
THERE AND BACK AGAIN: ON THE RELATION BETWEEN NOISES, IMAGES, AND THEIR INVERSIONS IN DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Denoising Diffusion Probabilistic Models (DDPMs) achieve state-of-the-art performance in synthesizing new images from random noise, but they lack meaningful latent space that encodes data into features. Recent DDPM-based editing techniques try to mitigate this issue by inverting images back to their approximated starting noise. In this work, we study the relation between the initial Gaussian noise, the samples generated from it, and their corresponding latent encodings obtained through the inversion procedure. First, we interpret their spatial distance relations to show the inaccuracy of the DDIM inversion technique by localizing latent representations manifold between the initial noise and generated samples. Then, we demonstrate the peculiar relation between initial Gaussian noise and its corresponding generations during diffusion training, showing that the high-level features of generated images stabilize rapidly, keeping the spatial distance relationship between noises and generations consistent throughout the training.

1 INTRODUCTION

Diffusion-based probabilistic models (DDPMs) (Sohl-Dickstein et al., 2015), have surpassed state-of-the-art solutions in many generative domains including image (Dhariwal & Nichol, 2021), speech (Popov et al., 2021), video (Ho et al., 2022), and music (Liu et al., 2021) synthesis. Nevertheless, one of the significant drawbacks that distinguishes diffusion-based approaches from other generative models like Variational Autoencoders (Kingma & Welling, 2014), Flows (Kingma & Dhariwal, 2018), or Generative Adversarial Networks (Goodfellow et al., 2014) is the lack of implicit latent space that encodes training data into low-dimensional, interpretable representations. Several works try to mitigate this issue, treating as latent features internal data representations extracted from the denoising UNet model (Kwon et al., 2022), combining diffusion models with additional external models (Preechakul et al., 2021), or by seeking structure in the starting noise used for generations (Song et al., 2020).

The last example, introduced by (Song et al., 2020) with Denoising Diffusion Implicit Models (DDIM), led to the proliferation of methods grouped under the name of inversion techniques (Garibi et al., 2024; Mokady et al., 2023; Huberman-Spiegelglas et al., 2024; Meiri et al., 2023; Hong et al., 2024). The main idea behind those approaches is to use a diffusion model to predict the noise that can be added to the original or generated image. Applying this procedure several times allows tracing back the backward diffusion process and approximating the initial noise that results in the starting image. However, due to approximation error and biases induced by a trained denoising model, there are inconsistencies between initial Gaussian noise and the reversed so-called *latent* representation. While recent works try to improve the approximation and mitigate discrepancies between noise and latent, in this work, we propose to closely study the DDIM inversion procedure and highlight the relation between Gaussian noise, generated samples, and their latents.

First, we emphasize the main differences between initial Gaussian noise and latent codes. We locate the latent space between the initial Gaussian noise and generated samples, showing that DDIM inversion does not properly turn the image into noise. We show how this relation changes with time, highlighting the importance of early training steps in the formation of sample-to-latent mappings. Second, we move to the analysis of the relation between noise and samples. We show that it is

possible to correctly assign initial Gaussian noise to the generated sample with Euclidean distance. We further deepen this analysis to a non-stationary setup, showing that the mapping between noises and generations emerges at the very beginning of the diffusion model training.

The main contribution of this work can be summarized as follows:

- We show that reverse DDIM produces latent representations that are not standard multivariate Gaussian, creating a gap between the diffusion models’ theory and practice.
- We study how this relation changes with training and show that improving the generative capabilities of the model does not improve the accuracy of reverse DDIM.
- We show that the relation between images and noises is defined by the simple L2 distance at the early stage of the training.

2 BACKGROUND

The training of diffusion models consists of forward and backward diffusion processes, where in the context of Denoising Diffusion Probabilistic Models (DDPMs), the former one with training image x_0 and $\{\beta_t\}_{t=1}^T$ being some variance schedule, can be expressed as

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, \quad (1)$$

with $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$, and $\epsilon_t \sim N(0, I)$.

In the backward process, the noise is gradually removed starting from a random noise $x_T \sim N(0, I)$,

$$x_{t-1} = \sqrt{\frac{\bar{\alpha}_{t-1}}{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_\theta(x_t, t, c)) + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(x_t, t, c) + \sigma_t z_t, \quad (2)$$

with $\sigma_t = \eta\beta_t(1 - \bar{\alpha}_{t-1})/(1 - \bar{\alpha}_t)$ being a variance schedule, $\epsilon_\theta(x_t, t, c)$ being the output of a network trained to remove noise, and $z_t \sim N(0, I)$.

While for $\eta = 1$, a non-deterministic DDPM model is used, setting $\eta = 0$ removes the random component from the equation 2, making it a Denoised Diffusion Implicit Model (DDIM), characterized by a deterministic mapping from noise space x_T to image space x_0 .

