

# ACCELERATING DIFFUSION MODEL WITH DYNAMIC ALIGNMENT

Anonymous authors

Paper under double-blind review

## ABSTRACT

Recent studies have shown improvements in both generation quality and training efficiency by constraining representations during the denoising process of generative diffusion models. While distilling simple visual representations is effective, it can lead to over-alignment issues. When the model achieves alignment early in training, these simple representations can become hindrance to training the generative capacity. Building upon prior efforts that addressed this problem from the perspectives of alignment objectives and training strategies, we introduce DyA. First, we incorporate richer alignment materials to address the problem of overly simplistic representations at the source. Second, we use the internal denoising time of the diffusion model as an indicator variable to dynamically adjust the constraint strength of different levels of information. Finally, we employ the Stochastic Dropout Strategy (SDS), which allows the model to emphasize generative capacity training while providing guidance throughout the entire process. Experiments have shown that this approach improves both generation quality and training efficiency. The DyA accelerates SiT training by approximately 20 times, achieving performance comparable to SiT-XL model trained for 7M steps in just around 350K steps.

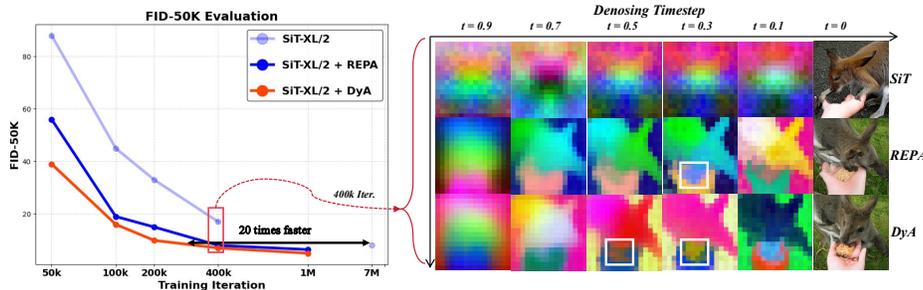


Figure 1: The **Dynamic Alignment** method has achieved further improvements compared to the REPA method, achieves convergence that is about  $20.0\times$  faster than the SiT. From the visualizations of the student block representations in the target model, the features of detail emerge at earlier denoising steps and are clearly disentangled from that of the background. DyA produces significantly sharper and more accurate details than both the SiT and REPA baselines, corroborating the quantitative improvements.

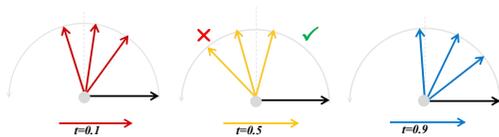
## 1 INTRODUCTION

Diffusion-based generative models (Ho et al., 2020; Song et al., 2020) and flow-based models (Albergo & Vanden-Eijnden, 2023; Lipman et al., 2023; Liu et al., 2023) have significant progress in generation task (Podell et al., 2023; Saharia et al., 2022a; Esser et al., 2024b; Min et al., 2023). Recent works have also explored the use of diffusion models as representation learners and produce discriminative features (Li et al., 2023; Xiang et al., 2023; Chen et al., 2024; Mukhopadhyay et al., 2023). Methods (Leng et al., 2025; Tian et al., 2025; Yu et al., 2025) use alignment technique between the middle layer of a diffusion model and representation from a pretrained teacher model, effectively accelerating model training and enhancing the quality of generated outputs. However,

054 the simple representations from the teacher model may become a hindrance to the model’s genera-  
 055 tive capacity training in the later stages. Seeing alignment as a variant of knowledge distillation(Fan  
 056 et al., 2024; Zhao et al., 2022) makes the problem clearer. This phenomenon can be attributed to **ca-**  
 057 **capacity mismatch**, which arises because the representations from the teacher model are too simplistic  
 058 to the complexity of the generative task of student model.

059 Essentially, pre-training visual encoders compress images into discriminative high-level repre-  
 060 sentations; their objective rewards invariance—robustness to translation, illumination, and occlu-  
 061 sion—while actively discarding task-irrelevant details. In contrast, diffusion models must recon-  
 062 struct the original image with pixel-level fidelity, requiring them to model the full high-order statis-  
 063 tics of the data distribution. Diffusion models retain copious high-frequency information through  
 064 multi-scale skip connections, yielding “draft-like” redundant features, whereas pre-trained features  
 065 behave like “summaries” that maximize inter-class separability in a low-dimensional semantic space.  
 066 The divergence between the two feature types is rooted in their objectives: discriminative modeling  
 067 seeks local invariance and semantic compression, whereas generative modeling seeks global fidelity  
 068 and detail preservation.

