

# Representation Alignment for Diffusion Transformers without External Components

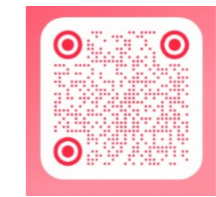
Dengyang Jiang<sup>2,1\*</sup> Mengmeng Wang<sup>3,1\*</sup> Liuzhuozheng Li<sup>1</sup> Lei Zhang<sup>2</sup> Haoyu Wang<sup>2</sup> Wei Wei<sup>2</sup> Guang Dai<sup>1</sup> Yanning Zhang<sup>2</sup> Jingdong Wang<sup>4†</sup>

1. SGIT AI Lab, State Grid Corporation of China 2. Northwestern Polytechnical University 3. Zhejiang University of Technology 4. Baidu Inc.

Project: <https://www.vjdy.github.io/sra>

Code: <https://github.com/vjdy/SRA>

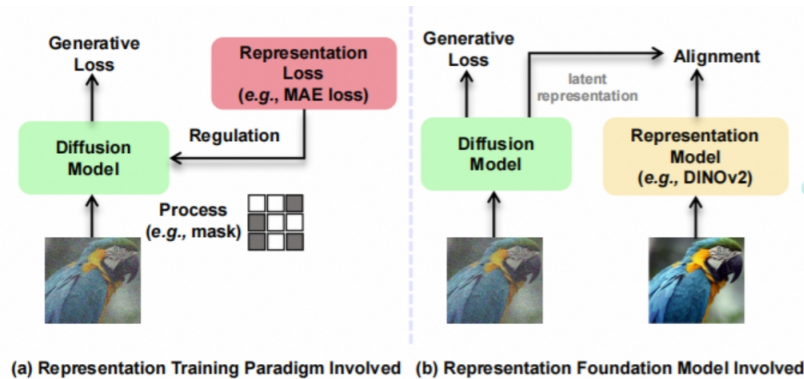
Arxiv: <https://arxiv.org/abs/2505.02831>



Contact us



## Background: External Representation Components Involved for Better Training Diffusion Transformer



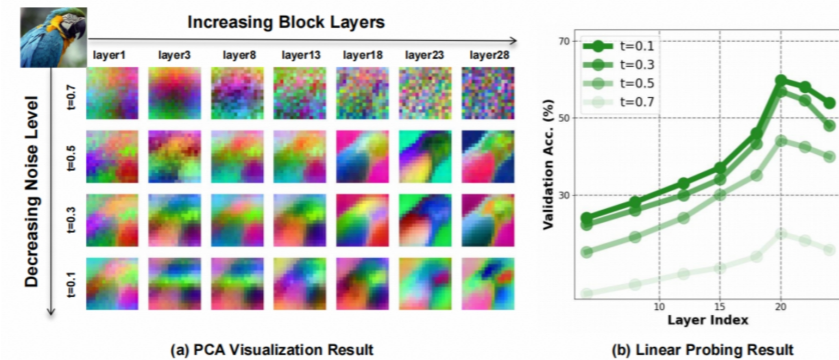
While recent studies show that learning a meaningful internal representation can accelerate generative training, existing approaches either introducing an auxiliary external representation task or relying on a large-scale pre-trained external encoder to provide representation guidance.

## Future Work: Scale-up and Efficiency

**Scaling up:** Applying SRA to large scale T2I and T2V domain where the data distribution is more complex and do not exist a dominant representation encoder.

**Efficiency optimizing:** SRA requires an additional forward pass. As model size increase, this incurs increasingly computational and memory cost. Developing more efficient variants of self-representation alignment constitutes a vital direction for future research.

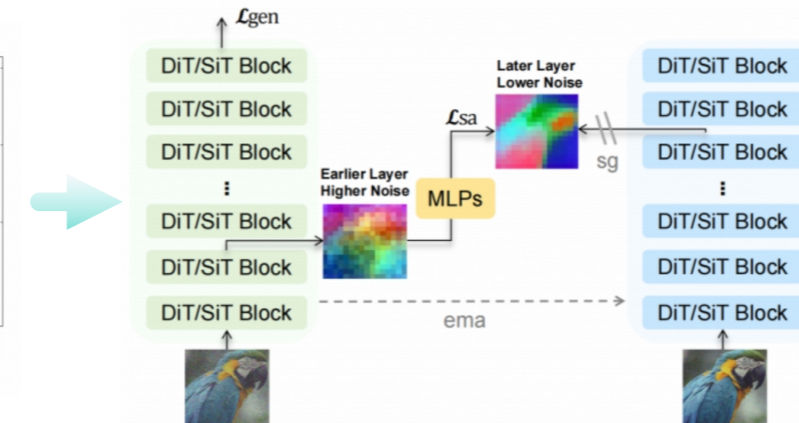
## Inspiration and Observation: Representations in Pre-trained Diffusion Transformer



**Inspiration:** The generative mechanism by which the diffusion model operates can be generally considered as a coarse to fine process, and how about the representations?

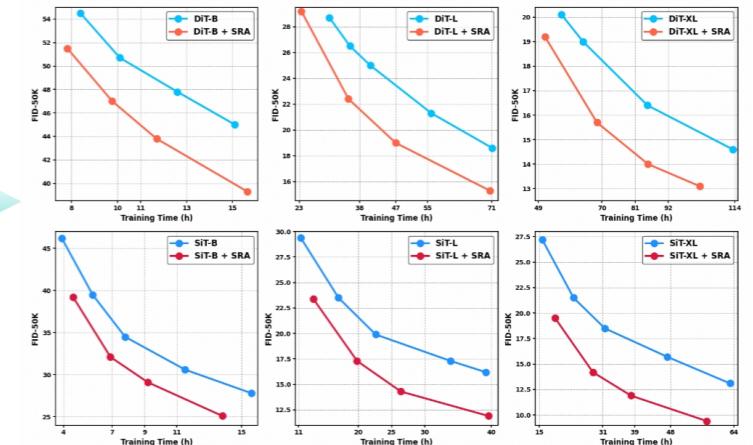
**Observation:** Diffusion transformer already got a roughly from bad to good discriminative process when only generative training is performed.

