

# SD-DiT: Unleashing the Power of Self-supervised Discrimination in Diffusion Transformer

Rui Zhu, Yingwei Pan, Yehao Li, Ting Yao, Zhenglong Sun, Tao Mei, Chang Wen Chen



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen

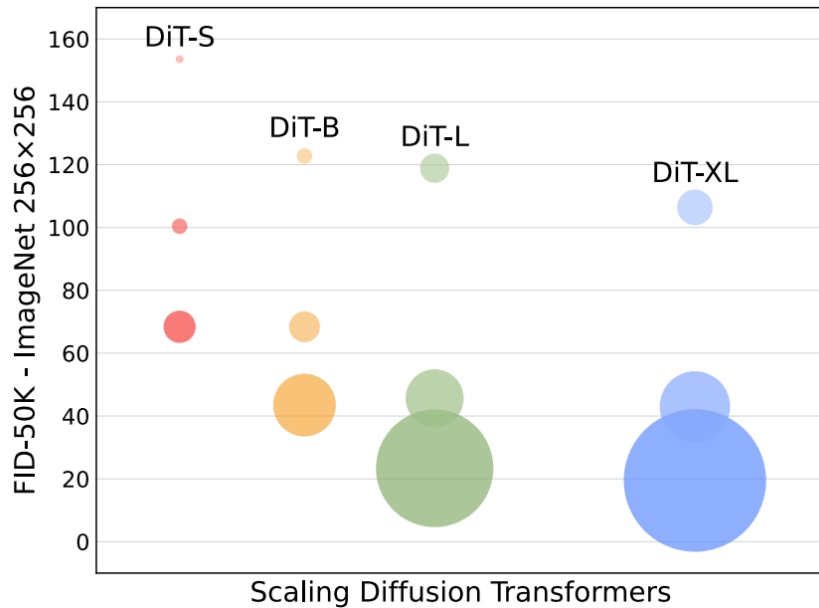


智象未来  
HiDream.ai

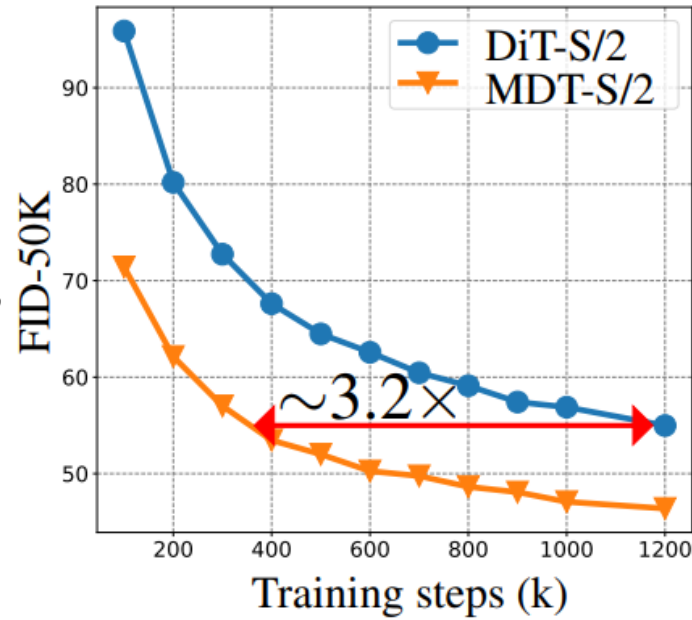


THE HONG KONG  
POLYTECHNIC UNIVERSITY  
香港理工大學

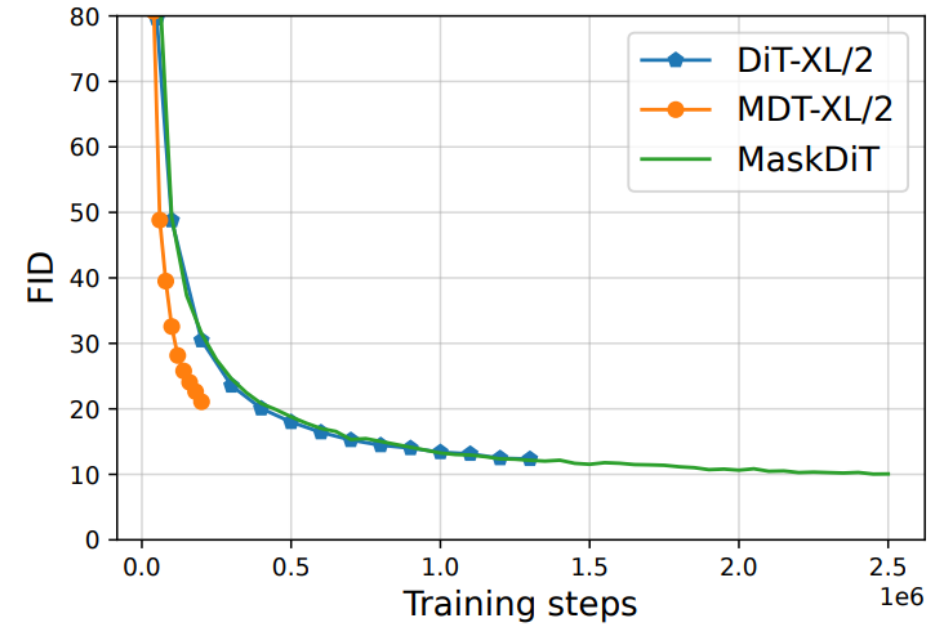
# Diffusion Transformer - Scalable but Slow Convergence