069 Drawing on previous studies that analyzed  
 070 model training from a gradient perspective  
 071 (Guo et al., 2024; Wang et al., 2025), the impact  
 072 of REPA loss on the model’s generative  
 073 capabilities varies over time. As shown in Fig-  
 074 ure 6, the gradient of auxiliary loss has almost  
 075 no effect on generative capabilities at the initial  
 076 training stage. As training progresses, it even  
 077 have a negative impact.



078 Figure 2: Black arrows indicate the gradient di-  
 079 rection of the Denoising Loss, and colored repre-  
 080 sent that of the Auxiliary Loss. The angle between  
 081 them increases as the training stage advances.

082 Suppose we redefine the core task as leveraging  
 083 simple representation to guide complex representations, rather than a generation problem, inspired  
 084 by the successes of knowledge distillation(Huang et al., 2022) and representation alignment based  
 085 generation(Yu et al., 2025; Wang et al., 2025), we recast the above issue as two sub-problems under  
 086 a unified generative framework. On the one hand, the complexity of the teacher representations  
 087 needs to match the complexity of the generative task more closely to prevent the issue of capacity  
 088 mismatch. On the other hand, the training objective of the model should emphasize the training of  
 089 generative capabilities, treating the alignment of representations merely as an auxiliary tool.

090 We propose Dynamic Alignment (DyA) to address these problems. First, the DyA introduces multi-  
 091 level information to address the issue of hindered generative capacity caused by overly simplistic  
 092 alignment targets in former methods.

093 Additionally, the DyA implement a temporal module, which uses the denoising time of the diffu-  
 094 sion model as an indicator variable to adjust the alignment constraints intensity of different level  
 095 information, bringing greater flexibility to the training process.

096 Second, we introduce the Stochastic Dropout Strategy (SDS), which addresses the second sub-  
 097 problem from the perspective of training methodology. In the SDS, the model halts alignment  
 098 behavior in every training iteration with a pre-set probability  $p$ , retaining only the denoising ob-  
 099 jective tied to generative capacity in the loss term. Unlike the truncation strategy (Wang et al.,  
 100 2025), which stops alignment from a pre-set step, SDS prioritizes generative capacity while contin-  
 101 uously guiding the model throughout whole training process. Experiments show that SDS brings a  
 102 greater enhancement to the model’s generative quality.

103 Our main contributions are:

- 104 • We propose Dynamic Alignment (DyA), a multi-level guidance strategy that dynamically  
 105 modulates alignment strength using diffusion timestep, mitigating representation collapse  
 106 in later training stages.
- 107 • We introduce Stochastic Dropout Strategy (SDS), a probabilistic training mechanism that  
 encourages generative capacity focus while retaining auxiliary guidance adaptively.
- Extensive experiments on ImageNet and ArtBench demonstrate  $20\times$  acceleration and su-  
 perior generation quality, surpassing the REPA and vanilla SiT baselines.

## 2 RELATED WORK

Recent developments in diffusion models have largely revolved around enhancements in learning methods, sampling strategies(Song et al., 2022; Song & Ermon, 2019; Lu et al., 2022; 2023), guidance techniques(Ho & Salimans, 2022; Nichol et al., 2022), latent representations(Rombach et al., 2022), and overall model structures(Ho et al., 2022; Peebles & Xie, 2023; Saharia et al., 2022b; Xue et al., 2023). Notably, innovations like DiT and U-ViT(Peebles & Xie, 2023; Bao et al., 2023) have introduced transformer-based architectures as alternatives or enhancements to the conventional U-Net. Transformer-based architectures introduce the attention mechanism into the diffusion framework, which endows the model with a more efficient capability to process global information. These advancements have significantly influenced the design of state-of-the-art image(Chen et al., 2023) and video(Gupta et al., 2024) synthesis systems, exemplified by Stable Diffusion 3.0(Esser et al., 2024a).

Many studies(Ye et al., 2025; Guo et al., 2025) have introduced or emphasized temporal information features in model training, thereby significantly enhancing the performance of the models. However, the focus is often placed on the temporal characteristics of external information, while the internal temporal features inherent to the model itself have not been sufficiently explored. Inspired by the adaLN module in DiT(Peebles & Xie, 2023), we introduce the denoising time steps within the diffusion model as indicator variables to dynamically adjust the constraint strength of different levels of information on the model at different times.

