

## Appendix: On Inductive Biases That Enable Generalization of Diffusion Transformers

This appendix is structured as follows: In Appendix A, we detail the structure of DiTs, training and sampling settings, as well as the PSNR computation. In Appendix B, we theoretically demonstrate that the attention locality correlates with a DiT’s generalization by connecting the discovered attention locality bias with existing inductive biases. In Appendix C, we provide additional quantitative results to evaluate both the generalization and generation quality of DiTs. For generalization, we report the PSNR gap with  $10^3$  training samples and introduce a new evaluation metric: cosine similarity to the nearest training image. For generation quality, we present results from both pixel-space and latent-space DiTs, evaluated using additional metrics including Inception Score and FD-DINOv2. In Appendix D, we qualitatively compare the generated images of a DiT with and without modifying its generalization. Appendix E analyzes how the input noise level influences the attention locality. Appendix F, Appendix G, and Appendix H further verify the existence and consistency of the attention map locality for a DiT, as well as the fact that a DiT lacks harmonic bases. These appendices provide additional support to our observations and analysis in the main paper. Appendix I presents the effects of reducing parameters as a complementary approach to reducing FLOPs. Appendix J discusses the limitations of this work, while Appendix K outlines its broader impact.

### A Experimental Settings

#### A.1 Model Definition

We verify the effectiveness of local attention in modifying the generalization of a DiT using 2 DiT backbones and 10 local attention variations. Tab. 6 provides more details about the DiT backbones and local attention configurations.

Table 6: DiT architectures and local attention settings. In the column titled ‘DiT Blocks’,  $G$  denotes global attention while a number  $k$  represents a local attention of window size  $k \times k$ .

Model	DiT Blocks	Hidden Size	Num Heads	Depth	Patch Size
DiT-XS/1	$(G, G, G, G, G, G, G, G, G, G, G)$	252	4	12	$1 \times 1$
w/ Local	$(3, 5, 7, 9, 11, 13, G, G, G, G, G)$	252	4	12	$1 \times 1$
w/ Local (mix)	$(3, G, 5, G, 7, G, 9, G, 11, G, 13, G)$	252	4	12	$1 \times 1$
w/ Local (tail)	$(G, G, G, G, G, G, 3, 5, 7, 9, 11, 13)$	252	4	12	$1 \times 1$
w/ Local (smaller win)	$(3, 5, 7, 9, 11, 13, 15, 17, 19, G, G, G)$	252	4	12	$1 \times 1$
w/ Local*	$(3, 3, 3, 5, 5, 5, 7, 7, 7, G, G, G)$	252	4	12	$1 \times 1$
w/ Local* (mix)	$(3, 3, 3, G, 5, 5, 5, G, 7, 7, 7, G)$	252	4	12	$1 \times 1$
w/ Local* (tail)	$(G, G, G, 3, 3, 3, 5, 5, 5, 7, 7, 7)$	252	4	12	$1 \times 1$
w/ Local* (larger win)	$(5, 5, 5, 7, 7, 7, 9, 9, 9, G, G, G)$	252	4	12	$1 \times 1$
w/ Local Attn ( $5^{*6}$ )	$(5, 5, 5, 5, 5, 5, G, G, G, G, G, G)$	252	4	12	$1 \times 1$
w/ Local Attn ( $3^{*2}, 5^{*2}, 7^{*2}$ )	$(3, 3, 5, 5, 7, 7, G, G, G, G, G, G)$	252	4	12	$1 \times 1$
DiT-S/1	$(G, G, G, G, G, G, G, G, G, G, G)$	384	6	12	$1 \times 1$
w/ Local	$(3, 5, 7, 9, 11, 13, G, G, G, G, G)$	384	6	12	$1 \times 1$
DiT-XXS/1	$(G, G, G, G, G, G, G, G, G, G, G)$	240	4	12	$1 \times 1$

To elaborate, we adopt two DiT backbones, DiT-XS/1 and DiT-S/1, to verify the effectiveness of attention window restrictions in modifying the generalization of a DiT. Both models have 12 DiT Blocks. We remove the auto-encoder and use a patch size of  $1 \times 1$ . Note, DiT-XS/1 has a hidden size of 252 and uses 4 attention heads. In contrast, DiT-S/1 has a hidden size of 384 and uses 6 attention heads. Regarding the local attention variations, the default setting *Local* combines 6 local attentions of window size  $(3, 5, 7, 9, 11, 13)$  and 6 global attentions. Meanwhile, *Local\** is a variant using 9 local attentions of window size  $(3, 3, 3, 5, 5, 5, 7, 7, 7)$ . For both *Local* and *Local\** settings, we place local attentions at the heading layers of a DiT. We also study interleaving the local and global attentions as well as placing local attentions at tailing layers of a DiT, leading to (mix) and (tail) variants. Additionally, to study the effects of modifying the attention window size, we decrease

the attention window size of the *Local* model and increase the attention window size of the *Local\** model, resulting in (smaller win) and (larger win) models in Tab. 6.

## A.2 Training and Sampling Setting

The implementation of the UNet<sup>4</sup> and the DiT<sup>5</sup> are based on the official repositories of Nichol & Dhariwal (2021) and Peebles & Xie (2023). Specifically, for the UNet, we use the architecture which has a 4-stage encoder with channel multipliers of (1, 2, 3, 4). For each stage, we include 3 ResBlocks. At the end of each stage, the resolution of the input tensor are down-sampled by a factor of 2. In the last stage, we use one layer of self-attention. The decoder mirrors the encoder layers and places them in the reverse order, replacing down-sampling layers with up-sampling ones. Between the encoder and decoder, there are 2 ResBlocks and 1 self-attention layer by default. The default skip connections are used between the encoder and decoder at the same resolution. Consequently, the UNet has 303.17G FLOPs and 109.55M parameters. The FLOPs of this UNet are nearly identical to the DiT-XS/1 model in Tab. 6.

All DiT and UNet models are trained with the same hyper-parameter settings. Concretely, we train each model in the pixel-space using a resolution of  $32 \times 32$ . All networks are using the same diffusion algorithm: diffusion steps of 1000 in training and 250 in sampling, and predicting the added noise and sigma simultaneously. To train a network, we use the random seed 43, learning rate  $1e^{-4}$  and an overall batch size of 64. All networks are trained with 8 or 4 A100/H100 GPUs, using the EMA checkpoint at train step 400k with EMA decay 0.9999. For each dataset, we first randomly shuffle the whole dataset. Then we choose the last  $N=10, 10^3, 10^4$  and  $10^5$  images as the training set. The train-test split of a dataset is kept consistent for different networks. When computing FID values for the model trained with  $N=10^4$  and  $10^5$  images, we randomly select  $M=\min\{N, 50k\}$  of the training images as the reference set. Regarding the sampling set, we generate 20k images for the comparably small LSUN Church dataset and 50k images for all other datasets. In Tab. 1 of the main manuscript, we present the results of DiT-XS/1<sup>†</sup> and DiT-S/1<sup>†</sup> models, with and without local attentions. For the <sup>†</sup> models, we use a different dataset shuffling, change the random seed to 143, and double the training batch size to 128. Notably, using local attention on both normal and <sup>†</sup> models can successfully modify the generalization of a DiT, confirming the effectiveness of the locality as the inductive bias of a DiT.