By removing the stochasticity of the DDPM sampling, we can additionally reverse deterministically the direction of the backward diffusion process and encode images back to the original noise space. The DDIM inversion is obtained by rewriting equation 2 as:

$$x_t = \sqrt{\frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}}}x_{t-1} + \left(\sqrt{1 - \bar{\alpha}_t} - \frac{\sqrt{\bar{\alpha}_t - \bar{\alpha}_t\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_{t-1}}}\right) \cdot \epsilon_\theta(x_t, t, c) \quad (3)$$

However, due to circular dependency on $\epsilon_\theta(x_t, t, c)$, DDIM inversion approximates this equation by assuming linear trajectory with direction to x_T in t -th step being same as in $(t - 1)$ -th step, i.e.,

$$\epsilon_\theta(x_t, t, c) \approx \epsilon_\theta(x_{t-1}, t, c). \quad (4)$$

While such approximation is often sufficient to obtain good reconstructions of images, it introduces the error that depends on the difference $(x_t - x_{t-1})$, which can be detrimental for models that leverage a few diffusion steps or use the classifier-free guidance Ho & Salimans (2021); Mokady et al. (2023). In this work, we empirically study the consequences of this approximation error.

We will be interested in studying the relations between the following three objects:

- Gaussian noise variable, \mathbf{x}^T , used to generate an image through a diffusion process.
- Image sample, \mathbf{x}^0 , the outcome of the generation process produced by the Diffusion Model, i.e., the result of going through denoising stages, starting from $t = T$ and ending on $t = 0$.
- Latent variable, $\hat{\mathbf{x}}^T$, the result of starting from \mathbf{x}^0 and applying T steps of a reversed DDIM generation process, see equation 4.

3 RELATED WORK

Diffusion models inversion techniques Thanks to the reversible sampling procedure, Denoising Diffusion Implicit Models are often employed in tasks such as inpainting (Zhang et al., 2023), image (Su et al., 2022; Kim et al., 2022; Hertz et al., 2022), video (Ceylan et al., 2023) or speech (Deja et al., 2023a) edition. However, the baseline DDIM approach is based on the assumption that the prediction of the noise removed from the image in the t -th backward diffusion step closely approximates the noise of the $(t - 1)$ -th step. This assumption is not always true, and recent works try to mitigate the discrepancies of such approximation. In particular, Renoise (Garibi et al., 2024) iteratively improves the prediction of added noise using the predictor-corrector technique. This allows to closely estimate the original Gaussian noise for a given image, even with fewer diffusion steps. Several works aim to reverse the diffusion process in text-to-image models. In such a case, the prompt selection significantly influences the final latent. To mitigate this issue, Null-text inversion method (Mokady et al., 2023) extends the DDIM inversion with optimized pivotal noise vectors and additional Null-text optimization technique, where unconditional textual embeddings employed by classifier free-guidance (Ho & Salimans, 2021) are optimized in order to reduce the reconstruction error. Other work (Huberman-Spiegelglas et al., 2024) proposes an alternative DDPM noise space, where noise maps do not have a normal distribution, yet it enables better image editing capabilities and perfect reconstructions. On the other hand, a regularized Newton-Raphson inversion method (Meiri et al., 2023) formulates the inversion process as solving an implicit equation using numerical techniques, achieving faster and higher-quality reconstructions with prompt-aware adjustments, enabling improvements in image interpolation, and boosting models’ diversity. Furthermore, exact inversion methods for DPM-solvers (Hong et al., 2024) mitigate the errors introduced by classifier-free guidance by leveraging gradients, boosting both the image and noise reconstruction.

Latent space in diffusion models One of the significant drawbacks of diffusion-based generative models compared to approaches such as VAEs, Flows, or GANs is the lack of meaningful latent space that encodes features of the generated samples. Several approaches try to mitigate this issue. Kwon et al. (2022) show that so-called *h-features*, which are the activations located inside of the U-Net model used as a diffusion decoder, can be used as meaningful representations providing space for semantically coherent image manipulation. This idea is further extended by Park et al. (2023), where the authors show that we can calculate the pullback metric that directly associates *h-features* with the original image space. Several works show that we can directly benefit from such features in downstream tasks such as image segmentation (Baranchuk et al., 2021; Tumanyan et al., 2023; Rosnati et al., 2023), image correspondence (Luo et al., 2024) or classification (Deja et al., 2023b).

Noise-to-image mapping in diffusion models Contrary to latent variable models such as VAEs, the forward diffusion process that maps images to the Gaussian noise is a parameter-free process that can be understood as hierarchical VAE with pre-defined unlearnable encoder Kingma et al. (2021). However, several works show interesting properties resulting from the training objective of DDPMs and score-based models. Kadkhodaie et al. (2024) show that due to inductive biases of denoising models, different DDPMs trained on similar datasets converge to almost identical solutions. This idea is further explored by Zhang et al. (2024), where the authors show that even models with different architectures converge to the same score function and, hence, the same noise-to-generations mapping. This mapping itself is further analyzed by Khrulkov & Oseledets (2022), where the authors show that the encoder map coincides with the optimal transport map for common distributions. In this work, we extend this analysis further by empirically validating that examples are generated from the closest random noise, even for more complex distributions.

