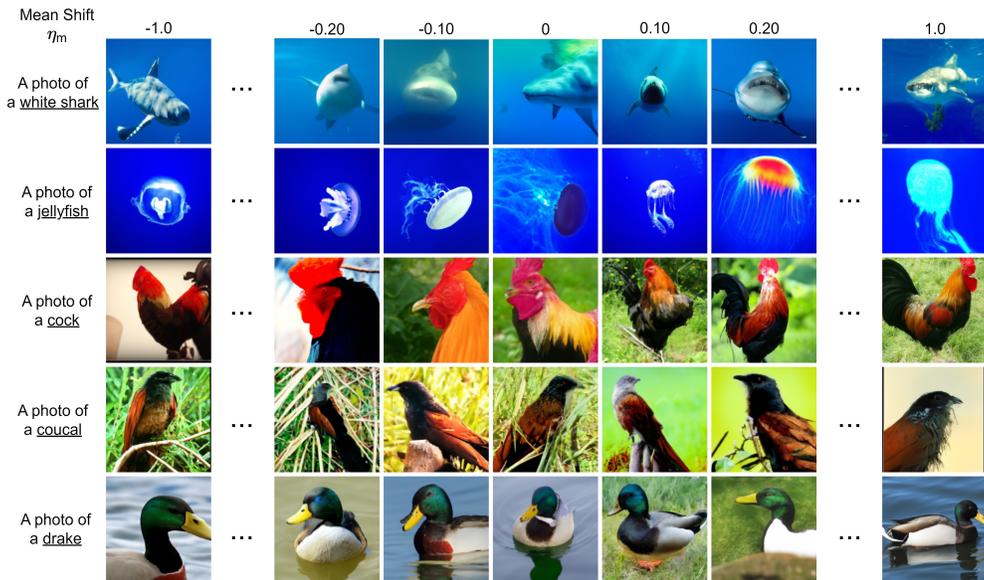


Figure 14: Visual Inspection of image generated with Mean Shift ( $\eta_m$ ) by GLIDE.

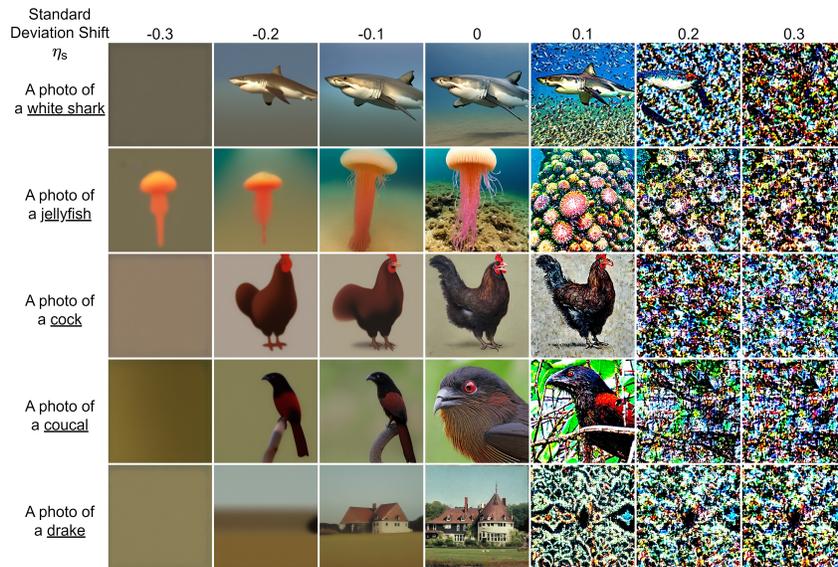


Figure 15: Visual Inspection of image generated with Standard Deviation Shift ( $\eta_s$ ) by Stable Diffusion v2.1.