## 3 DYNAMIC ALIGNMENT AND STOCHASTIC DROPOUT STRATEGY

In this section, we will first present the formal definition of the task. Then, we will introduce the overall framework of DyA. Finally, we will provide detailed explanations of each component within the Dynamic Alignment and the Stochastic Dropout Strategy.

### 3.1 TASK DEFINITION

We provide a concise overview of flow-based and diffusion-based models from the unified viewpoint of stochastic interpolants(Albergo et al., 2023; Ma et al., 2024).

We consider  $p(x)$  to be an unknown distribution of data  $x \in X$ . We use a model distribution to approximate  $p(x)$ , with a dataset drawn from  $p(x)$ . We define our task as learning a latent distribution  $p(E(x))$  through a diffusion model, where  $E$  represents an encoder from a pretrained auto-encoder(Rombach et al., 2022), with  $x \sim p_{data}(x)$ .

We consider a continuous time-dependent process with a data  $x_* \sim p(x)$  and a Gaussian noise  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  on  $t \in [0, T]$ :

$$\mathbf{x}_t = \alpha_t \mathbf{x}_* + \beta_t \epsilon, \alpha_0 = \beta_T = 1, \alpha_T = \beta_0 = 0, \quad (1)$$

where  $\alpha_t$  is a decreasing function of  $t$  and  $\beta_t$  is an increasing one. Considering the described process, there exists a Probability Flow Ordinary Differential Equation characterized by a velocity field.

$$\dot{\mathbf{x}}_t = \mathbf{v}(\mathbf{x}_t, t), \quad (2)$$

where the distribution of this ODE at  $t$  can be seen as the marginal  $p_t(x)$ . Target can be sampled through the solution of Eq.(2) through existing ODE samplers starting from a random Gaussian noise  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ (Ma et al., 2024; Lipman et al., 2023).

This velocity  $v(x, t)$  can be described through the following sum of two conditional expectations

$$\mathbf{v}(\mathbf{x}, t) = \mathbb{E}[\dot{\mathbf{x}}_t | \mathbf{x}] = \dot{\alpha}_t \mathbb{E}[\mathbf{x}_* | \mathbf{x}] + \dot{\sigma}_t \mathbb{E}[\epsilon | \mathbf{x}], \quad (3)$$

which can be approximated by training the model  $v_\theta(x_t, t)$  to minimize the following training objective:

$$\mathcal{L}_{\text{velocity}}(\theta) := \mathbb{E}_{\mathbf{x}_*, \epsilon, t} \left[ \|\mathbf{v}_\theta(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x}_* - \dot{\sigma}_t \epsilon\|^2 \right], \quad (4)$$

Following (Ma et al., 2024), we primarily employ a straightforward linear interpolant with restricting  $T = 1 : \alpha_t = 1 - t$  and  $\beta_t = t$ .

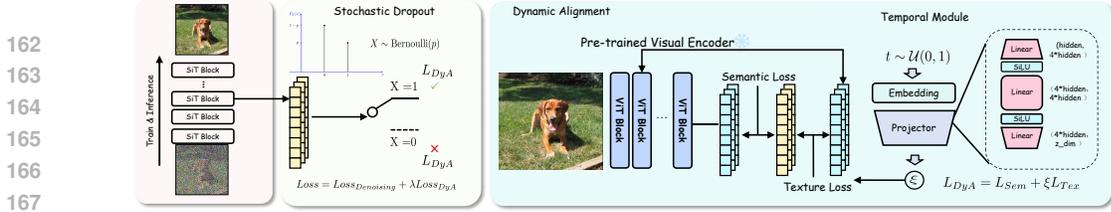


Figure 3: DyA employs visual representations at different levels to guide the training of the generation process and adjusts the guidance strength based on the denoising time as an indicator variable. This guidance is further regulated by the SDS strategy: at every training step, a Bernoulli sample decides whether to enable the DyA loss, thereby preventing over-alignment.

### 3.2 FRAMEWORK

As shown in Figure 3 DyA extracts the output of a certain layer of the diffusion model for alignment. After the output is resized by Multilayer Perceptron Layer, it is then compared with pre-trained visual representations of different level to calculate similarity and use them as Semantic Loss and Texture Loss respectively. The strength of the Texture Loss is dynamically adjusted based on the denoising time  $t$  as an indicator, with its intensity varying according to the changes in denoising time.

