

# Unrevealing Hidden Relations Between Latent Space and Image Generations in Diffusion Models

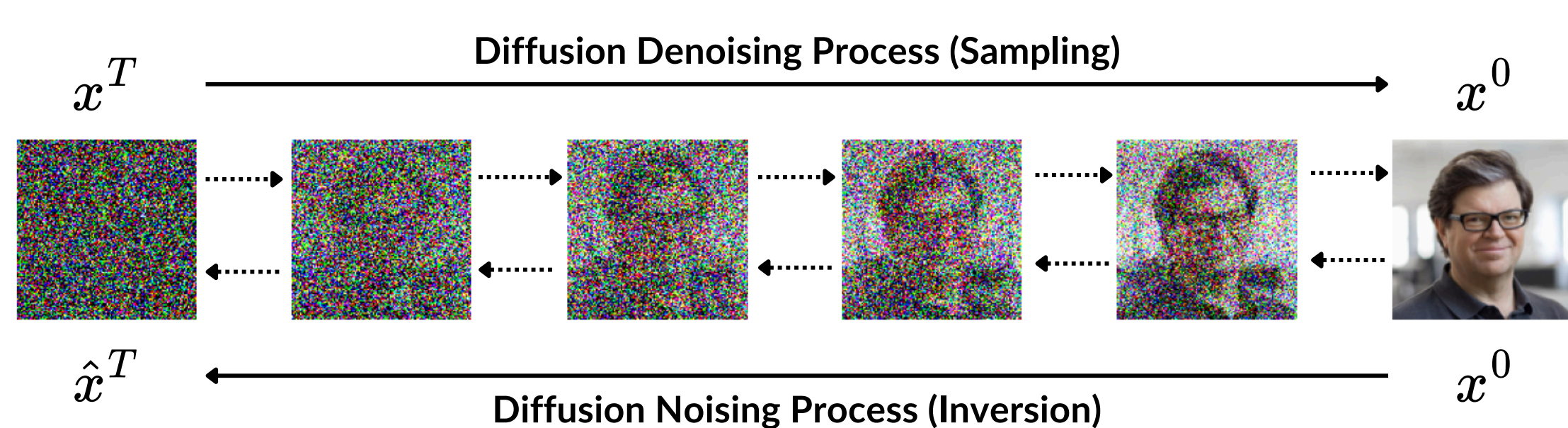
Łukasz Staniszewski<sup>1</sup> Łukasz Kuciński<sup>2,3,4</sup> Kamil Deja<sup>1,2</sup>

<sup>1</sup>Warsaw University of Technology <sup>2</sup>IDEAS NCBR <sup>3</sup>University of Warsaw <sup>4</sup>Polish Academy of Sciences

## #TLDR

- We study relations between Gaussian noises  $\mathbf{x}^T$ , image samples  $\mathbf{x}^0$  and their latent encodings  $\hat{\mathbf{x}}^T$  from the DDIM inversion procedure.
- We show that those encodings  $\hat{\mathbf{x}}^T$  manifold is between initial noise  $\mathbf{x}^T$  and image generations  $\mathbf{x}^0$ .
- We show that noise  $\mathbf{x}^T$  to image  $\mathbf{x}^0$  mapping may be defined using the smallest  $L_2$  distance and that DMs learn important image features at the beginning of the fine-tuning.

## Background



We can inverse the standard diffusion denoising procedure into the noising procedure:

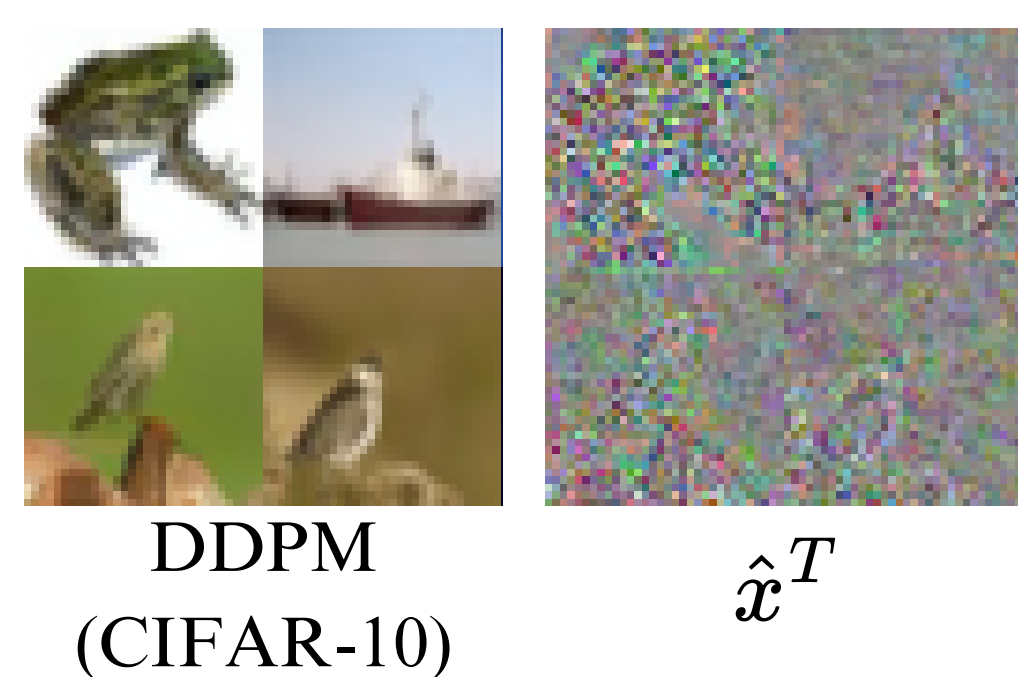
$$x_t = \gamma \cdot x_{t-1} + \eta \cdot \epsilon_{\theta}(x_t, t, c)$$

Due to circular dependency on  $\epsilon_{\theta}(x_t, t, c)$ , DDIM inversion approximates it:

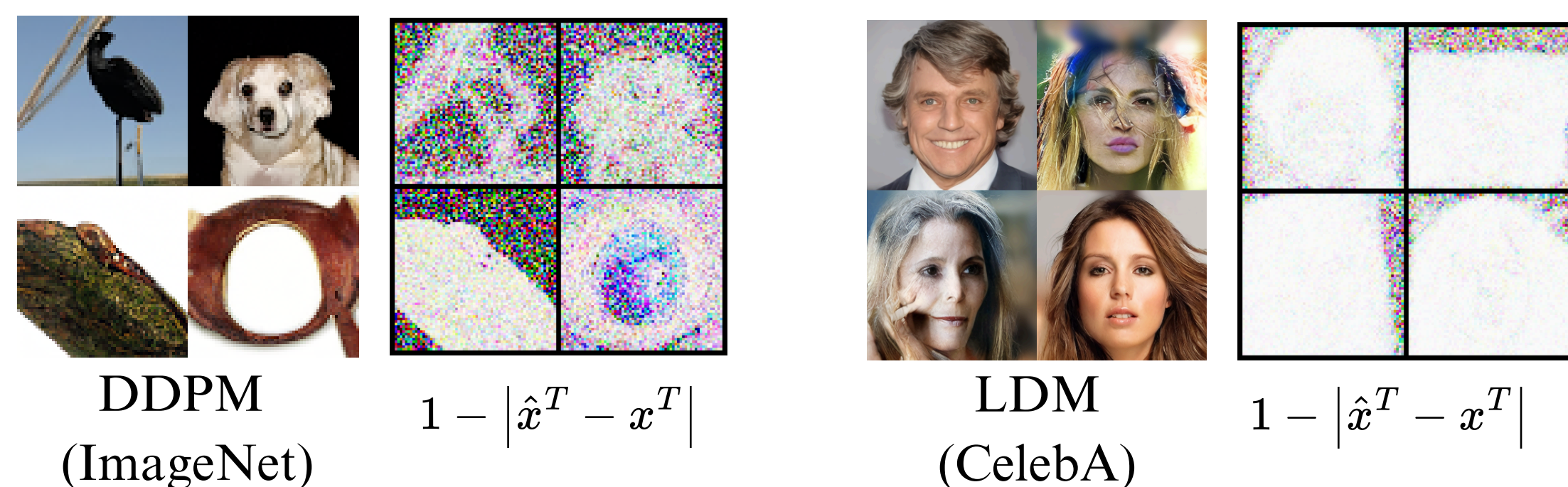
$$\epsilon_{\theta}(x_t, t, c) \approx \epsilon_{\theta}(\mathbf{x}_{t-1}, t, c).$$

## Latent $\neq$ Noise

We can observe clear structures of original images  $\mathbf{x}^0$  in the inverted latents  $\hat{\mathbf{x}}^T$ ...