## A.3 PSNR Computation

We compute the PSNR based on a training or testing subset of 300 images following Kadkhodaie et al. (2024). For each image, we re-space the diffusion steps from 1000 to 50, and compute the train and test PSNR on each step. Specifically, we first perform the noising step of the diffusion model to get the noisy image at a diffusion step  $t$ . Next, we feed the noisy image into the diffusion model backbone and get the estimation of the added noise, which is then used to recover the clean image from the training or testing subset, *i.e.*, performing a one-step denoising. The final PSNR at step  $t$  is obtained using the estimated clean image and the ground truth. Consequently, the PSNR value can estimate a diffusion model’s accuracy at each diffusion step. Therefore, the PSNR gap between the training and testing subsets can measure a diffusion model’s generalization: when a diffusion model has good generalization, its prediction accuracy should be comparably between the training and testing set, resulting in a small PSNR gap.

## B Connection to Theoretical Results

In this section we’ll provide connections to theoretical work (De Wolf, 2008; Yang & Salman, 2019; Vasudeva et al., 2024) that can be used to explain our empirical findings.

We will start by discussing preliminaries from prior work on the simplicity bias of transformers (Vasudeva et al., 2024) in Appendix B.1. We’ll subsequently connect this work to our results in Appendix B.2 and show that local attention encourages the low sensitivity bias of a transformer. Finally, in Appendix B.3, we demonstrate low sensitivity of a transformer is connected to better generalization.

<sup>4</sup><https://github.com/openai/improved-diffusion>  
<sup>5</sup><https://github.com/facebookresearch/DiT>



## B.1 Preliminaries

Prior work (Vasudeva et al., 2024) showed that attention modules learn simpler features more quickly, which implies that the transformer is biased towards simple functions and lower sensitivity. To show this, Vasudeva et al. (2024) considered a model with at least one self-attention layer. To simplify the analysis, prior work removes the non-linear Softmax function from a standard self-attention layer and focuses on linear attention of the form

$$\Phi = \frac{\mathbf{x}W_q \cdot W_k^\top \mathbf{x}^\top}{\sqrt{\dim}} \cdot \mathbf{x}W_v, \quad (8)$$

with input  $\mathbf{x} \in \mathbb{R}^{T \times \tilde{d}}$  and  $\dim$  the scaling dimension of the attention layer. Further,  $W_q$ ,  $W_k$ , and  $W_v$  are trainable parameters that map the input  $\mathbf{x}$  to query, key, and value, respectively. Below, we use  $d = T\tilde{d}$ .

Under the assumption that a transformer with linear self-attention layers works in a boolean space  $\{0, 1\}^d$ , Vasudeva et al. (2024) showed the following main results: A transformer model  $f(\mathbf{x})$  that contains at least one self-attention layer can be represented by the linear combination of a set of orthonormal monomial terms

$$f(\mathbf{x}) = \sum_{U \subseteq [d]} \hat{f}(U) \chi_U(\mathbf{x}), \quad \chi_U := \prod_{i \in U} x_i, \quad \forall U \subseteq [d], \quad (9)$$

where set  $U \subseteq [d] = \{1, \dots, d\}$ , and term  $\hat{f}(U)$  is the coefficient for a monomial term. For an input sequence  $\mathbf{x}$ , these orthonormal monomial terms, under their specific assumptions (Yang & Salman 2019), form a set of Fourier bases (De Wolf 2008).

In addition, at an input location  $\mathbf{x}$ , Vasudeva et al. (2024) compute the eigenvalues of those orthonormal monomial terms  $\chi_U$ , which form an eigenfunction, via

$$\mu_{|U|} := \mathbb{E}_{\mathbf{x} \sim \{0,1\}^d} [\chi_U K(\mathbf{x}, \mathbf{1})]. \quad (10)$$

In Eq. (10),  $\mu_{|U|}$  is the eigenvalue for monomial  $\chi_U$ ,  $|U|$  denotes the size of  $U$ ,  $\mathbf{1}$  is a vector of all ones in space  $\{0, 1\}^d$ , and  $K(\cdot)$  represents a neural kernel (Yang & Salman, 2019; Hron et al., 2020), e.g., conjugate kernel (CK) or neural tangent kernel (NTK). Based on theorems discussed in prior work (Yang & Salman 2019; Hron et al., 2020), Vasudeva et al. (2024) theoretically prove that eigenvalues  $\mu_{|U|}$  for  $U \subseteq [d]$  satisfy

$$\begin{aligned} \mu_0 &\geq \mu_2 \geq \dots \geq \mu_{2k} \geq \dots, \\ \mu_1 &\geq \mu_3 \geq \dots \geq \mu_{2k+1} \geq \dots \end{aligned} \quad (11)$$

This result is important because it explains why attention modules learn simpler features more quickly: the eigenvalues of monomial terms with lower degree are larger as shown in Eq. (11). This indicates that transformers are biased towards polynomials with lower orders. Considering that a low-degree polynomial tends to have low sensitivity, this result also implies that transformers are biased toward low sensitivity functions.

## B.2 Relation to Inductive Biases in Diffusion Transformers

This result is relevant because it provides a theoretical foundation for our work. Concretely, when using global attention,  $U \subseteq [d]$  is not restricted in any form. This hence means that any elements in the input tensor  $\mathbf{x} \in \mathbb{R}^{T \times \tilde{d}}$  can interact with each other.

In contrast, using local attention restricts the interaction between elements in the input tensor as illustrated in Figs. 4 and 5 of our main paper. This implies that  $U$  now only represents a subset of the possible interactions, which reduces the order of the highest degree monomial  $\chi_U$  significantly.

Because the highest degree monomials are of much lower order, local attention lowers the sensitivity of the transformer *w.r.t.* data perturbations.

Table 7: PSNR gap↓ comparison with  $10^3$  training data between a DiT with and without local attention for two architectures: DiT-XS/1 and DiT-S/1. *Local* denotes applying local attention with window sizes (3, 5, 7, 9, 11, 13) to the first six layers of the DiT.

Dataset	CelebA	ImageNet	MSCOCO	LSUN Church	LSUN Bedroom	LSUN Bridge	LSUN Tower
DiT-XS/1	7.49	7.77	7.36	6.76	7.45	6.79	7.39
w/ Local	6.56 −0.12	6.76 −0.13	6.36 −1.00	5.71 −0.16	6.20 −0.17	5.77 −0.15	6.15 −0.17
DiT-S/1	9.89	8.71	8.86	8.45	8.98	9.39	9.78
w/ Local	8.76 −0.11	7.22 −0.17	7.35 −1.51	7.96 −0.06	7.50 −0.16	7.95 −0.15	8.42 −0.14

### B.3 Lower Sensitivity Leads to Better Generalization

Under the linear self-attention assumption, Vasudeva et al. (2024) also demonstrate that sensitivity of a transformer of data perturbation is connected with the sharpness of the minima, *i.e.*, the sensitivity of the loss for small changes of the network weight near minima of the parameter space. The low sharpness of the minima is a widely accepted indicator of model generalization (Keskar et al., 2016; Neyshabur et al., 2017; Jiang et al., 2019) and has been empirically verified for transformers (Hahn & Rofin, 2024). Considering a linear model  $\Phi(\theta; x) = \theta^\top x$ , where  $\theta$  is the weight of the linear layer and  $x$  is the input data. Adding a small perturbation  $\Delta x$  to input  $x$  is equivalent to perturbing the layer weight  $\theta$

$$\Phi(\theta; x + \Delta x) = \theta^\top (x + \Delta x) = \Phi(\theta; x) + \theta^\top \Delta x = \Phi(\theta; x) + \Delta \theta^\top x = \Phi(\theta + \Delta \theta; x), \quad (12)$$

where  $\Delta \theta = \frac{\theta^\top \Delta x}{\|x\|_2^2} x$ .