DiT: Scalable  
Arxiv 2022.12



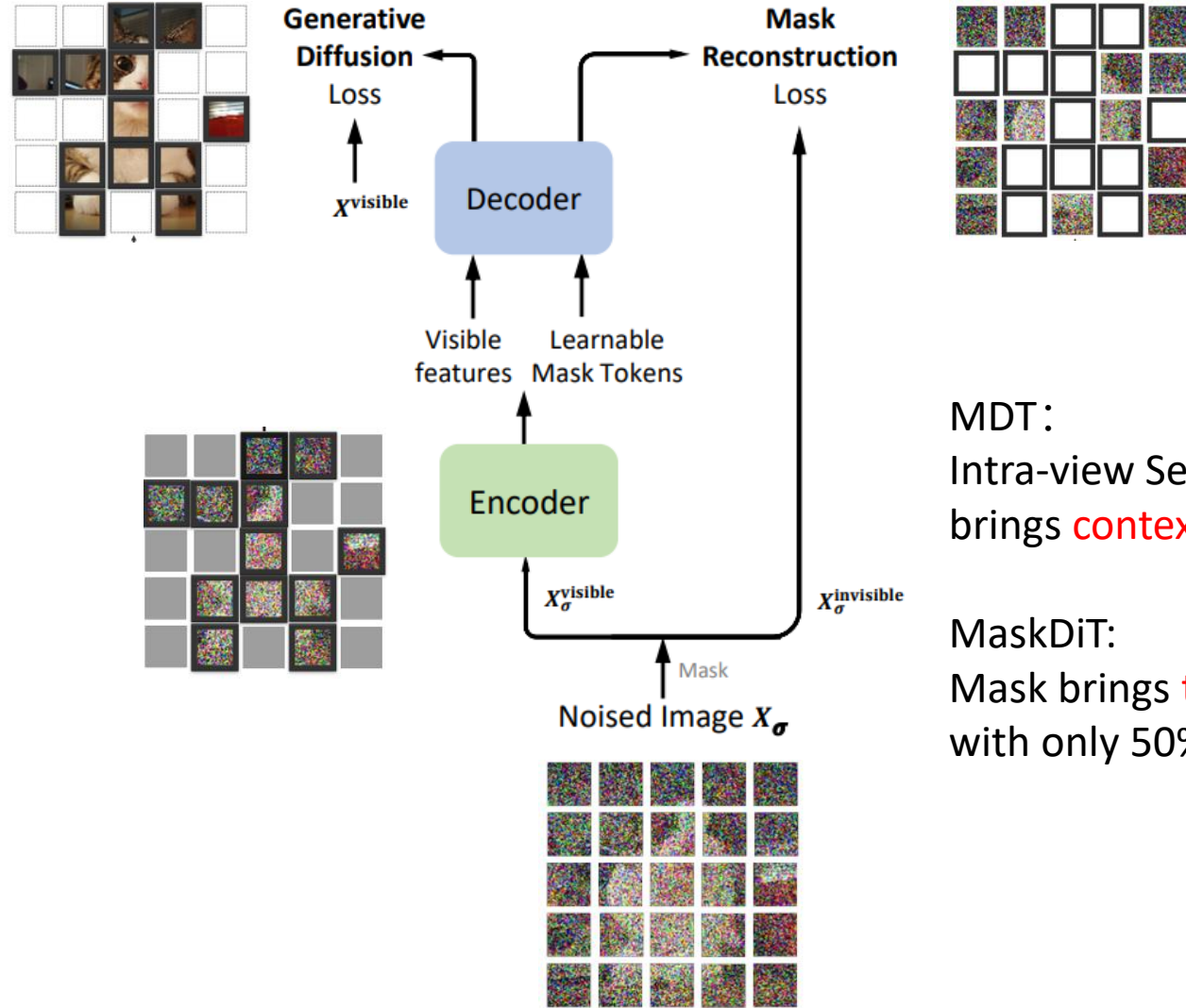
MDT: Fast Convergence  
Arxiv 2023.03



MaskDiT: Efficiency  
Arxiv 2023.06

[1] Peebles, William, and Saining Xie. "Scalable diffusion models with transformers." ICCV 2023.  
[2] Gao, Shanghua, et al. "Masked diffusion transformer is a strong image synthesizer." ICCV 2023.  
[3] Zheng, Hongkai, et al. "Fast training of diffusion models with masked transformers." TMLR 2024.

# Why Bring Mask to DiT – Contextual Relation inside View



MDT:  
Intra-view Self-Reconstruction via mask  
brings **contextual relation learning**

MaskDiT:  
Mask brings **training efficiency**  
with only 50% input

[1] Gao, Shanghua, et al. "Masked diffusion transformer is a strong image synthesizer." ICCV 2023.

[2] Zheng, Hongkai, et al. "Fast training of diffusion models with masked transformers." TMLR 2024.