...or by showing the image difference between the latent  $\hat{\mathbf{x}}^T$  and the noise  $\mathbf{x}^T$ .



Latent encodings  $\hat{\mathbf{x}}^T$  have correlated pixels.

	DDPM (CIFAR-10)	DDPM (ImageNet)	LDM (CelebA)
Noise ( $\mathbf{x}^T$ )	0.159 $\pm$ 0.003	0.177 $\pm$ 0.007	
Latent ( $\hat{\mathbf{x}}^T$ )	0.462 $\pm$ 0.009	0.219 $\pm$ 0.006	0.179 $\pm$ 0.008
Sample ( $\mathbf{x}^0$ )	0.986 $\pm$ 0.001	0.966 $\pm$ 0.001	0.904 $\pm$ 0.005

Table 1. Top-10 correlation coefficients in random Gaussian noise vs. latent encoding.

## If you enjoy this work...

See the full paper for more details!



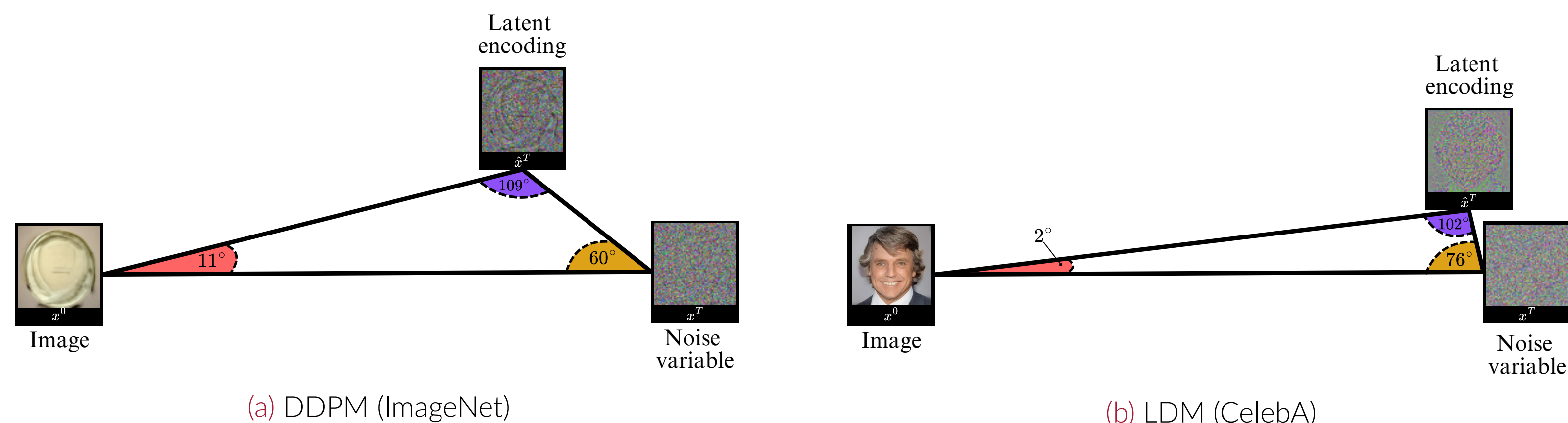
arXiv paper



Łukasz Staniszewski

## Where are the latents $\hat{\mathbf{x}}^T$ located?

Latent encodings ( $\hat{\mathbf{x}}^T$ ) manifold is between random Gaussian noises ( $\mathbf{x}^T$ ) and their corresponding samples ( $\mathbf{x}^0$ ) manifolds.