For a more complex model like a transformer, Vasudeva et al. (2024) empirically verified that the connection between the low sensitivity and flat minima still holds. Taking both Eq. (11) and Eq. (12) into consideration, we can draw the conclusion that using local attention reduces the sensitivity of a transformer, resulting in flatter minima, which leads to improved generalization.

## C Additional Quantitative Results

Here, we first present additional quantitative results on diffusion model generalization, including PSNR gap evaluations with  $10^3$  training samples and an alternative generalization metric based on cosine similarity to the nearest training image. Next, we provide results using additional image quality metrics (Inception Score and FD-DINOv2), as well as evaluations on UNet, various DiT backbones, and latent-space diffusion models.

### C.1 PSNR Gap Comparisons with $10^3$ Training Images

In Tab. 2 of the main paper, we compare PSNR gaps using  $10^4$  and  $10^5$  training images across seven datasets. Additionally, Tab. 7 reports results with  $10^3$  training images on all evaluated datasets. Tab. 7 demonstrates that encouraging the attention locality can reduce the PSNR gap of models trained with  $10^3$  images. These results show that encouraging attention locality consistently reduces the PSNR gap, even with limited training data, further confirming its robustness across diverse datasets and varying data scales.

### C.2 More Quantitative Evaluation of Diffusion Model Generalization

Besides PSNR gap, cosine similarity between a generated image and its nearest neighbor in the training set is a metric for evaluating model generalization, as adopted by Kadkhodaie et al. (2024) and Zhang et al. (2024). Following Zhang et al. (2024), we compare the cosine similarity of standard DiT and DiT with local attention in Tab. 8. The reduced similarity with local attention indicates improved generalization.

### C.3 Inception Score Results

To further compare the generation quality of a model with and without enforcing attention locality, we present Inception Scores (IS) in Tab. 9. These results complement the FID results presented in Tab.

Table 8: The cosine similarity $\downarrow$  to the nearest training image with  $10^4$  training data between a standard DiT and a DiT equipped with local attention. In this setting, local attention with window sizes of (3, 5, 7, 9, 11, 13) is applied to the first six layers of the DiT. The best results are highlighted in **bold** font.

Model	CelebA	LSUN Church
DiT-XS/1	0.5939	0.5002
w/ Local	<b>0.5901</b>	<b>0.4988</b>

Table 9: Inception Score (IS) $\uparrow$  comparison between a standard DiT and a DiT equipped with local attention. In this setting, local attention with window sizes of (3, 5, 7, 9, 11, 13) is applied to the first six layers of the DiT. The best results are highlighted in **bold** font.

Model	CelebA		ImageNet		MSCOCO		LSUN Church		LSUN Bedroom		LSUN Bridge		LSUN Tower	
Train Set Size	$10^4$	$10^5$	$10^4$	$10^5$	$10^4$	$10^5$	$10^4$	$10^5$	$10^4$	$10^5$	$10^4$	$10^5$	$10^4$	$10^5$
DiT-XS/1	2.34	2.54	6.49	<b>8.50</b>	<b>6.46</b>	<b>6.85</b>	2.96	<b>3.15</b>	<b>3.41</b>	<b>3.27</b>	4.19	<b>4.53</b>	3.61	<b>3.91</b>
w/ Local	<b>2.38</b>	<b>2.55</b>	<b>6.96</b>	8.33	6.40	6.58	<b>2.99</b>	3.14	3.33	3.24	<b>4.23</b>	4.49	<b>3.62</b>	3.89
DiT-S/1	2.13	<b>2.58</b>	7.42	<b>7.99</b>	6.41	6.57	2.90	3.11	<b>3.37</b>	<b>3.20</b>	3.60	<b>4.57</b>	3.34	<b>4.00</b>
w/ Local	<b>2.15</b>	2.56	<b>7.52</b>	7.95	<b>6.46</b>	<b>6.58</b>	<b>2.93</b>	<b>3.14</b>	3.31	<b>3.20</b>	<b>3.81</b>	4.54	<b>3.38</b>	3.97

3. When trained with  $10^4$  images, using local attention mostly leads to improved Inception Scores, demonstrating the successful modification of a DiT’s generalization. Meanwhile, the Inception Scores of a DiT trained with  $10^5$  images tend to be comparable or biased towards the vanilla DiT, which also aligns with the FID results presented in the main paper.

#### C.4 FD-DINOv2 Results

In addition to FID and IS scores, we report the FD-DINOv2 score (Oquab et al., 2023) in Tab. 10 as an additional metric for evaluating image quality. When trained with  $10^4$  images. Incorporating local attention consistently improves the FD-DINOv2 score, indicating a successful enhancement of DiT’s generalization capabilities.

#### C.5 More Quantitative Results with Pixel-Space Diffusion Models

To confirm our findings in the main paper that attention map locality is an inductive bias that drives the generalization of a DiT, we present more quantitative comparisons. We provide the PSNR gap results in Tab. 11 and the FID in Tab. 12.

In the main paper, we find that a UNet has a worse generalization ability than a DiT with the same FLOPs, when measured by the PSNR gap. Tab. 11 and Tab. 12 demonstrate that the PSNR gap and FID of a UNet are worse than those of a DiT when the training image number  $N$  is insufficient and are better than those of a DiT when the training image number  $N$  is sufficient. This observation aligns well with our findings in the main paper.

Reducing the complexity of a neural network is a well-known way to improve a model’s generalization when the dataset is small. In Tab. 11 and Tab. 12, we compare the PSNR gap and FID of DiT-XS/1 and DiT-XXS/1. The latter is a smaller model with fewer hidden dimensions. A smaller DiT can reduce the PSNR gap and the FID when the training image number  $N$  is small. We recognize reducing network complexity as an orthogonal way to improve a DiT’s generalization. However, importantly, Tab. 11 and Tab. 12 show that it is less effective than using local attention. We present an additional study to reducing a DiT’s parameters in Appendix I.

#### C.6 Quantitative Results with Latent Diffusion Model

The attention map locality is identified to be the inductive bias that drives the generalization of a pixel-space DiT. In addition, we study whether the latent-space DiT also demonstrates such an inductive bias. To clarify, we use the pre-trained VAE from the official repository of DiT (Peebles & Xie, 2023). Then we train a DiT on the latent space of the pre-trained VAE, where the training images have a shape of  $256 \times 256$  with the corresponding latent code being of resolution  $32 \times 32$ .