# Can we impose **Inter-view Discrimination** to DiT?

Representation Learning

Generative Modeling

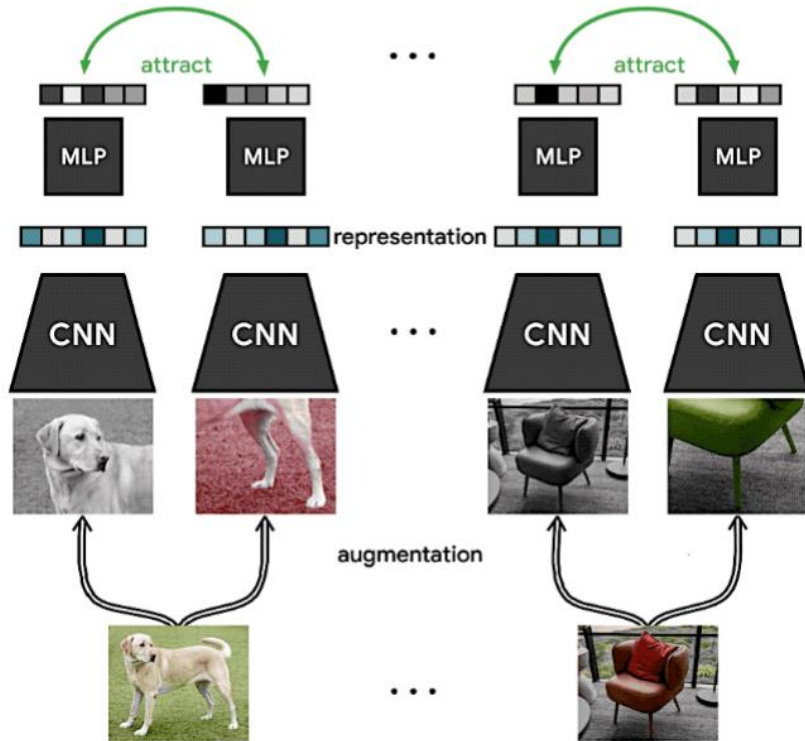
How to construct discriminative view pair  
for Generative Diffusion Transformer?

Contrastive Self-supervised Learning relies on  
Data Augmentation for positive pair

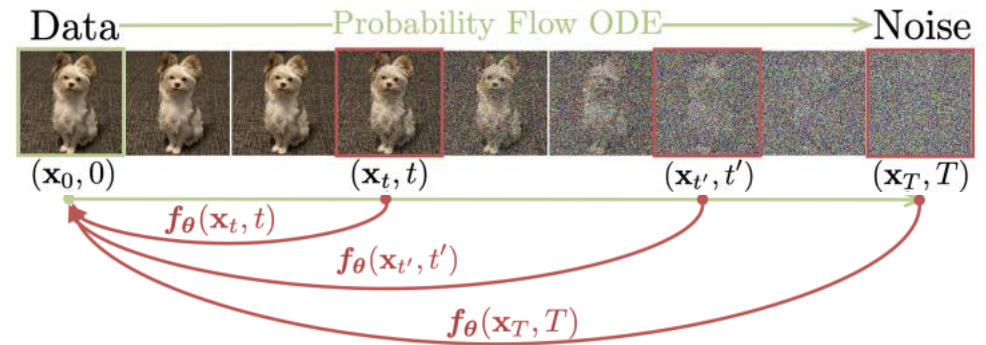
# Can we impose **Inter-view Discrimination** to DiT?

Representation Learning

Generative Modeling



$$p_{data \odot Aug} \rightarrow p_{data}$$



Inspired by Consistency models, whose outputs of the points on the same PF-ODE trajectory are **consistent**

$$f(x_\sigma, \sigma) = f(x_{\sigma'}, \sigma'), \quad \sigma, \sigma' \in [\sigma_{\min}, \sigma_{\max}].$$

$$f : (x_\sigma, \sigma) \mapsto x_{\sigma_{\min}}$$

We construct discriminative pair by adding noise  $(x_{\sigma_S}, x_{\sigma_T})$  on the same PF-ODE

$$p_{\sigma_S} \rightarrow p_{\sigma_T}$$

[1] <https://github.com/google-research/simclr>

[2] Song, Yang, Dhariwal Prafulla, Chen Mark, Sutskever Ilya. "Consistency models." ICML 2023.

# Preliminary

## PF-ODE

$$d\mathbf{x}_t = [\boldsymbol{\mu}(\mathbf{x}, t) - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)] dt.$$