Diffusion model denoising trajectory is aligned with linear interpolation path between the Gaussian noise  $\mathbf{x}^T$  and latent encoding  $\hat{\mathbf{x}}^T$ .

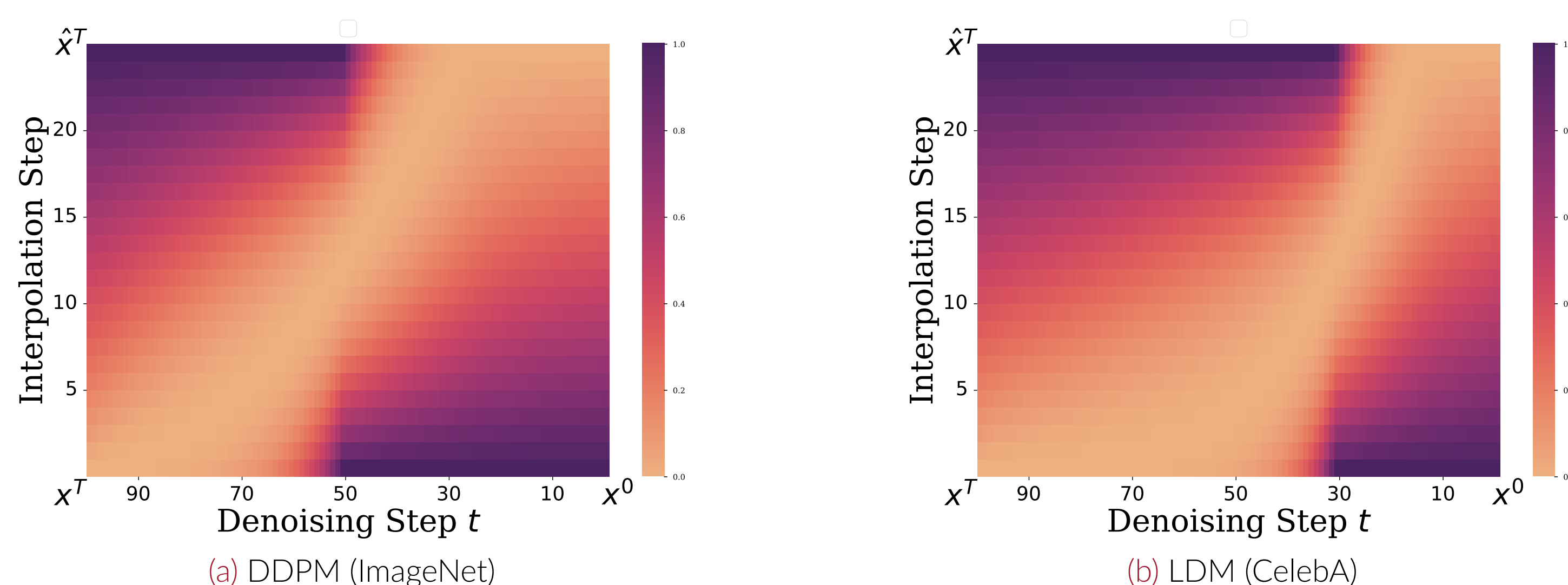


Figure 5. Distances between next denoising steps and the  $\mathbf{x}^T \rightarrow \hat{\mathbf{x}}^T$  interpolation points. Intermediate generations along the sampling trajectory initially get closer to the latent variable, and after approximately 50-70% of the path, they pass the latent.

## Noise-to-Sample mapping

The mapping between initial Gaussian noise  $\mathbf{x}^T$  and its corresponding generation  $\mathbf{x}^0$  is secretly a  $L_2$ -based nearest neighbor mapping.