4 EXPERIMENTS

4.1 EXPERIMENTAL DETAILS

For the training, we follow Nichol & Dhariwal (2021) and train two unconditional pixel-space diffusion models on ImageNet and CIFAR10 datasets. In our studies we use three diffusion models, a DDPM-based Dhariwal & Nichol (2021) model trained on CIFAR10 and ImageNet, and a Latent Diffusion Model (LDM) trained on CelebA dataset. The CIFAR10 model was trained for 500K steps, while the ImageNet one for 1.5M steps. Both DDPM models use 4K diffusion steps during

their training. For both diffusion models trained on the CIFAR10 and ImageNet datasets, we investigate the noising and denoising process with a number of diffusion steps T varying from 50 up to 4K. For our experiments, we average our metrics over 1K samples generated from 3 models trained with 3 random training seeds.

4.2 NOISE \neq LATENT

Numerous methods employ reversed DDIM in applications such as image editing or inpainting. In this case, the underlying assumption is that by *encoding* the image back with a denoising decoder, we can obtain the original *noise* that can be used to reconstruct the original image. However, as noticed by recent works Garibi et al. (2024); Parmar et al. (2023); Hong et al. (2024), the latents created with DDIM inversion (\hat{x}^T) often do not follow standard multivariate Gaussian distribution, creating a gap between the promise of diffusion models’ theory and practice. We study the implications of this fact and show its far-reaching consequences.

We visualize this phenomenon in Figure 1 for three models, showing the latent or its difference from the generated image. For simpler ones trained with smaller datasets, such as CIFAR10, we can observe clear structures of original images in the inverted latents, as presented in Figure 1 (C). However, even for more complex datasets, we can highlight the inversion error by plotting the image difference between the latent and the noise, as presented for DDPM trained on ImageNet and LDM trained on CelebA (Figure 1 (A, B)). We can observe that the highest inversion error can be observed in the areas of large monotonic surfaces, where more high-frequency noise needs to be added.

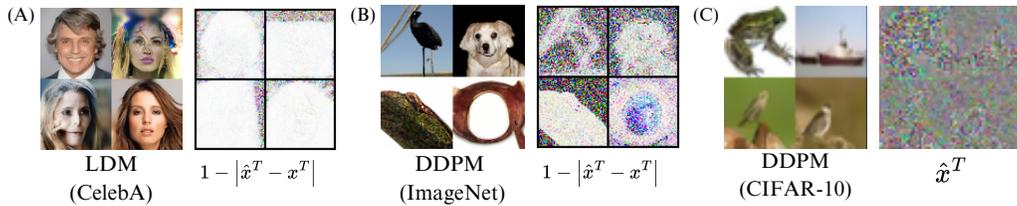


Figure 1: Visualization of samples and latent images from three diffusion models. For smaller DDPM model trained on CIFAR10 (C) we can observe original image structure in the latent calculated with reverse DDIM. We can observe similar structures on image differences from larger models (A and B).

To measure the extent of this effect, in Table 1, we show the values of the top-10 Pearson correlation coefficients measured between individual pixels of either initial noise (x^T), latent (\hat{x}^T), or samples (x^0). We can observe correlated pixels in the latents, which further supports the claim that they do not come from the standard multivariate Gaussian distribution. This is especially true for models trained on smaller datasets such as CIFAR-10 and almost invisible for latent models.

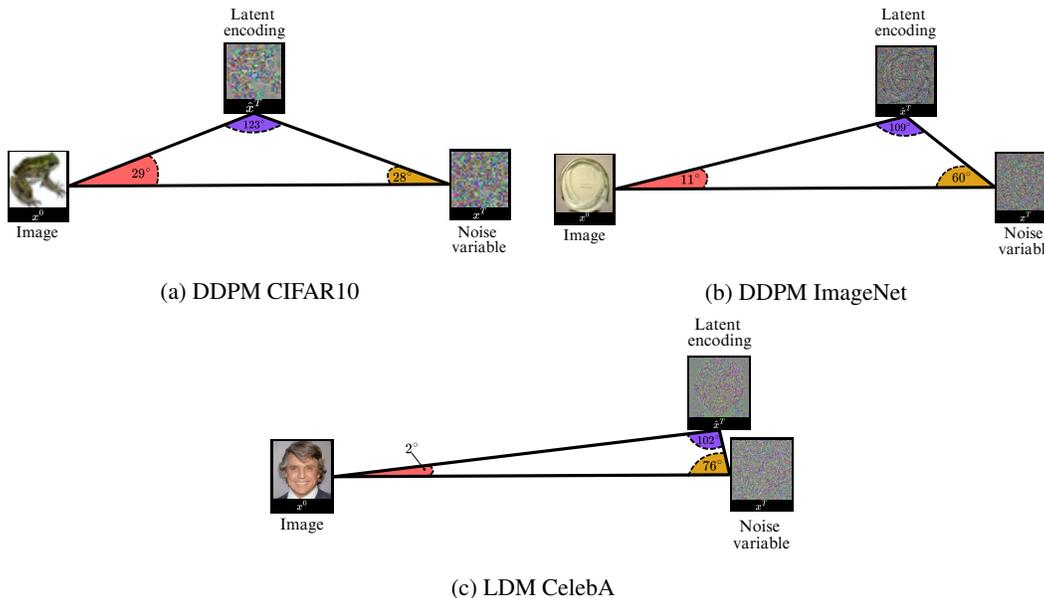
	DDPM (CIFAR10)	DDPM (ImageNet)	LDM (CelebA)
Noise (x^T)	0.159 ± 0.003	0.177 ± 0.007	
Latent (\hat{x}^T)	0.462 ± 0.009	0.219 ± 0.006	0.179 ± 0.008
Sample (x^0)	0.986 ± 0.001	0.966 ± 0.001	0.904 ± 0.005

Table 1: Top-10 correlation coefficients in random Gaussian noise vs. latent. We can observe that latents created with reverse DDIM have correlated pixel values. This is especially visible for smaller DDPM models and almost not noticeable for LDM.