EDM utilizes  $p_{\sigma}(\mathbf{x})$  instead of  $p_t(\mathbf{x})$

$$\boldsymbol{\mu}(\mathbf{x}, t) := \mathbf{0} \text{ and } g(t) := \sqrt{2t}$$

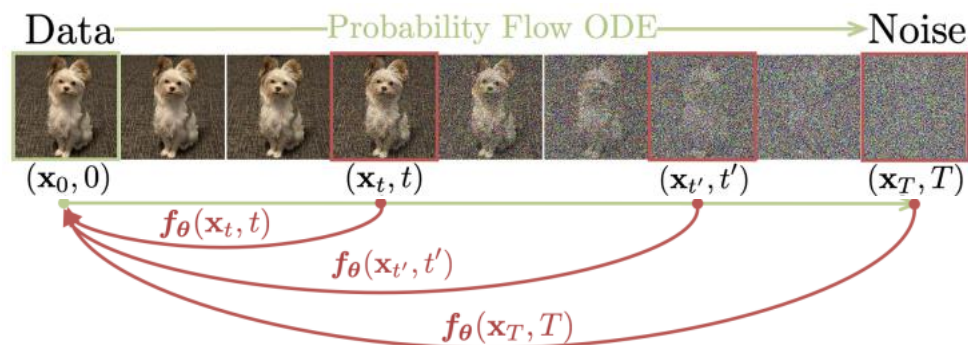
## PF-ODE in EDM

$$d\mathbf{x} = -\sigma \nabla_{\mathbf{x}} \log p_{\sigma}(\mathbf{x}) d\sigma, \quad \sigma \in [\sigma_{\min}, \sigma_{\max}],$$

$$p_{\sigma}(\mathbf{x}) = p_{\text{data}}(\mathbf{x}) * \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\mathbf{x}_{\sigma} = \mathbf{x}_0 + \mathbf{n}, \quad \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

## Generative Modeling



Consistency models, whose outputs of the points on the same PF-ODE trajectory are **consistent**

$$\mathbf{f}(\mathbf{x}_{\sigma}, \sigma) = \mathbf{f}(\mathbf{x}_{\sigma'}, \sigma'), \quad \sigma, \sigma' \in [\sigma_{\min}, \sigma_{\max}].$$

$$\mathbf{f} : (\mathbf{x}_{\sigma}, \sigma) \mapsto \mathbf{x}_{\sigma_{\min}}$$

We construct discriminative pair by adding noise  $(\mathbf{x}_{\sigma_S}, \mathbf{x}_{\sigma_T})$  on the same PF-ODE

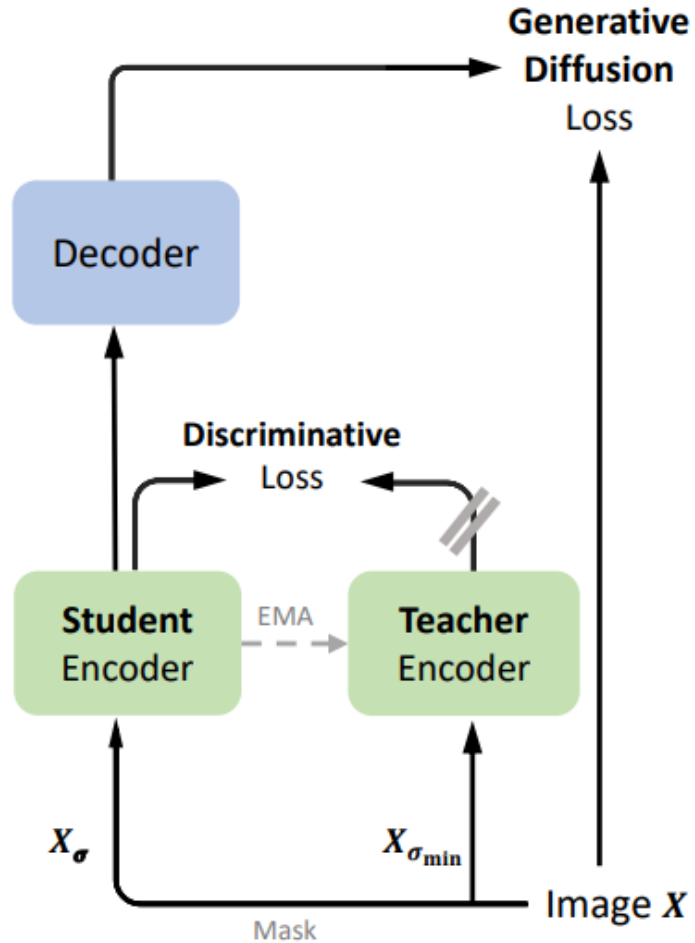
$$p_{\sigma_S} \rightarrow p_{\sigma_T}$$

[1] Song, Yang, Dhariwal Prafulla, Chen Mark, Sutskever Ilya. "Consistency models." ICML 2023.

[2] Song, Yang, et al. "Score-based generative modeling through stochastic differential equations." ICLR2021.

[3] Karras, Tero, et al. "Elucidating the design space of diffusion-based generative models." NeurIPS 2022.

# SD-DiT: Discriminative Objective



(a) SD-DiT

$$\mathbf{x}_{\sigma_S} = \mathbf{x}_0 + \mathbf{n}, \quad \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_S^2 \mathbf{I}), \quad \sigma_S \in [\sigma_{\min}, \sigma_{\max}]$$

$$\mathbf{x}_{\sigma_T} = \mathbf{x}_0 + \mathbf{n}, \quad \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_{\min}^2 \mathbf{I})$$