T	ImageNet (DDPM)		CelebA (LDM)	
	$x^0 \rightarrow x^T$	$x^T \rightarrow x^0$	$x^0 \rightarrow x^T$	$x^T \rightarrow x^0$
10	99.4 $\pm$ 0.0	100 $\pm$ 0.0	100 $\pm$ 0.0	100 $\pm$ 0.0
100	100 $\pm$ 0.0	59.0 $\pm$ 7.1	100 $\pm$ 0.0	100 $\pm$ 0.0
1000	99.8 $\pm$ 0.2	44.6 $\pm$ 6.3	100 $\pm$ 0.0	100 $\pm$ 0.0
4000	99.5 $\pm$ 0.3	43.3 $\pm$ 6.7	-	-

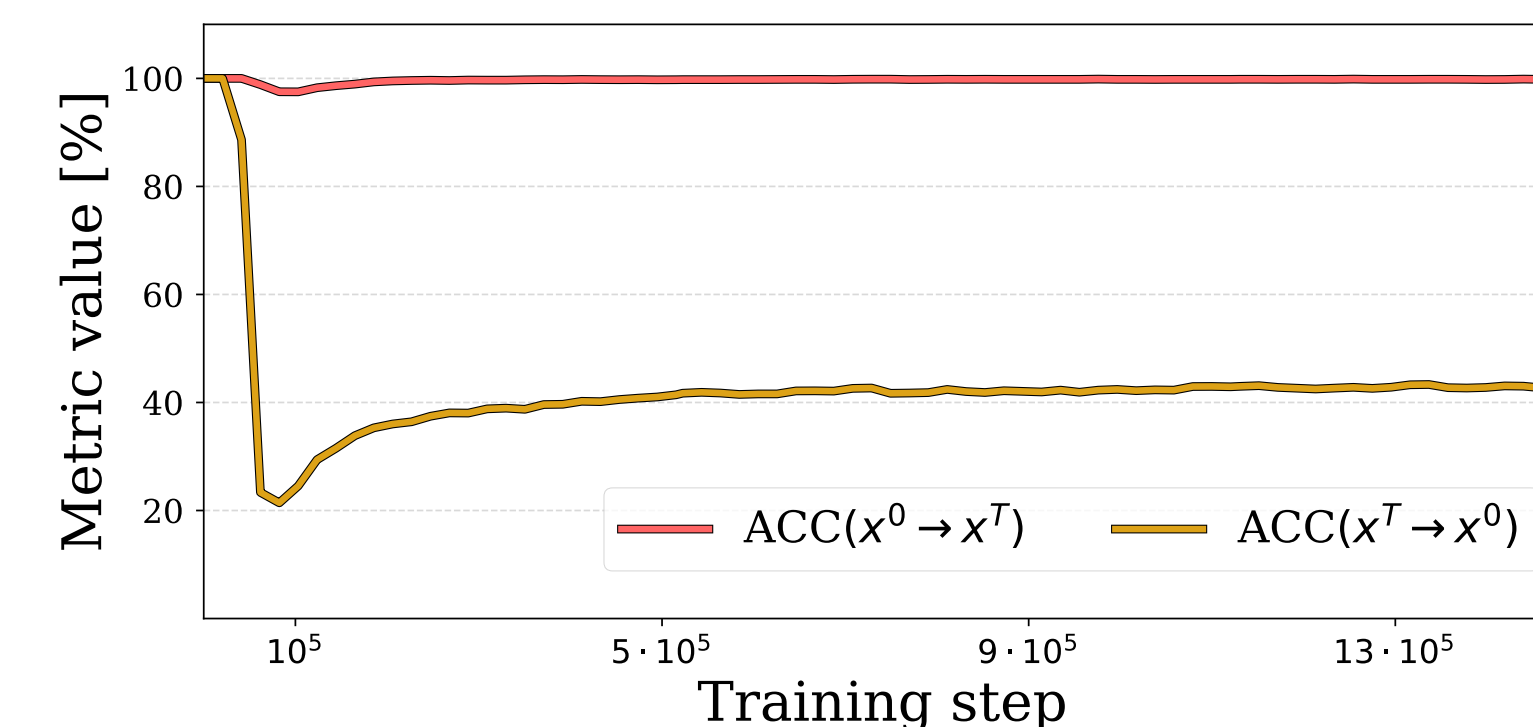


Figure 6. We are able to correctly select the original noise ( $\mathbf{x}^T$ ) for a given sample ( $\mathbf{x}^0$ ) by indicating the one with the closest  $L_2$  distance (left). Moreover, we show (right) that this mapping is established at the beginning of fine-tuning.

DMs generate the most important image features right at the beginning of fine-tuning, with only small details added further.