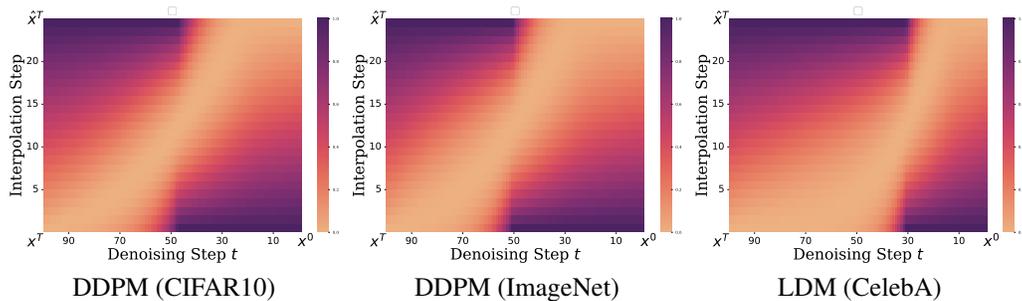
4.3 WHAT IS THE LOCATION OF THE LATENT VARIABLE?

Given the results of the previous analysis, we pose the next question: *what is the location of the latent variables calculated with the reverse DDIM procedure?* We show that the DDIM-generated latent (\hat{x}^T) is located along the trajectory x_t between the Gaussian noise (x^T) and the generated sample (x^0), and discuss how this relation changes with different characteristics of diffusion models.

216 In the first experiment, shown in Figure 2, we calculate the angles between noises, samples, and
 217 latents. We average the angles across 1K samples from three different models and plot the triangles
 218 in 2D space. We can observe that in each case, the angle located at the image’s vertex, $\angle x^0$, is
 219 acute but never zero. On the other hand, the angle located at the vertex representing latents, $\angle \hat{x}^T$, is
 220 always obtuse. This leads to the conclusion that due to the imperfect approximation of the reverse
 221 DDIM procedure, latents are located along the trajectory x_t of the generated image.



242 Figure 2: Visualisation of the most probable relation between random Gaussian noises, their corre-
 243 sponding samples, and latents recovered with reverse DDIM procedure for three different models.
 244



256 Figure 3: Distances between intermediate steps of backward diffusion process, and the interpolated
 257 points between initial noise x^T and the inverted latents \hat{x}^T . We can observe that independently
 258 of the model, the consecutive intermediate generations along the sampling trajectory initially get
 259 closer to the latent variable up to the point where after approximately 50-70% the generations pass
 260 the latent. This aligns with the visualization of the most probable relations in Figure 2.

261 To approximate the exact location of the latents more closely, we analyze the distance between
 262 different steps of the backward diffusion process and the Noise-Latent side of the triangle, i.e., the
 263 interval $x^T \hat{x}^T$, see Figure 3. Each pixel, with coordinates (t, λ) , is colored according to the L2
 264 distance between the intermediate step of trajectory x_t and the corresponding interpolation step, i.e.
 265 $\|(1 - \lambda)x^T + \lambda\hat{x}^T - x_t\|_2$. Figure 3 shows that while moving from the random Gaussian noise
 266 (x^T) towards the final sample (x^0), the intermediate steps x_t are getting closer to the latent
 267 (\hat{x}^T), which is visually represented by the light-colored area in the upper right corner of the plot. This
 268 also demonstrates that from some time onwards, say t_0 , the trajectory x_t closest point to the interval
 269 $x^T \hat{x}^T$ is the right endpoint \hat{x}^T . This has non-trivial consequences in the light of \hat{x}^T featuring a lot of
 structure (see Section 4.2), in the form of an unprecedented amount of structure between x_t, x^0, \hat{x}^T

unaccounted by theory. Failing to realize this implication could lead to incorrect reasoning and spurious discoveries. The same trend can be observed across all three evaluated models.

Finally, we show that the situation is persistent across the training process; see Figure 4. Hypothetically, the value of all three triangle angles shown in Figure 2 could fluctuate during the training process. However, both the angle adjacent to the image $\angle x^0$ and the distance between the image and noise quickly converge to a certain value that remains constant through the rest of the training. This observation brings two main conclusions: (1) The relation between noises, latents, and samples is defined at the early stage of the training, and (2) The inverse DDIM method does not benefit from the prolonged training time of the diffusion model.