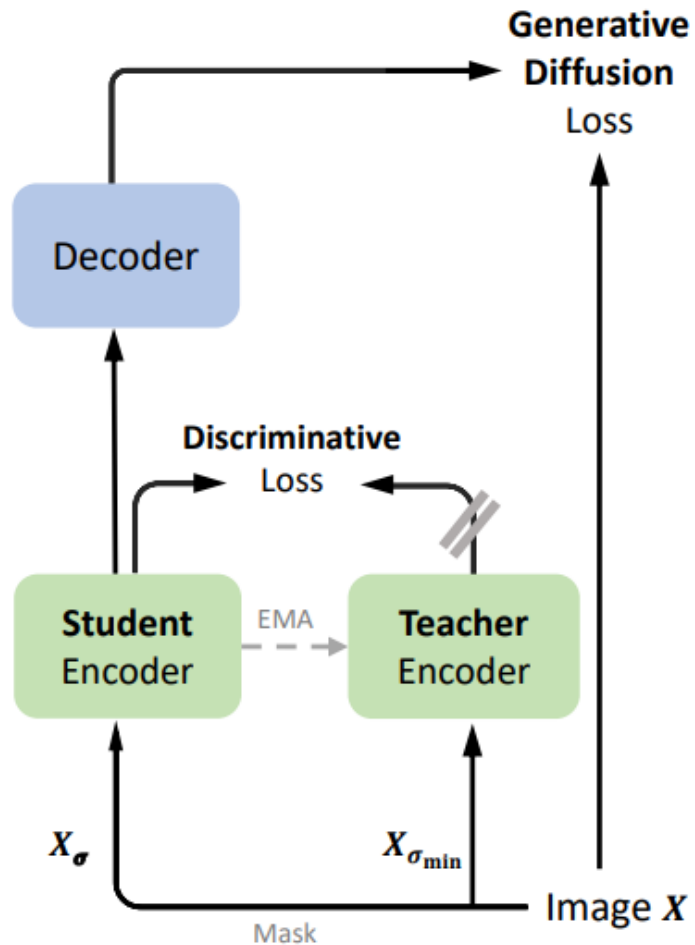
$$P_{S_i} = \frac{\exp(j_{\theta}(\mathbf{e}_{S_i})/\tau_S)[k]}{\sum_{k=1}^K \exp(j_{\theta}(\mathbf{e}_{S_i})/\tau_S)[k]},$$

$$\mathcal{L}_D(i) = - \sum_k P_{T_i} \log(P_{S_i}).$$

Loss on visible tokens and CLS token:

$$\mathcal{L}_D = \frac{1}{(1 - \mathcal{M})} \sum_{i \in (1 - \mathcal{M})} \mathcal{L}_D(i) + \mathcal{L}_D([\text{CLS}]).$$

# Various Teacher Noise in Discriminative Pair



(a) SD-DiT

$$\mathbf{x}_{\sigma_S} = \mathbf{x}_0 + \mathbf{n}, \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_S^2 \mathbf{I}), \sigma_S \in [\sigma_{\min}, \sigma_{\max}]$$

$$\mathbf{x}_{\sigma_T} = \mathbf{x}_0 + \mathbf{n}, \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_{\min}^2 \mathbf{I})$$