### 3.3 DYNAMIC ALIGNMENT

Simple teacher representations for distillation can lead to capacity mismatch issues in the later stages of training. The pre-trained model provides a clear guidance for the information flow in generation tasks, but due to the inherent differences in tasks, the information contained in the generation task is far more complex than the guidance provided by the pre-trained model. After a certain stage of training, the simple guidance provided by the pre-trained model can suppress the training of generative capabilities.

Therefore, we adopt a very concise and straightforward approach to address this issue: increasing the richness of the guidance representations. We extract information from the shallow layer of the encoder to serve as the alignment target. Providing the model with more complex guidance but also alleviates the issue of over-alignment during training to a certain extent.

As can be seen from Figure 4, the shallow representations show a more distinct separation between objects with different semantics in adjacent areas. In contrast, deeper representations group them more roughly into one category. However, the even shallower representations are too vague to serve as a guiding factor for convergence. Such observations are consistent with the results of the following ablation experiments on the guidance layer.

Subsequently, we introduce the denoising time steps of the diffusion model as an indicator variable, dynamically adjusting the constraint strength of different levels of information on the model according to the varying denoising times. The mapping from temporal information to variables is accomplished by the **Temporal Module**, taking  $t \in [0, 1]$  as input and returning a parameter  $\xi \in [0, 1]$ .

### 3.4 STOCHASTIC DROPOUT STRATEGY

To further address the issue of model over-alignment, we propose the Stochastic Dropout Strategy (SDS), a training scheme that employs probabilistic dropout of the alignment method. To be specific, SDS will cease to align representation with a probability  $p$  (where  $p$  is a pre-set probability value) in a single training iteration. As

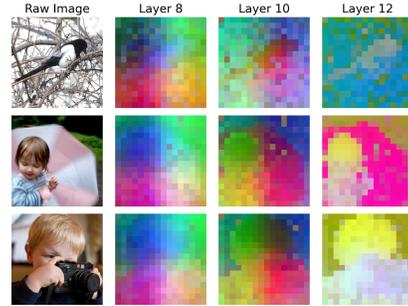


Figure 4: Visualization of intermediate representations of pre-trained visual encoder.(Oquab et al., 2023).

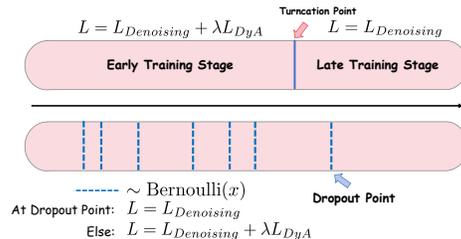


Figure 5: Visualization of SDS and Truncation.

a result, the model maintains the denoising objective as the sole loss value in that iteration, while remaining unchanged in other circumstances.

Compared to truncation-based approaches, the Stochastic Dropout Strategy (SDS) not only allows the model to emphasize its generative capabilities but also provides guidance throughout the entire training process. Experimental results demonstrate that the SDS approach leads to significantly better model performance.

### 3.5 LOSS DESIGN

Considering  $E$  to be a pretrained encoder and  $x_*$  a clean image. Let  $z_*^k = E(x_*) \in \mathbb{R}^{N \times D}$  to be the output of the encoder, where  $N, D > 0$  are the number of patches and the embedding dimension of  $E$  respectively, and  $k$  stands for the index of the output layer of encoder. DyA aligns  $h_\phi(\mathbf{h}_t) \in \mathbb{R}^{N \times D}$  with  $z_*^k$  and  $z_*^{final}$ , where  $h_\phi(\mathbf{h}_t)$  represents the intermediate output of a diffusion transformer at a specific time step  $t$ .

DyA achieves alignment through a maximization of patch-wise similarities between texture representation  $z_*^k$ , semantic representation  $z_*^{final}$  and the hidden state  $\mathbf{h}_t$ , with  $\gamma$  as a time-related variable used for adjusting the regularization intensity of texture:

$$\mathcal{L}_{\text{DyA}}(\theta, \phi, t) := -\mathbb{E}_{\mathbf{x}_*, \epsilon, t} [\xi A(z_*^k, h_\phi(\mathbf{h}_t) + A(z_*^{final}, h_\phi(\mathbf{h}_t))], \quad (5)$$

where  $A$  means alignment and

$$\xi = T_\sigma(t), \quad (6)$$

$T_\sigma$  is a temporal module that takes a time embedding as input and produces the variable  $\xi$  via a linear transformation. The first and second alignment equal to

$$A(z_*^k, h_\phi(\mathbf{h}_t)) := \frac{1}{N} \sum_{n=1}^N \text{sim}(z_*^{k[n]}, h_\phi(\mathbf{h}_t^{[n]})), \quad (7)$$

$$A(z_*^{final}, h_\phi(\mathbf{h}_t)) := \frac{1}{N} \sum_{n=1}^N \text{sim}(z_*^{final[n]}, h_\phi(\mathbf{h}_t^{[n]})), \quad (8)$$

where  $n$  is a patch index and  $\text{sim}(\cdot, \cdot)$  is cosine similarity function.