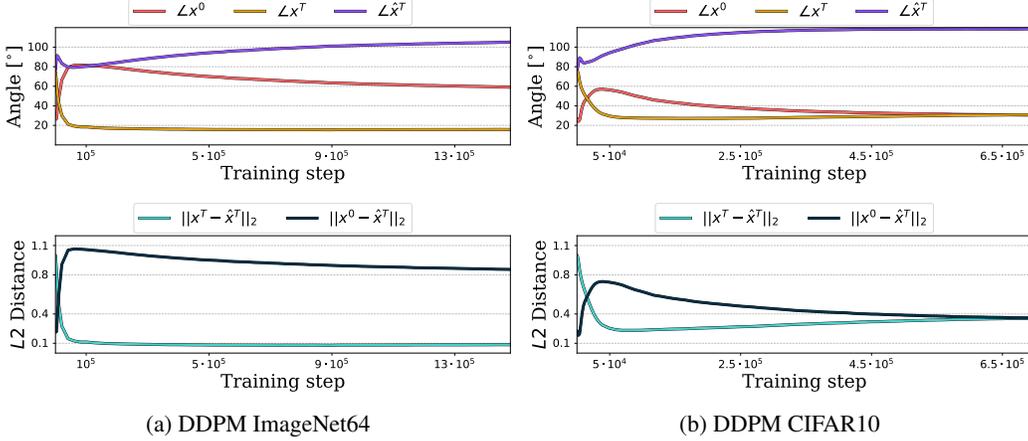


Figure 4: We generate 1000 samples from the final model and revert them to the corresponding latents using an intermediate training checkpoint saved after a given number of training steps.

4.4 NOISE-TO-SAMPLE MAPPING

As indicated by previous works (Kadkhodaie et al., 2023; Zhang et al., 2024), diffusion models converge to the same mapping between the random Gaussian noise (x^T) and the generated images (x^0) independently on the random seed, parts of the training dataset, or even the model architecture. In this work, we investigate this phenomenon further and study the nature of mapping between noises and samples and how it changes during training.

T	CIFAR10 (DDPM)		ImageNet (DDPM)		CelebA (LDM)	
	$x^0 \rightarrow x^T$	$x^T \rightarrow x^0$	$x^0 \rightarrow x^T$	$x^T \rightarrow x^0$	$x^0 \rightarrow x^T$	$x^T \rightarrow x^0$
10	90.3 \pm 6.3	94.0 \pm 2.6	99.4 \pm 0.0	100 \pm 0.0	100 \pm 0.0	100 \pm 0.0
100	98.9 \pm 1.2	50.4 \pm 1.9	100 \pm 0.0	59.0 \pm 7.1	100 \pm 0.0	100 \pm 0.0
1000	99.1 \pm 1.0	46.8 \pm 3.0	99.8 \pm 0.2	44.6 \pm 6.3	100 \pm 0.0	100 \pm 0.0
4000	99.1 \pm 1.0	46.4 \pm 3.0	99.5 \pm 0.3	43.3 \pm 6.7	-	-

Table 2: Accuracy on assigning noise to the corresponding generated image ($x^0 \rightarrow x^T$) and vice-versa ($x^T \rightarrow x^0$). In DDPMs, we are able to correctly select the original noise for a given sample by calculating the L2 distances, while the reverse assignment is only valid for a small number of diffusion steps. In LDM, we can correctly predict assignments in both directions.

To that end, we first generate 1K samples from random Gaussian noises and try to predict which noise (x^T) was used to generate which sample (x^0) and vice-versa. As presented in Table 2, we show that we can accurately assign the images to the corresponding noises ($x^0 \rightarrow x^T$) according to the smallest L2 distance criterion. This is especially true for the higher number of diffusion timesteps, where for all models, we achieve over 99% accuracy. The situation changes when trying to assign the noise to the corresponding generated image ($x^T \rightarrow x^0$). We can observe high accuracy with a low number of generation timesteps ($T = 10$), but the results deteriorate quickly with the

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

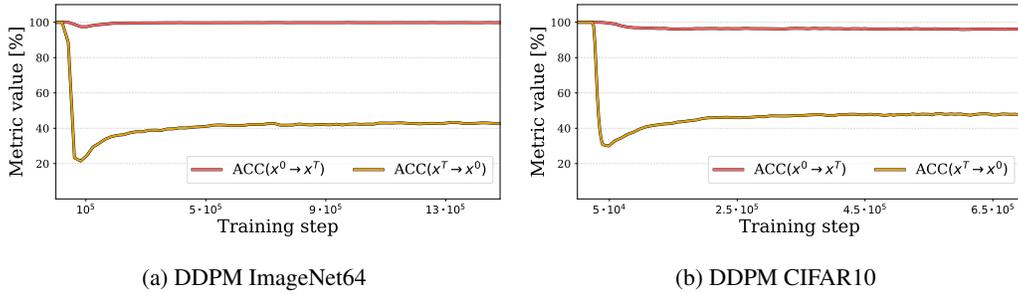


Figure 5: Accuracy of assigning initial noise given the generated sample ($\mathbf{x}^0 \rightarrow \mathbf{x}^T$) and sample given the initial noise ($\mathbf{x}^T \rightarrow \mathbf{x}^0$) when training the diffusion model. We can observe that from the very beginning of training, we can assign initial noise with a simple L2 distance, while the accuracy of the reverse assignment rapidly drops.

increase of this parameter. The reason for this is that greater values of T allow the generation of a broader range of images, including the ones with large plain areas of low pixel variance. Such generations turn out to be close to the majority of initial Gaussian noises. We provide examples of such generations in the Appendix A.1. Notably, we do not observe such behavior in the LDM model, where final generations are well normalized through the KL-divergence applied to the latent space of the LDM’s autoencoder.