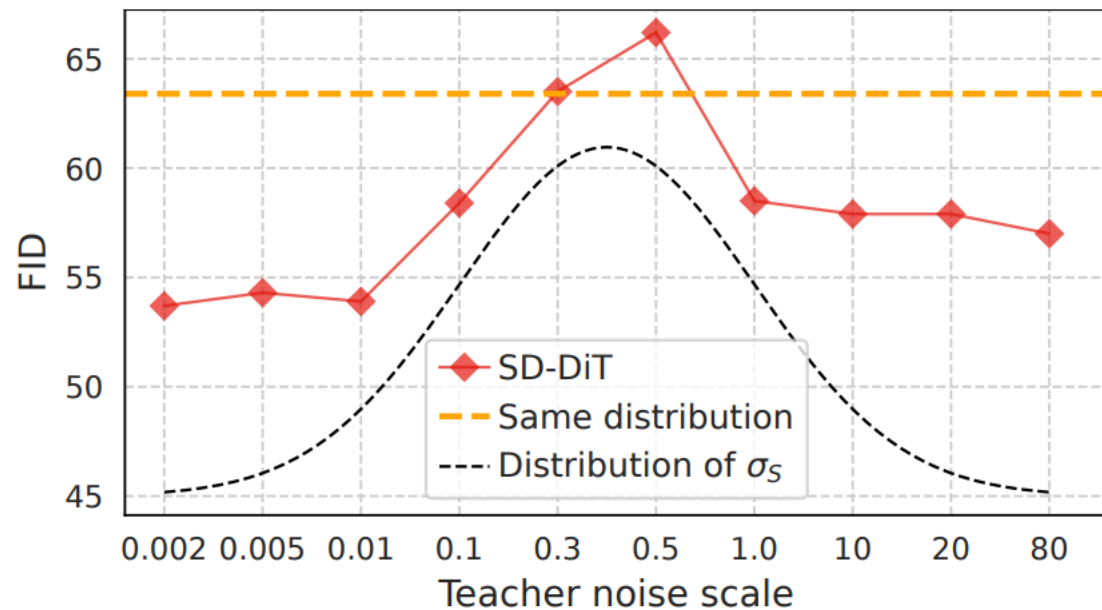
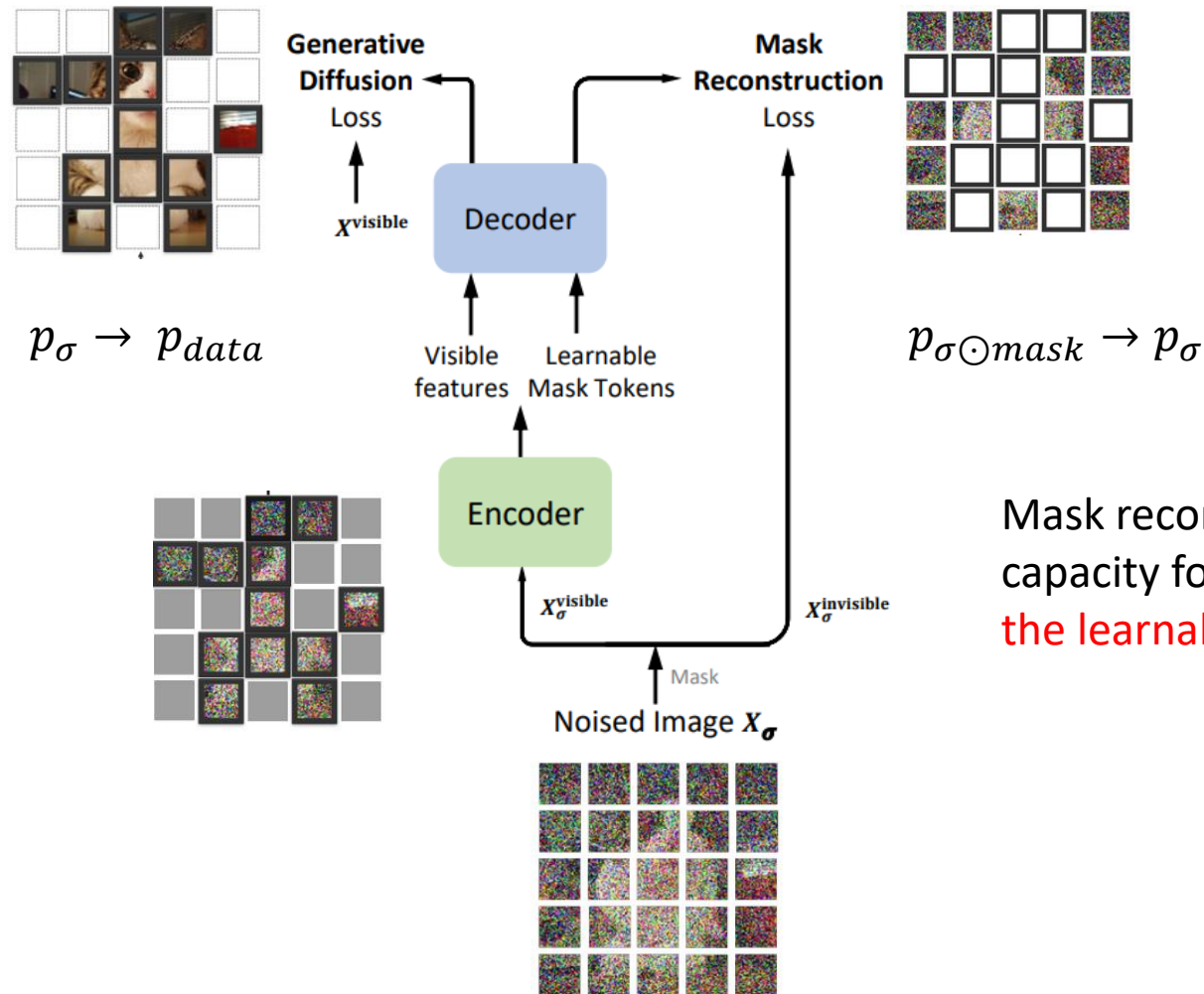


Figure 6. FID vs. teacher noise on SD-DiT-S/2 with 400k steps.