Tab. 13 presents the PSNR gap and FID results comparing UNet, DiT-XS/1 and DiT-XS/1 with local attention, using CelebA and MSCOCO data. Comparing UNet and DiT-XS/1, DiT-XS/1 has a smaller

Table 10: FD-DINOv2 $\downarrow$  comparison with  $10^4$  training data between a standard DiT and a DiT equipped with local attention. In this setting, local attention with window sizes of (3, 5, 7, 9, 11, 13) is applied to the first six layers of the DiT. The best results are highlighted in **bold font**.

Model	CelebA	ImageNet	LSUN Church	LSUN Bedroom	LSUN Bridge	LSUN Tower
DiT-XS/1	182	868	614	747	709	451
w/ Local	<b>176</b>	<b>824</b>	<b>569</b>	<b>731</b>	<b>643</b>	<b>424</b>
DiT-S/1	292	789	590	659	863	566
w/ Local	<b>258</b>	<b>782</b>	<b>565</b>	<b>652</b>	<b>767</b>	<b>497</b>

Table 11: PSNR gap $\downarrow$  comparison based on pixel-space diffusion model. The training images have a resolution of  $32 \times 32$ .

Model	CelebA			ImageNet			MSCOCO		
	$N=10^3$	$N=10^4$	$N=10^5$	$N=10^3$	$N=10^4$	$N=10^5$	$N=10^3$	$N=10^4$	$N=10^5$
UNet	13.86	5.53	0.06	13.39	4.84	0.05	13.65	5.20	0.13
DiT-XXS/1	7.40	0.71	0.01	7.13	0.43	0.05	7.40	0.52	0.13
DiT-XS/1	7.49	0.80	0.01	7.77	1.08	0.05	7.36	0.60	0.13
DiT-XS/1 w/ Local	6.56	0.57	0.01	6.76	0.74	0.05	6.36	0.41	0.13

PSNR gap and FID when the training image number  $N$  is small, reconfirming the observation of the pixel-space experiments.

Comparing DiT-XS/1 with and without local attention, we observe the use of local attention to reduce the PSNR gap. However, we do not observe a smaller FID value when  $N$  is small, making it different from the pixel-space DiT. To investigate this further, we compare the attention map between pixel-space and latent-space DiTs in Fig. 6. We observe that the attention map locality gap between  $N=10^3$  and  $N=10^5$  is larger in pixel-space DiT than in latent-space DiT. We speculate that this is because larger training images ( $256 \times 256$  for the latent DiT compared with  $32 \times 32$  for pixel DiT), coupled with the VAE encoder, create more diverse information that enables a latent DiT to more easily achieve good generalization (reflected by attention map locality). Because of this, it is hard to improve a DiT’s FID further when  $N$  is small.

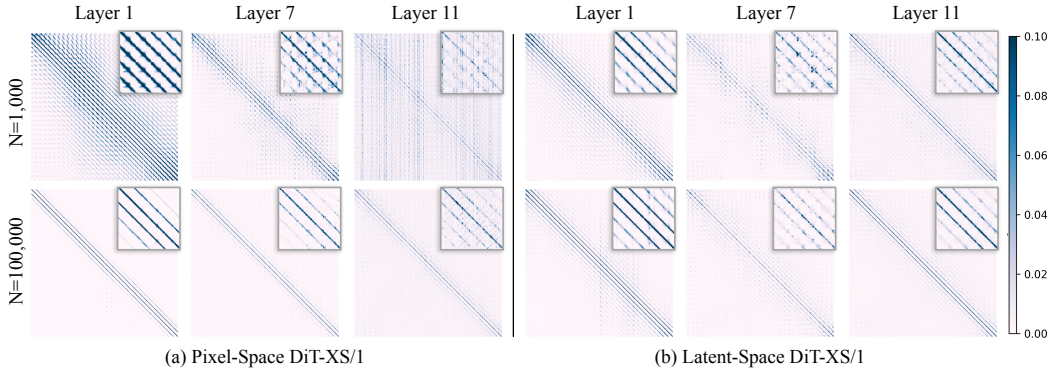


Figure 6: Attention map comparison between pixel-space and latent-space DiTs. The attention maps of latent-space DiT demonstrate smaller gaps between  $N=10^3$  and  $N=10^5$  in terms of the attention map locality.

## D Qualitative Results

In addition to the quantitative comparison, we present qualitative results with LSUN Bridge, LSUN Church, and ImageNet datasets in Fig. 7. With and without using local attention, two DiTs generate images of similar quality. Taking a closer look, for some samples, we find that DiTs trained with  $10^4$  and  $10^5$  images while using local attention produces images that are more like each other than DiTs trained without using local attention. For example, in each subfigure, the samples highlighted with **green** boxes are more similar, irrespective of whether the model was trained with  $N=10^4$  and

Table 12: FID $\downarrow$  comparison based on a pixel-space diffusion model. The training images have a resolution of  $32\times 32$ .

Model	CelebA		ImageNet		MSCOCO	
	$N=10^4$	$N=10^5$	$N=10^4$	$N=10^5$	$N=10^4$	$N=10^5$
UNet	9.8136	3.3871	61.3965	13.1302	58.4580	7.0214
DiT-XXS/1	9.0085	2.5749	33.2946	20.3075	26.0462	13.6076
DiT-XS/1	9.6932	2.6303	52.5650	17.3114	28.3496	12.9695
DiT-XS/1 w/ Local	8.4258	2.4988	43.8687	18.0671	24.4308	13.4735

Table 13: PSNR gap $\downarrow$  and FID $\downarrow$  comparison based on latent diffusion model with CelebA dataset. All models are trained with  $256\times 256$  images, where the corresponding latent codes have a resolution of  $32\times 32$ .

Model	CelebA					MSCOCO				
	PSNR Gap			FID		PSNR Gap			FID	
	$N=10^3$	$N=10^4$	$N=10^5$	$N=10^4$	$N=10^5$	$N=10^3$	$N=10^4$	$N=10^5$	$N=10^4$	$N=10^5$
UNet	8.69	3.35	0.11	36.5862	6.4805	5.74	1.84	0.12	159.7947	41.9908
DiT-XS/1	3.36	1.19	0.07	12.7045	8.5751	2.08	0.14	0.11	72.3063	68.1251
DiT-XS/1 w/ Local	2.21	0.17	0.07	13.6513	9.2621	1.19	0.13	0.10	78.8400	74.1681

$10=10^5$  images, than the samples surrounded with red boxes. This phenomenon aligns with our finding that the use of local attention can improve a DiT’s generalization.

To further verify the above observation we verify that the use of local attention makes images generated by a DiT trained with  $10^4$  images more similar to images generated by a DiT trained with  $10^5$  images. For this we randomly sample 50,000 images with each model, using a fixed random noise, so that the two models generate the same image content. Then we compute the average pixel intensity difference between the two generations. We show the results on LSUN Church, LSUN Bridge, and ImageNet data in Tab. 14.

## E Attention Locality Analysis Across Noise Levels

In addition to generalization, noise level represents an orthogonal factor that may affect the locality of attention maps. To explore this relationship, we analyze how attention locality varies with different input noise levels. In Tab. 15, we report the deviation score, a metric introduced in the main paper to quantify attention locality, under varying training data sizes and noise levels. The results show that: (1) attention locality strengthens as the model is trained with more data, reflecting improved generalization; and (2) the deviation score increases slightly with higher noise levels, suggesting that attention locality becomes marginally less pronounced under inputs with higher noise levels, though the overall locality remains strong.