To further analyze the nature of this property, we show how those metrics change when sampling generations from intermediate checkpoints of the diffusion models’ training; see Figure 5. We can observe that the distance between noises and latents accurately defines the assignment of initial noises given the generated samples ($\mathbf{x}^0 \rightarrow \mathbf{x}^T$) from the beginning of the training till the end. At the same time, the accurate reverse assignment ($\mathbf{x}^T \rightarrow \mathbf{x}^0$) can only be observed at the beginning of the training, when generations have not yet been properly formed.

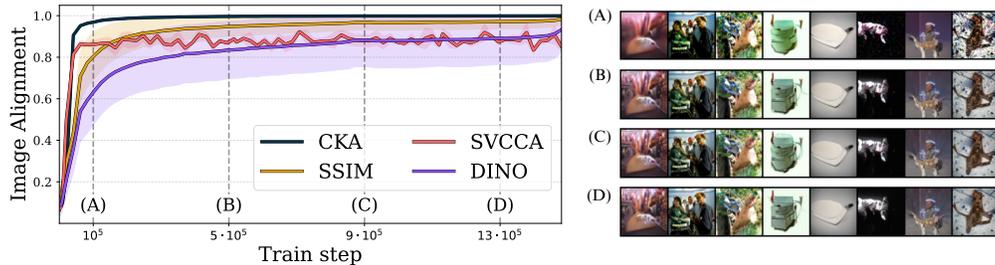


Figure 6: Similarity of the generations sampled from the same random noise at different stages of diffusion model’s training to the final outputs for DDPM and ImageNet. Already after a few epochs, model learns the mapping between Gaussian noise and generations. Prolonged training improves the quality of samples, adding high-frequency features, without changing their content. This can be observed through different image alignment metrics (left) and visual inspection (right).

To further measure how the relation between noises and samples change when training the model, for each training step $n \in \{1 \dots 500K\}$ for CIFAR10 and $\{1 \dots 1.5M\}$ for ImageNet we generate 1K samples $\{\mathbf{x}_{i,n}^0\}_{i=1}^{1000}$ from the same random noise $\mathbf{x}_{\text{fixed}}^T \sim N(0, I)$, and compare them with generations obtained for the fully trained model. We present the visualization of this comparison in Figure 6 using CKA, DINO, SSIM, and SVCCA as metrics. We notice that image features rapidly converge to the level that persists until the end of the training. This means that prolonged learning does not significantly alter how the data is assigned to the Gaussian noise after the early stage of the training. It is especially visible when considering the SVCCA metric, which measures the average correlation of top-10 correlated data features between two sets of samples.

We can observe that this quantity is high and stable through training, showing that generating the most important image concepts from a given noise will not be affected by a longer learning pro-

378 cess. For visual comparison, we plot the generations sampled from the model trained with different
379 numbers of training steps in Figure 6 (right).
380

381 5 CONCLUSIONS 382

383 In this work, we empirically study the relation between initial Gaussian noise, generated samples,
384 and their latent representations calculated with the inverse DDIM technique. First, we show that the
385 error in the approximation of the previous noise in DDIM leads to representations located next to the
386 generation trajectory between starting samples and their true noises. Moreover, prolonged diffusion
387 training does not affect this property, as the accuracy of DDIM inversion does not improve in time.
388 Then, studying the relation between the generated samples and Gaussian noise, we show that we
389 can accurately assign the initial noise of the given generation with a simple L2 distance. We also
390 demonstrate that this behavior emerges at the very beginning of the diffusion models training. Our
391 experiments lead to the conclusion that the initial part of the diffusion model’s training is responsible
392 for building the relation between the initial Gaussian noise, final generations, and their inverted
393 representations.
394

395 REFERENCES 396

- 397 Dmitry Baranchuk, Ivan Rubachev, A. Voynov, Valentin Khulkov, and Artem Babenko. Label-
398 efficient semantic segmentation with diffusion models. *International Conference on Learning*
399 *Representations*, 2021.
- 400 Duygu Ceylan, Chun-Hao P. Huang, and Niloy J. Mitra. Pix2video: Video edit-
401 ing using image diffusion. In *Proceedings of the IEEE/CVF International Confer-*
402 *ence on Computer Vision (ICCV)*, pp. 23206–23217, October 2023. URL https://openaccess.thecvf.com/content/ICCV2023/html/Ceylan_Pix2Video_Video_Editing_using_Image_Diffusion_ICCV_2023_paper.html.
- 406 Kamil Deja, Georgi Tinchev, Marta Czarnowska, Marius Cotescu, and Jasha
407 Droppo. Diffusion-based accent modelling in speech synthesis. In *Interspeech*
408 *2023*, 2023a. URL <https://www.amazon.science/publications/diffusion-based-accent-modelling-in-speech-synthesis>.
- 410 Kamil Deja, Tomasz Trzcinski, and Jakub M Tomczak. Learning data representations with joint
411 diffusion models. *arXiv preprint arXiv:2301.13622*, 2023b.
- 413 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *Ad-*
414 *vances in Neural Information Processing Systems*, 34, 2021.
- 415 Daniel Garibi, Or Patashnik, Andrey Voynov, Hadar Averbuch-Elor, and Daniel Cohen-Or. Renoise:
416 Real image inversion through iterative noising. *arXiv preprint arXiv: 2403.14602*, 2024. URL
417 <https://arxiv.org/abs/2403.14602v1>.
- 419 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
420 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Infor-*
421 *mation Processing Systems*, pp. 2672–2680, 2014.
- 422 Amir Hertz, Ron Mokady, J. Tenenbaum, Kfir Aberman, Y. Pritch, and D. Cohen-Or. Prompt-to-
423 prompt image editing with cross attention control. *International Conference on Learning Repr-*
424 *esentations*, 2022. doi: 10.48550/arXiv.2208.01626.
- 426 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on*
427 *Deep Generative Models and Downstream Applications*, 2021. URL <https://openreview.net/forum?id=qw8AKxfYbI>.
- 430 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
431 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