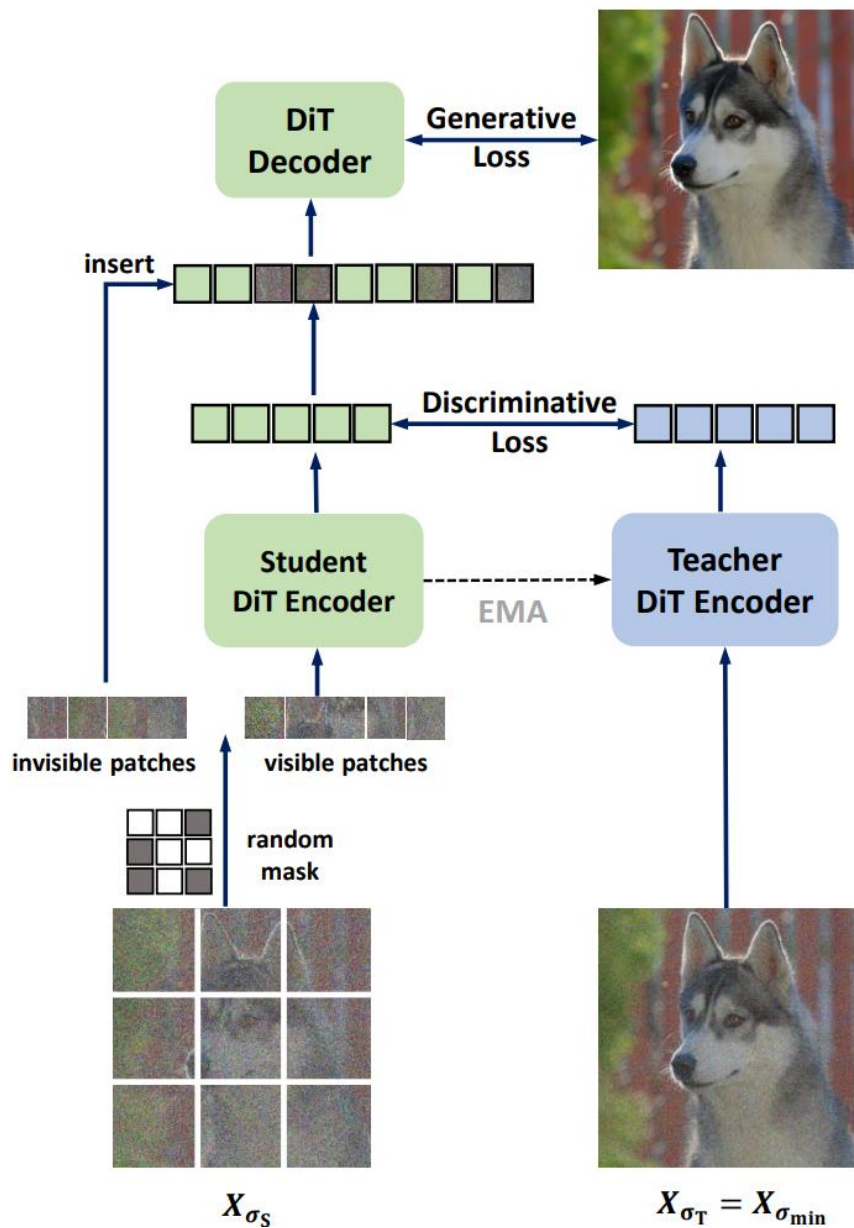


# Fuzzy relations: Mask Reconstruction vs. Generative Diffusion



Mask reconstruction loss wastes model capacity for representation learning and **the learnable mask tokens**.

# SD-DiT: Decoupled Encoder-decoder w/o mask tokens



- Decoder for generative loss:  $p_{\sigma} \rightarrow p_{data}$
- Encoder for discriminative loss:  $p_{\sigma} \rightarrow p_{min}$
- Keep masks for training efficiency and location contextual awareness.
- Remove the mask reconstruction loss  $p_{\sigma \odot mask} \rightarrow p_{\sigma}$   
(which wastes model capacity for representation learning)

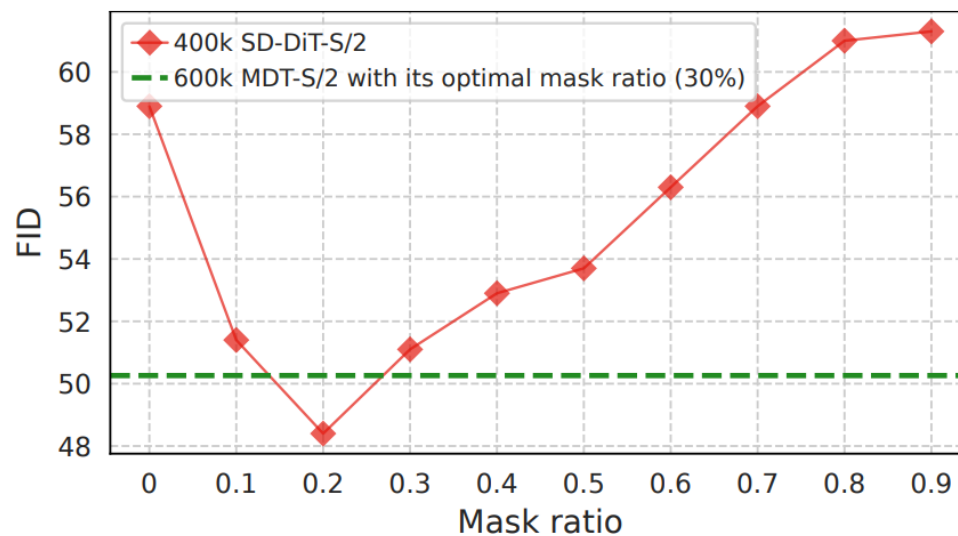


Figure 5. FID vs. mask ratio on SD-DiT-S/2 with 400k steps.

# Experiments on ImageNet: Fast Convergence

Method	Training Steps(k)	FID-50K↓
DiT-S/2 [45]	400	68.40
MDT-S/2 [19]	400	53.46
<b>SD-DiT-S/2</b>	400	<b>48.39</b>
DiT-B/2 [45]	400	43.47
MDT-B/2 [19]	400	34.33
<b>SD-DiT-B/2</b>	400	<b>28.62</b>
DiT-XL/2 [45]	7000	9.62
MaskDiT-XL/2 [73]	1300	12.15
MDT-XL/2 [19]	1300	9.60
<b>SD-DiT-XL/2</b>	1100	9.66
<b>SD-DiT-XL/2</b>	1300	<b>9.01</b>

Table 1. Performance comparison with state-of-the-art DiT-based approaches under various model sizes on ImageNet  $256 \times 256$  for class-conditional image generation (batch size: 256).