## F Attention Map Consistency

To verify the robustness of the discovered inductive bias of a DiT, *i.e.*, the locality of attention maps, we obtain the attention maps corresponding to distinct input images and compare them visually. Specifically, we show the attention maps of the 1<sup>st</sup>, 6<sup>th</sup>, and 12<sup>th</sup> self-attention layers in Fig. 8, Fig. 9 and Fig. 10, respectively, using randomly selected 6 input images from the CelebA (Liu et al., 2015) dataset. In these figures, from top to bottom, each row is the attention map of a DiT model trained with  $N=10, 10^3, 10^4$ , and  $10^5$  images. Meanwhile, each column is related to an input image. For a better visualization, we use the logarithm normalization on attention maps before applying a colormap. For the same DiT Block, attention maps of different images demonstrate a similar pattern. Interestingly, we find that the attention maps of a DiT’s self-attention layers demonstrate a consistent pattern among different input images, suggesting that the attention maps of a DiT, after training, are part of its inductive biases rather than being mostly governed by a specific input image.





Figure 7: Visual comparison between DiT-XS/1 and DiT-XS/1 w/ Local Attention. As highlighted by red and green boxes, using local attention results in images from models trained with  $N=10^4$  and  $N=10^5$  images to be closer to each other.

## G Jacobian Eigenvector Analysis of Using Local Attention

The geometry-adaptive harmonic bases extracted via Jacobian eigenvectors is the inductive bias that drives the generalization of UNet-based diffusion models. Our analysis shows that these harmonic bases do not exist in a DiT. To further verify our finding, we extract the Jacobian eigenvectors of a DiT equipped with local attention. Fig. 11 compares the Jacobian eigenvectors of a DiT-XS/1 with and without using local attention. We follow [Kadkhodaie et al. \(2024\)](#) to extract the Jacobian eigenvectors and perform the analysis discussed in the main paper. The Jacobian eigenvectors of both DiTs demonstrate similar sparse patterns, showing no harmonic bases similar to the one observed in simplified [\(Kadkhodaie et al., 2024\)](#) and normal UNets. This observation corroborates our finding

Table 14: Averaged pixel intensity difference between generations of models trained with  $N=10^4$  and  $N=10^5$  images using pixel-space DiT-XS/1 with and without local attention. The generated images have a resolution of  $32 \times 32$ . Using local attention reduces the averaged pixel intensity difference.

Model	LSUN Bridge	LSUN Church	ImageNet
DiT-XS/1	13.0078	12.9620	17.5844
DiT-XS/1 w/ Local	11.0645	10.9402	15.2034

Table 15: Deviation scores under varying training data sizes and input noise levels, with the total number of diffusion steps fixed at 50. All results are based on a DiT-XS/1 trained with the CelebA dataset.

Diffusion Step	10	20	30	40
Training set size: $10^3$	0.236	0.210	0.208	0.210
Training set size: $10^4$	0.041	0.051	0.622	0.072
Training set size: $10^5$	0.032	0.038	0.042	0.046

that harmonic bases are not the driving factor of DiT’s generalization, regardless of the employed attention type.

## H Additional Visualizations of DiT’s Attention Maps

In Fig. 4 of the main paper, we visualize the attention maps of a DiT trained with  $N=10^3$ ,  $10^4$ , and  $10^5$  images, applying a colormap to the interval  $[0, 0.1]$ . More specifically, we first normalize an attention map to the range of  $[0, 1.0]$ . Then we apply the colormap to attention maps with an upper bound of 0.1, meaning that all values larger than 0.1 are colored identically. We choose 0.1 as the upper bound to make sure patterns of attention maps are easy to read. To further demonstrate the importance of attention map locality for the generalization of a DiT, we use different upper bounds for the colormap. Fig. 12 and Fig. 13 show attention maps with colormap upper bound 0.3 and 0.5, respectively. The stronger attention map locality can still be observed when increasing training image number  $N$  in both figures, confirming that attention map locality is an inductive bias of a DiT rather than being caused by a specific colormap upper bound.

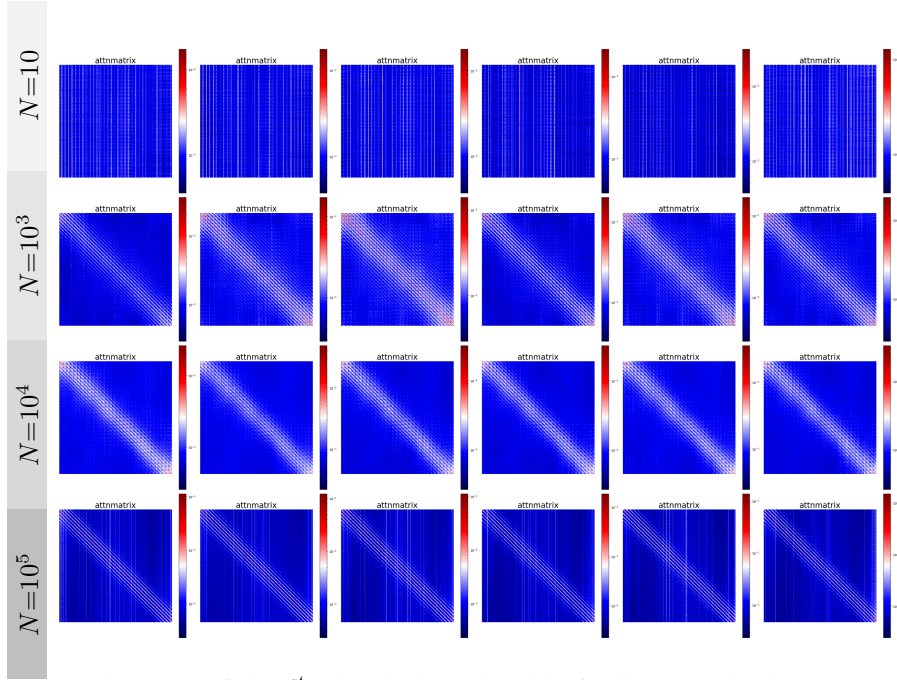


Figure 8: Attention maps of the 1<sup>st</sup> DiT Block produced by feeding 6 random images to DiT-XS/l models.