-
- 432 Seongmin Hong, Kyeonghyun Lee, Suh Yoon Jeon, Hyewon Bae, and Se Young Chun. On exact
433 inversion of dpm-solvers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
434 *Pattern Recognition*, pp. 7069–7078, 2024.
- 435 Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddp noise
436 space: Inversion and manipulations. In *Proceedings of the IEEE/CVF Conference on Computer*
437 *Vision and Pattern Recognition*, pp. 12469–12478, 2024.
- 438 Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization
439 in diffusion models arises from geometry-adaptive harmonic representation. *arXiv preprint*
440 *arXiv:2310.02557*, 2023.
- 441 Zahra Kadkhodaie, Florentin Guth, Eero P. Simoncelli, and Stéphane Mallat. Generalization in
442 diffusion models arises from geometry-adaptive harmonic representations. In *The Twelfth Inter-*
443 *national Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.*
444 OpenReview.net, 2024. URL <https://openreview.net/forum?id=ANvmVS2Yr0>.
- 445 Valentin Khruikov and I. Oseledets. Understanding ddp latent codes through optimal transport.
446 *International Conference on Learning Representations*, 2022. URL <https://arxiv.org/abs/2202.07477v2>.
- 447 Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models
448 for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
449 *and Pattern Recognition (CVPR)*, pp. 2426–2435, June 2022.
- 450 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Ad-*
451 *vances in neural information processing systems*, 34:21696–21707, 2021.
- 452 Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *International Confer-*
453 *ence on Learning Representations*, 2014.
- 454 Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 con-
455 volutions. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
456 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran
457 Associates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/](https://proceedings.neurips.cc/paper_files/paper/2018/file/d139db6a236200b21cc7f752979132d0-Paper.pdf)
458 [paper/2018/file/d139db6a236200b21cc7f752979132d0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/d139db6a236200b21cc7f752979132d0-Paper.pdf).
- 459 Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent
460 space. *arXiv preprint arXiv:2210.10960*, 2022.
- 461 Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, and Zhou Zhao. Diffsinger: Singing voice synthesis
462 via shallow diffusion mechanism. *AAAI Conference on Artificial Intelligence*, 2021. doi: 10.
463 1609/aaai.v36i10.21350.
- 464 Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. Diffusion
465 hyperfeatures: Searching through time and space for semantic correspondence. *Advances in*
466 *Neural Information Processing Systems*, 36, 2024.
- 467 Barak Meiri, Dvir Samuel, Nir Darshan, Gal Chechik, Shai Avidan, and Rami Ben-Ari. Fixed-point
468 inversion for text-to-image diffusion models. *arXiv preprint arXiv:2312.12540*, 2023.
- 469 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion
470 for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF*
471 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6038–6047, June
472 2023. URL [https://openaccess.thecvf.com/content/CVPR2023/html/](https://openaccess.thecvf.com/content/CVPR2023/html/Mokady_NULL-Text_Inversion_for_Editing_Real_Images_Using_Guided_Diffusion_Models_CVPR_2023_paper.html)
473 [Mokady_NULL-Text_Inversion_for_Editing_Real_Images_Using_Guided_](https://openaccess.thecvf.com/content/CVPR2023/html/Mokady_NULL-Text_Inversion_for_Editing_Real_Images_Using_Guided_Diffusion_Models_CVPR_2023_paper.html)
474 [Diffusion_Models_CVPR_2023_paper.html](https://openaccess.thecvf.com/content/CVPR2023/html/Mokady_NULL-Text_Inversion_for_Editing_Real_Images_Using_Guided_Diffusion_Models_CVPR_2023_paper.html).
- 475 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
476 In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.
- 477 Yong-Hyun Park, Mingi Kwon, Jaewoong Choi, Junghyo Jo, and Youngjung Uh. Understanding the
478 latent space of diffusion models through the lens of riemannian geometry. *Advances in Neural*
479 *Information Processing Systems*, 36:24129–24142, 2023.