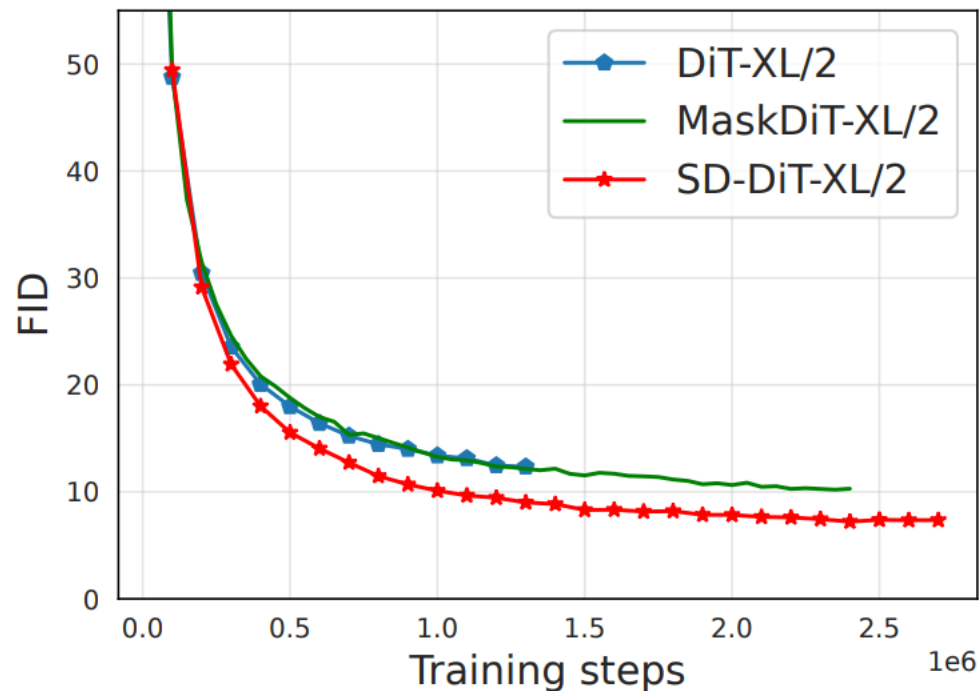


Figure 4. Comparison of convergence speed with SOTA DiT-based approaches in DiT-XL backbone (batch size: 256). The results of DiT and MaskDiT are directly cited from MaskDiT [81]. Our SD-DiT-XL/2 consistently outperforms DiT-XL/2 and MaskDiT-XL/2 across training steps, leading to better training convergence.

# Experiments on ImageNet: Compare with SOTAs

Method	Cost(Iter $\times$ BS)	FID $\downarrow$	sFID $\downarrow$	IS $\uparrow$	Prec. $\uparrow$	Rec. $\uparrow$
VQGAN [16]	-	15.78	78.3	-	-	-
BigGAN-deep [5]	-	6.95	7.36	171.4	0.87	0.28
StyleGAN [57]	-	2.30	4.02	265.12	0.78	0.53
I-DDPM [43]	-	12.26	-	-	0.70	0.62
MaskGIT [9]	1387k $\times$ 256	6.18	-	182.1	0.80	0.51
CDM [29]	-	4.88	-	158.71	-	-
ADM [14]	1980k $\times$ 256	10.94	6.02	100.98	0.69	0.63
ADM-U [14]	—	7.49	5.13	127.49	0.72	0.63
LDM-8 [50]	4800k $\times$ 64	15.51	-	79.03	0.65	0.63
LDM-4 [50]	178k $\times$ 1200	10.56	-	103.49	0.71	0.62
MaskDiT-XL/2 [73]	2000k $\times$ <b>1024</b>	5.69	10.34	177.99	0.74	0.60
DiT-XL/2 [45]	7000k $\times$ 256	9.62	6.85	121.50	0.67	<b>0.67</b>
MDT-XL/2 [19]	2500k $\times$ 256	7.41	<b>4.95</b>	121.22	0.72	0.64
SD-DiT-XL/2	2400k $\times$ 256	<b>7.21</b>	5.17	<b>144.68</b>	<b>0.72</b>	0.61

Table 2. Performance comparison with state-of-the-art methods on ImageNet  $256\times 256$  for class-conditional image generation. Similar to most DiT-based approaches, here we report the results of our SD-DiT in DiT-XL backbone with 256 batch size, while MaskDiT reports results with the largest batch size (1024).

# Thanks for Listening!



**Rui Zhu**  
The Chinese University  
of Hong Kong, Shenzhen



**Yingwei Pan**  
HiDream.ai



**Yehao Li**  
HiDream.ai



**Ting Yao**  
HiDream.ai



**Zhenglong Sun**  
The Chinese University of Hong Kong, Shenzhen



**Tao Mei**  
HiDream.ai



**Chang Wen Chen**  
The Hong Kong Polytechnic University

Any question: [ruizhu@link.cuhk.edu.cn](mailto:ruizhu@link.cuhk.edu.cn)