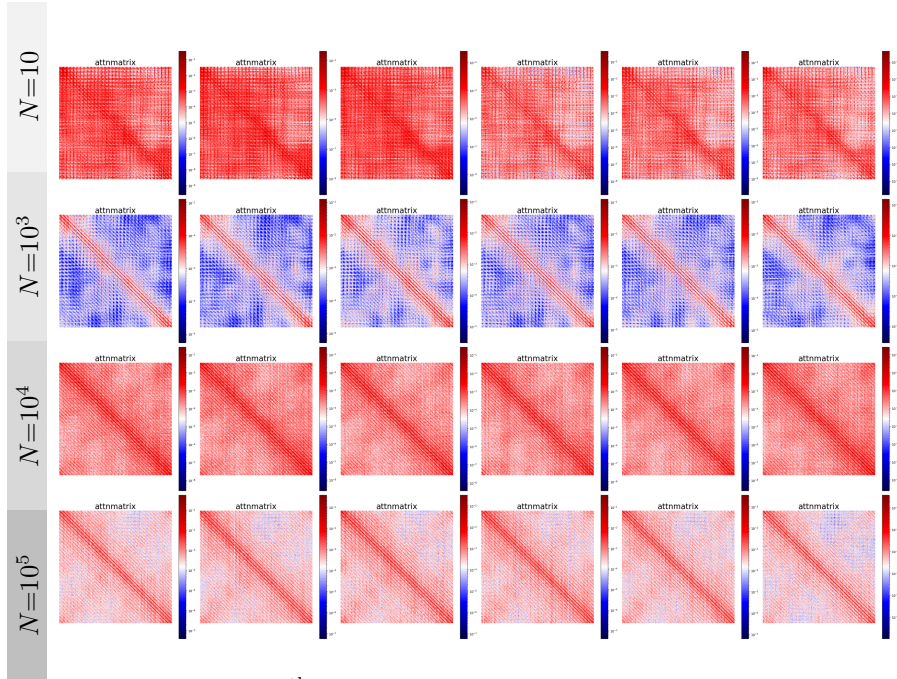


Figure 9: Attention maps of the 6<sup>th</sup> DiT Block produced by feeding 6 random images to DiT-XS/l models.



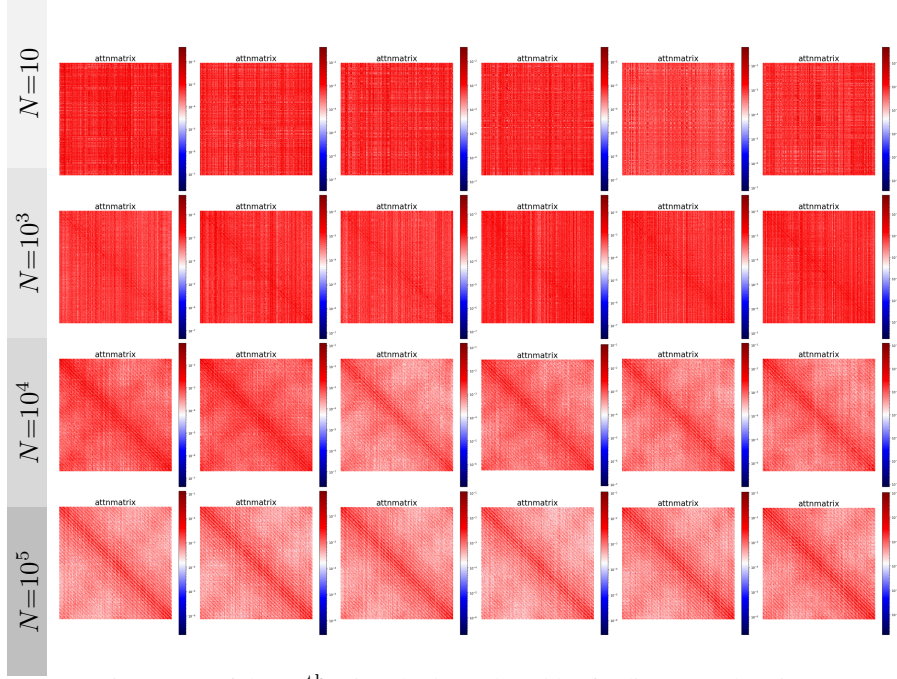


Figure 10: Attention maps of the 12<sup>th</sup> DiT Block produced by feeding 6 random images to DiT-XS/1 models.

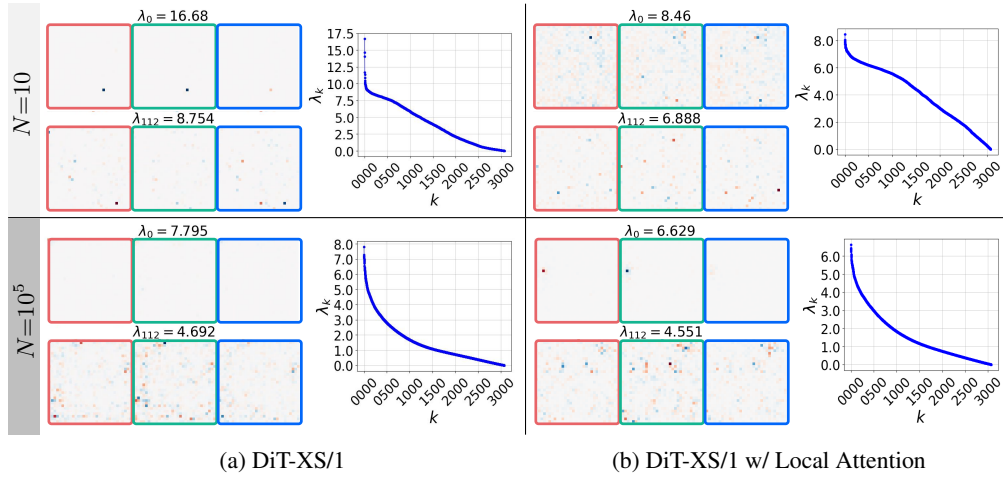


Figure 11: Jacobian eigenvector comparison between DiT-XS/1 w/ and w/o using local attention. In both cases, their Jacobian eigenvectors do not exhibit the harmonic bases observed in a UNet.

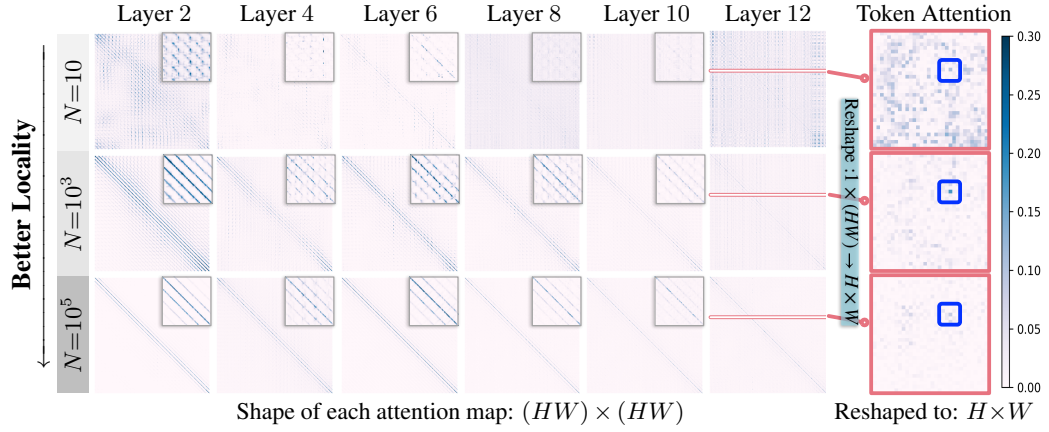


Figure 12: Attention maps of DiTs trained with 10,  $10^3$ , and  $10^5$  images. All attention maps are linearly normalized to the range  $[0, 1]$ , with a colormap applied to the interval  $[0, 0.3]$ .