## Method: Self-Representation Alignment (SRA)



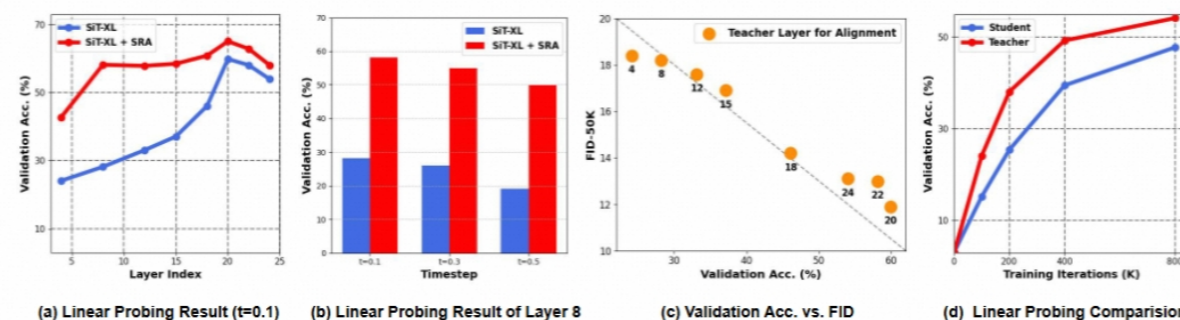
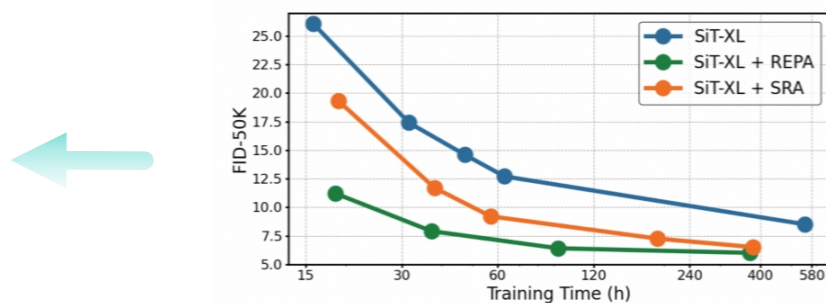
SRA aligns the output latent representation in earlier layer with higher noise to that in later layer with lower noise to progressively enhance the overall representation learning during only generative process.

## Results: Fast Convergence and Better Scalability



SRA shows benefit on different baselines across different model size. Meanwhile, As the size of the model increases, the gain also increases.

## Results: Improved Generation and Representation Capacity



Model	Epochs	FID↓	sFID↓	IS↑	Pre.↑	Rec.↑
<i>Pixel diffusion</i>						
VDM++	-	2.65	-	278.1	-	-
Simple diffusion	800	4.28	-	171.0	-	-
<i>Latent diffusion, Transformer</i>						
DiT-XL/2	600	3.04	5.02	240.8	0.84	0.54
SiT-XL/2	600	2.62	4.18	252.2	<b>0.84</b>	0.57
MaskDiT	600	2.50	5.10	256.3	0.84	0.56
SiT + REPA	200	<b>2.08</b>	4.19	274.6	0.83	0.58
SiT + SRA (ours)	<b>200</b>	<b>2.17</b>	<b>4.15</b>	<b>279.3</b>	<b>0.83</b>	<b>0.59</b>

Method	FID↓	PickScore↑
<i>ODE, NFE=50, Trained for 150K iter</i>		
MMDiT	5.86	20.05
MMDiT + REPA	4.60	20.88
MMDiT + SRA	4.85	21.14

Model	Epochs	Tokenizer	FID↓	sFID↓	IS↑	Pre.↑	Rec.↑
<i>Pixel diffusion</i>							
ADM-U	400	-	3.94	6.14	186.7	0.82	0.52
VDM++	560	-	2.40	-	225.3	-	-
Simple diffusion	800	-	2.77	-	211.8	-	-
CDM	2160	-	4.88	-	158.7	-	-
<i>Latent diffusion, U-Net</i>							
LDM-4	200	LDM-VAE	3.60	-	247.7	<b>0.87</b>	0.48
<i>Latent diffusion, Transformer</i>							
DiT-XL/2	1400	SD-VAE	2.27	4.60	278.2	0.83	0.57
SiT-XL/2	1400	SD-VAE	2.06	<b>4.50</b>	270.3	0.82	0.59
SD-DiT	480	SD-VAE	3.23	-	-	-	-
MaskDiT	1600	SD-VAE	2.28	5.67	276.6	0.80	0.61
DiT + TREAD	740	SD-VAE	1.69	4.73	292.7	0.81	0.63
SiT + REPA	800	SD-VAE	<b>1.42</b>	4.70	305.7	0.80	<b>0.65</b>
SiT + MAETok	800	MAE-Tok	1.67	-	311.2	-	-
SiT + SRA (ours)	<b>400</b>	SD-VAE	1.85	<b>4.50</b>	297.2	0.82	0.61
SiT + SRA (ours)	<b>800</b>	SD-VAE	1.58	4.65	<b>311.4</b>	0.80	0.63

SRA shows comparable or superior performance against other methods that leverage either auxiliary representation training task or external representation encoder on ImageNet 256/512 and COCO T2I. The representation capacity of the model is also by SRA during training.