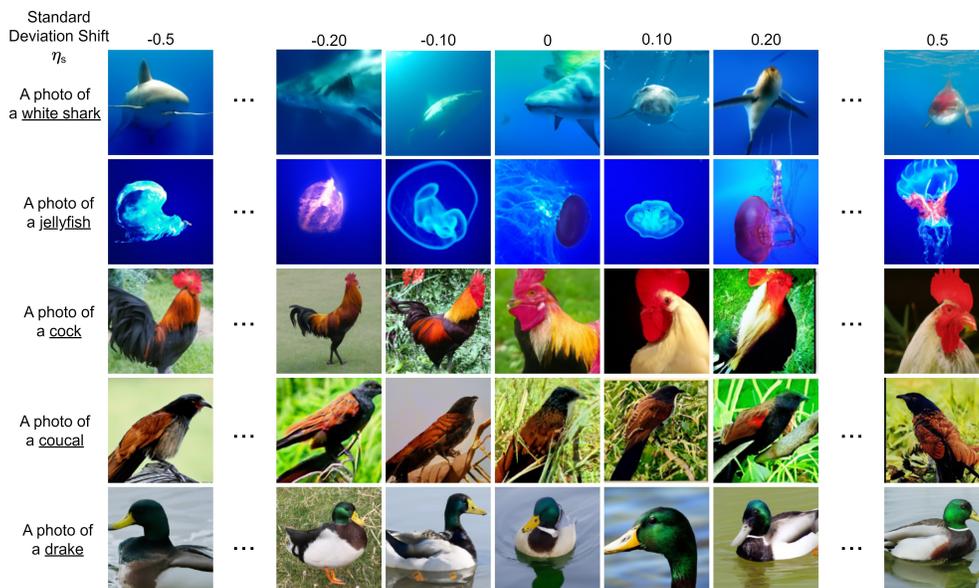


Figure 16: Visual Inspection of image generated with Standard Deviation Shift ( $\eta_s$ ) by GLIDE.

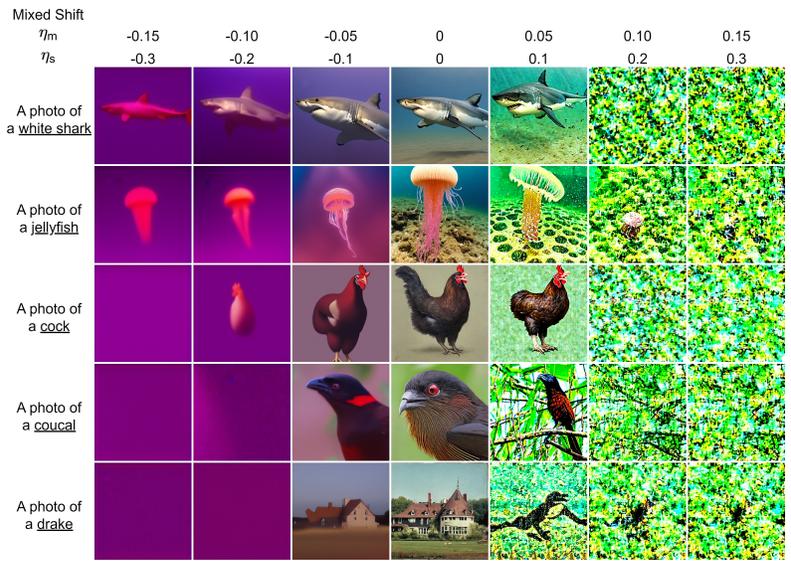


Figure 17: Visual Inspection of image generated with Mixed Shift by Stable Diffusion v2.1.