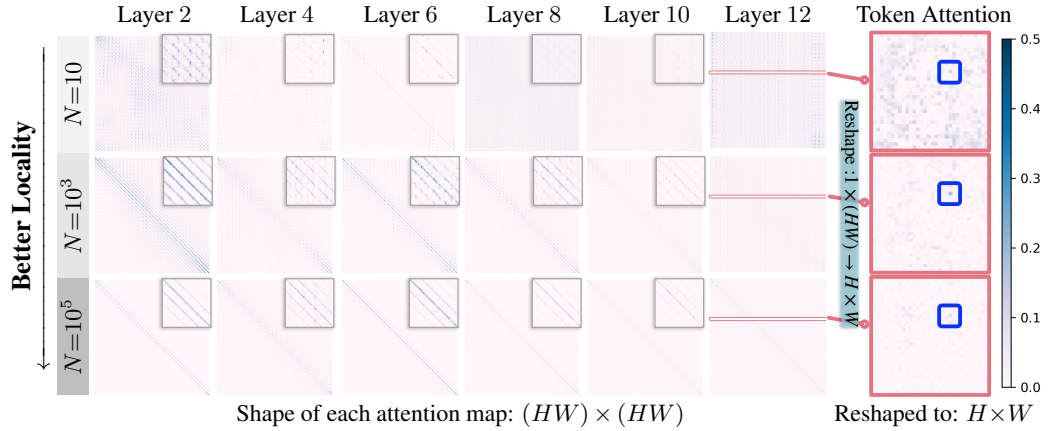


Figure 13: Attention maps of DiTs trained with 10,  $10^3$ , and  $10^5$  images. All attention maps are linearly normalized to the range  $[0, 1]$ , with a colormap applied to the interval  $[0, 0.5]$ .



## I Reducing Parameters of a DiT

Table 16: FID $\downarrow$  comparison between reducing DiT FLOPs and parameter size. Using parameter sharing is not as effective in reducing FID as the local attention. Using PCA will make the FID worse when  $N=10^4$ . Both are not as effective as using local attention in improving the FID with  $10^4$  training images.

Model	CelebA	
	$N=10^4$	$N=10^5$
DiT-XS/1	9.6932	2.6303
w/ Local	8.4580 <small>-1.2352</small>	2.5469 <small>0.0834</small>
w/ Weight Sharing	8.7819 <small>-0.9113</small>	2.5802 <small>-0.0501</small>
w/ PCA	11.3872 <small>+1.6940</small>	2.5482 <small>-0.0821</small>

We study two parameter reduction approaches for a DiT: parameter sharing and composing attention maps with PCA.

**Parameter Sharing.** For this approach, the 2<sup>nd</sup> and 4<sup>th</sup> DiT Blocks reuse the parameters of the 1<sup>st</sup> and 3<sup>rd</sup> DiT Blocks, respectively, leading to a DiT model with the same FLOPs but fewer parameters. Fig. 14 (row 3) shows the PSNR and PSNR gap of the DiT using this parameter sharing approach. Meanwhile, Tab. 16 demonstrates the FID of the same model trained with  $N=10^4$  and  $10^5$  images. Notably, sharing parameters in a DiT can reduce the PSNR gap, resulting in FID improvement when  $N=10^4$ . However, it is not as effective as using local attention in modifying the generalization of a DiT.

**Composing Attention Maps with PCA.** Another parameter reduction approach we explore is composition of attention maps of a DiT with PCA. To elaborate, we first collect the attention maps of 2048 images. Particularly, we noise each image with 10, 25 and 40 steps and obtain the attention maps of all attention heads corresponding to these noisy images, where the sampling diffusion step is set to 50. Taking the DiT-XS/1 model as an example, we collect a total of  $24,576 = 2048 \text{ (images)} \times 3 \text{ (diffusion steps)} \times 4 \text{ (attention heads)}$  attention maps. We use the DiT-XS/1 model trained with the ImageNet (Deng et al., 2009) dataset and collect attention maps of a DiT’s first three self-attention layers using randomly selected 2048 images from the testing set of the same dataset. Next, we compute the principal components of each self-attention layer from the corresponding attention maps. We use the low rank PCA function<sup>6</sup> of PyTorch and obtain the first 50 principal components, where each principal component has the same size as the attention map. Fig. 15 shows the principal components and the corresponding coefficients for the first three DiT Blocks. Notably, we find that PCA is effective in capturing the dominant diagonal patterns that indicate the locality in a DiT’s attention map. Finally, we use the principal components of the attention maps to reduce the parameters of a DiT. Concretely, we replace the two MLP layers that map the input tensor to query and key matrices by a smaller MLP mapping the input tensor to 50 coefficients for each principal component (PC). Then the new attention map is obtained as follows:

$$\text{Attention Map} = \text{Coefficients} \odot \text{PCs} + \delta, \quad (13)$$

where  $\odot$  denotes matrix multiplication while  $\delta = \frac{1.0}{1024}$  is used to ensure that the attention weights for a specific token sum to 1. We replace the normal attention maps of a DiT’s first three self-attention layers with the attention maps composed with principal components following Eq. (13). According to the PSNR and PSNR gap comparison in Fig. 14 (row 4) as well as the FID comparison in Tab. 16, reducing parameters of DiT by composing its first three attention maps with principal components cannot reduce a DiT’s PSNR gap, leading to a worse FID when  $N=10^4$ .

<sup>6</sup>[https://pytorch.org/docs/stable/generated/torch.pca\\_lowrank.html](https://pytorch.org/docs/stable/generated/torch.pca_lowrank.html)