-
- 486 Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu.
487 Zero-shot image-to-image translation. In *ACM SIGGRAPH 2023 Conference Proceedings*, pp.
488 1–11, 2023.
- 489 Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-
490 tts: A diffusion probabilistic model for text-to-speech. In *International Conference on Machine*
491 *Learning*, pp. 8599–8608. PMLR, 2021.
- 493 Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Dif-
494 fusion autoencoders: Toward a meaningful and decodable representation. *Computer Vision And*
495 *Pattern Recognition*, 2021. doi: 10.1109/CVPR52688.2022.01036. URL <https://arxiv.org/abs/2111.15640v3>.
- 497 Margherita Rosnati, Mélanie Roschewitz, and Ben Glocker. Robust semi-supervised segmentation
498 with timestep ensembling diffusion models. In *Machine Learning for Health (ML4H)*, pp. 512–
499 527. PMLR, 2023.
- 501 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
502 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learn-*
503 *ing*, pp. 2256–2265. PMLR, 2015.
- 504 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
505 *preprint arXiv:2010.02502*, 2020.
- 507 Xu Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. Dual diffusion implicit bridges for image-
508 to-image translation. *International Conference on Learning Representations*, 2022. doi: 10.
509 48550/arXiv.2203.08382.
- 510 Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for
511 text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Com-*
512 *puter Vision and Pattern Recognition*, pp. 1921–1930, 2023.
- 513 Guanhua Zhang, Jiabao Ji, Yang Zhang, Mo Yu, Tommi S. Jaakkola, and Shiyu Chang. Towards
514 coherent image inpainting using denoising diffusion implicit models. In Andreas Krause, Emma
515 Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *In-*
516 *ternational Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii,*
517 *USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 41164–41193. PMLR, 2023.
518 URL <https://proceedings.mlr.press/v202/zhang23q.html>.
- 519 Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Liyue Shen, and Qing Qu. The emergence of
520 reproducibility and consistency in diffusion models, 2024. URL <https://openreview.net/forum?id=UkLSvLqi07>.

524 A APPENDIX

526 A.1 MISTAKES IN NOISE TO SAMPLE MAPPING

527
528 In Figures 7 and 8 we show examples of images that are not properly assigned to their initial noises.
529 We can observe that those images are characterised by low variance of pixels.
530
531
532
533
534
535
536
537
538
539

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

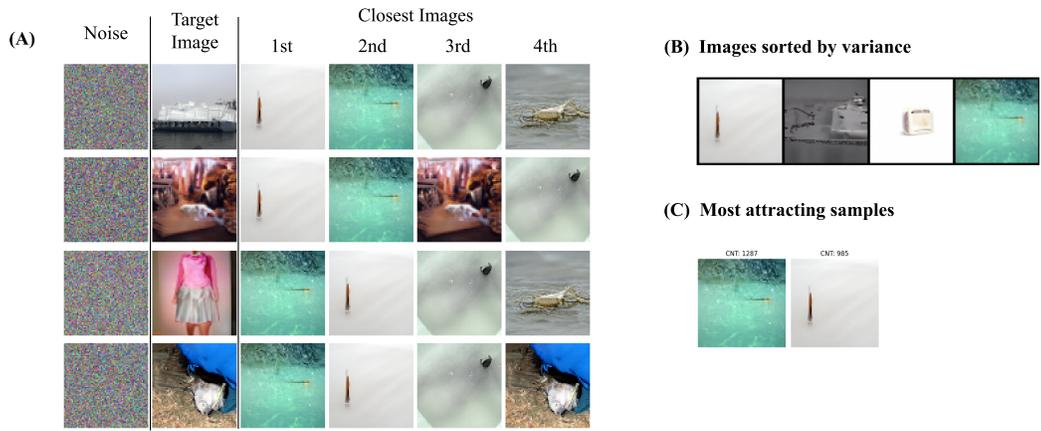


Figure 7: Examples assigned to the wrong initial noises for CIFAR10 datasets

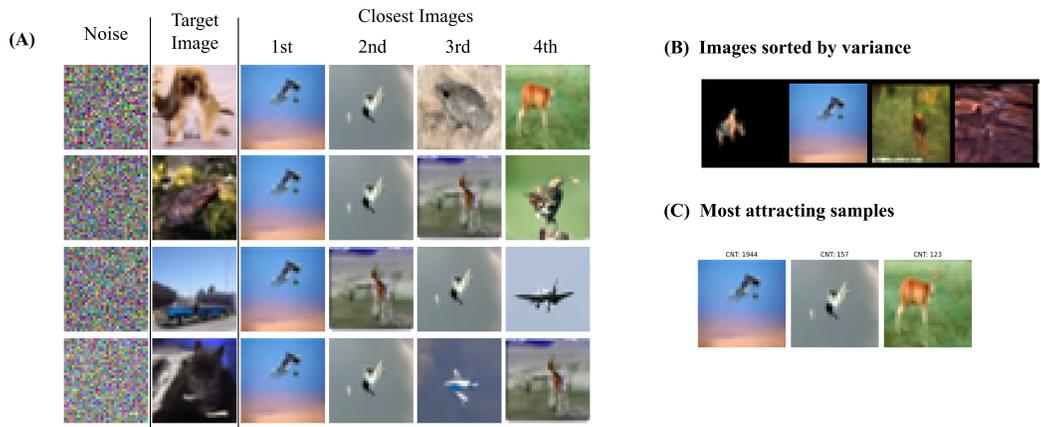


Figure 8: Examples assigned to the wrong initial noises for CIFAR10 datasets