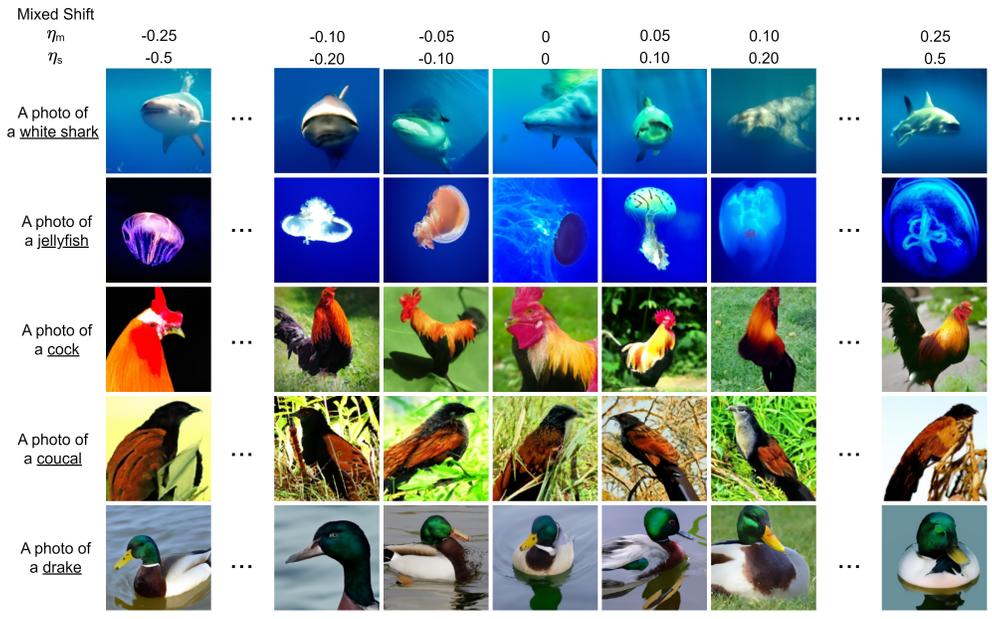


Figure 18: Visual Inspection of image generated with Mixed Shift by GLIDE.

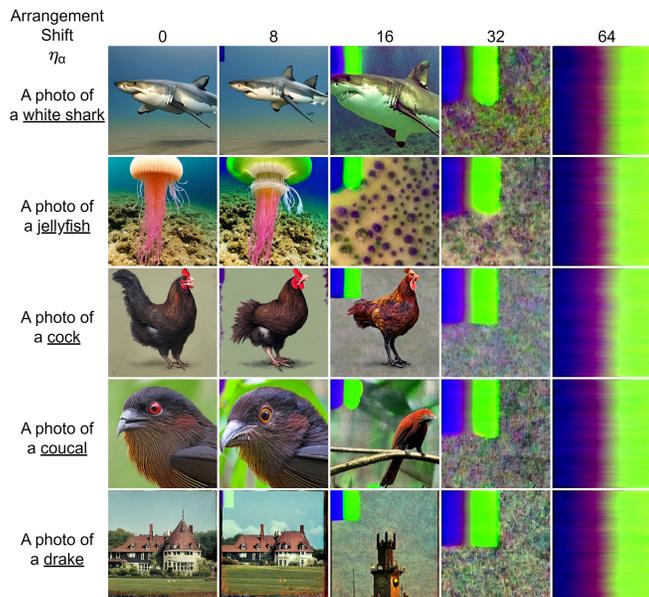


Figure 19: Visual Inspection of image generated with Arrangement Shift by Stable Diffusion v2.1.

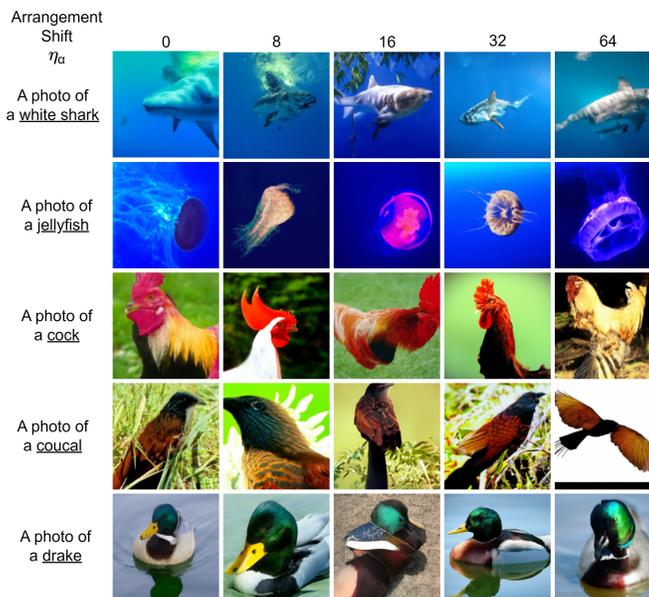


Figure 20: Visual Inspection of image generated with Arrangement Shift by GLIDE.