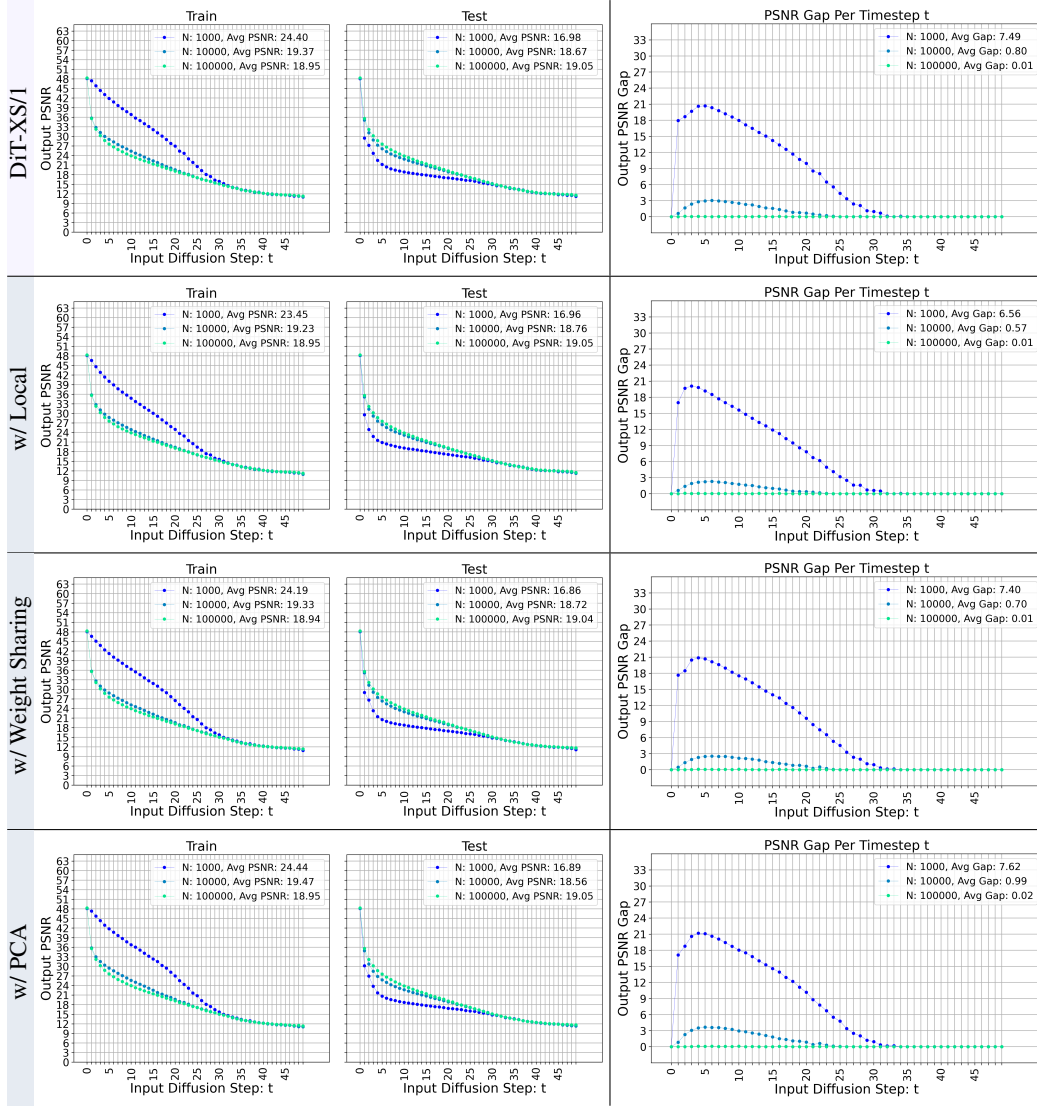
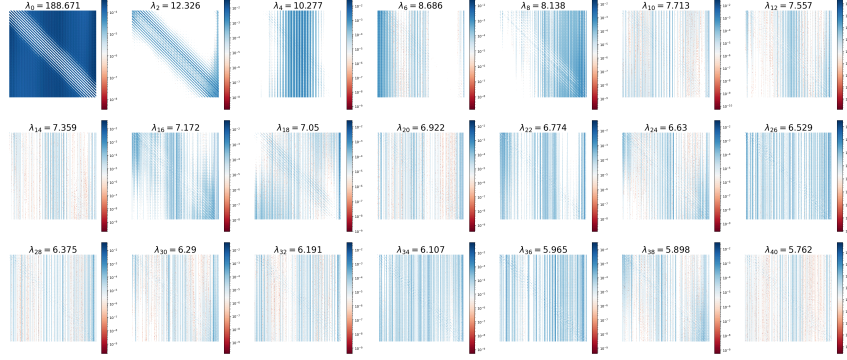
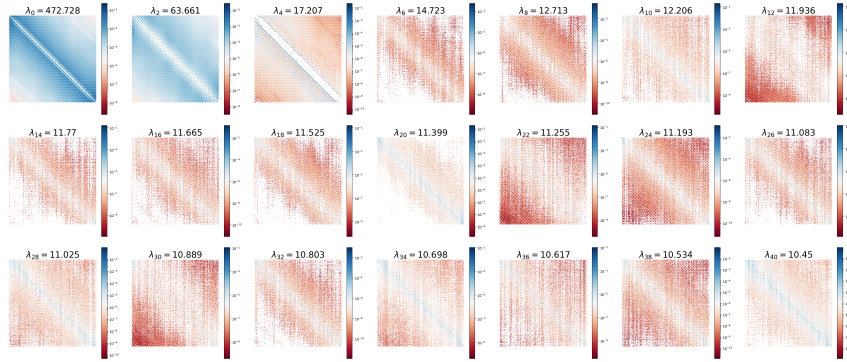


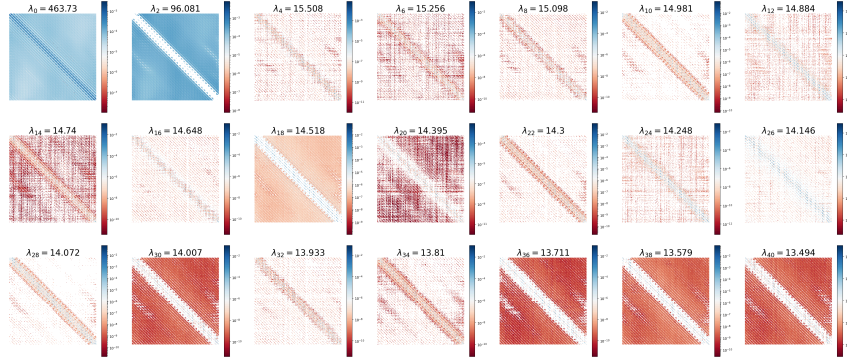
Figure 14: The PSNR (a) and PSNR gap (b) comparison. Taking DiT-XS/1 as the baseline (row 1), using local attention (row 2) can achieve decent PSNR gap improvement. In contrast, using parameter reduction approaches: weight sharing (row 3) and PCA (row 4), hardly achieve a PSNR gap improvement (row 3) or even make it worse then the baseline (row 4).



(a) Principal components from attention maps of the 1<sup>st</sup> DiT Block.



(b) Principal components from attention maps of the 2<sup>nd</sup> DiT Block.



(c) Principal components from attention maps of the 3<sup>rd</sup> DiT Block.

Figure 15: Principal components extracted from attention maps of different DiT Blocks. Based on a DiT-XS/1 model trained with  $N=10^5$  data from ImageNet, we perform PCA on attention maps of its first three layers, using 2048 images, resulting in a total of 24576 attention maps. For a better visualization, we adopt the logarithm normalization to principal components before applying colormaps.

## J Limitations

While this paper identifies the locality of attention maps as a key inductive bias contributing to the generalization of a DiT, and verifies this inductive bias by incorporating local attention windows into early layers of a DiT, this work has the following limitations:

**Formal Theoretical Proof.** In Appendix [B](#), we provide a theoretical analysis linking attention locality to model generalization. Specifically, we show that local attention encourages a simplicity bias, which in turn reduces sensitivity to data perturbations. This reduced sensitivity is then connected to flatter minima—an established indicator of better generalization. While our analysis is as rigorous as existing theoretical studies on diffusion model generalization, we acknowledge a limitation: it does not constitute a formal proof that attention locality directly leads to generalization. We note that developing formal generalization proofs for modern transformers remains a significant challenge for the research community.

**Encouraging Attention Locality in the Large Data Regime.** This paper reveals the attention locality bias of DiTs and demonstrates that promoting attention locality improves both generalization and generation quality in low-data regimes. However, we acknowledge a limitation: when sufficient training data is available (that is, when the DiT already generalizes well), further encouraging attention locality has limited effect on generalization or performance. While some may argue that the large-data regime is more common and practically relevant, the contribution of this paper lies in providing insights into the generalization behavior of DiTs, rather than focusing solely on performance improvement.

## K Broader Impact

This paper analyzes the inductive bias of diffusion transformers. The goal of this work is to advance the field of generative modeling by providing insights regarding diffusion model generalization. There are many potential societal consequences of a generative model. However, as an analysis paper that does not introduce new methods, we feel no societal impact must be specifically highlighted.