We define a random variable  $D \sim \text{Bernoulli}(p)$  to indicate whether to use the DyA regularization in a particular iteration, where  $P(D = 1) = p$  and  $P(D = 0) = 1 - p$ .

We present the complete algorithm in Appendix A and ultimately obtain the DyA loss term. In practice, we add this term to the original diffusion-based objectives which give us the final loss function:

$$L := L_{\text{Diff}} + \lambda L_{\text{DyA}}, \quad (9)$$

where  $\lambda > 0$  is a hyperparameter that controls the trade off between denoising object and flexible regularization intensity noted as  $L_{\text{DyA}}$  in the above equation.

Table 1: Model configuration details.

Config	Layers	HiddenDim	Heads	Student Layer	Alignment Target	Params
B/2	12	768	12	L6	L9 + L11	142M
L/2	24	1024	16	L8	L9 + L11	470M
XL/2	28	1152	16	L8	L9 + L11	687M

## 4 EXPERIMENT

**Implement details.** We follow the setup in SiT(Ma et al., 2024) and REPA(Yu et al., 2025). We use ImageNet(Deng et al., 2009) and Artbench(Liao et al., 2022) to carry out our experiment, and each image is processed to resolution  $256 \times 256$ . We follow ADM(Dhariwal & Nichol, 2021) for data preprocessing protocols and evaluation metrics. Each image is encoded into a compressed vector  $z \in \mathbb{R}^{32 \times 32 \times 4}$  using Stable Diffusion VAE(Rombach et al., 2022). Model details can be seen in Table 1. Training the L size DyA model for 400K iterations with 2 A100 GPUs and in a batch size of 256 takes about 97 hours.

Table 2: FID comparisons with vanilla SiTs, DiTs and REPA on ImageNet 256x256. We do not use classifier-free guidance (CFG).  $\downarrow$  denotes lower values are better. Iter. indicates the training iteration.

Model	FLOPs	Iter.	FID $\downarrow$
DiT-B/2	23.0	400K	43.5
<b>+DyA(ours)</b>	23.0	400k	<b>34.7</b>
DiT-L/2	80.7	400K	23.3
+REPA	80.7	400K	15.6
<b>+DyA(ours)</b>	80.7	400K	<b>12.8</b>
SiT-B/2	21.8	400K	33.0
+REPA	21.8	400K	24.4
<b>+DyA(ours)</b>	21.8	400K	<b>22.6</b>
SiT-L/2	77.5	400K	18.8
+REPA	77.5	400K	10.0
<b>+DyA(ours)</b>	77.5	400K	<b>9.0</b>
SiT-XL/2	117.7	400K	17.2
+REPA	117.7	400K	7.9
<b>+DyA(ours)</b>	117.7	400K	<b>7.5</b>

#### 4.1 SETUP

**Evaluation metric.** We implement Frechet Inception Distance (Heusel et al., 2017), sFID (Nash et al., 2021), Inception Score (Salimans et al., 2016), Precision (Pre.) and Recall (Rec.)(Kynkäänniemi et al., 2019) on 50,000 samples.

**Sampler.** We follow SiT(Ma et al., 2024) and use the SDE Euler-Maruyama sampler (for SDE with  $\omega_t = \sigma_t$ ) and set the number of function evaluations (NFE) as 250 by default.

**Baselines.** We take several recent diffusion-based generation methods into consideration, each employing different inputs and network architectures. Four types of approaches are included: **Pixel diffusion**, **Latent diffusion with U-Net**, **Latent diffusion with transformer+U-Net hybrid models** and **Latent diffusion with transformers**.

#### 4.2 SYSTEM-LEVEL COMPARISON

We conduct a systematic comparison between recent state-of-the-art diffusion model approaches, diffusion transformers with REPA, and diffusion transformers with DyA. First, we compare the FID values between vanilla SiT, vanilla DiT, SiT with REPA, DiT with REPA and the same models trained with DyA. As shown in Table 2, DyA achieves comprehensive improvements over vanilla SiT, vanilla DiT and REPA. Specifically, SiT achieves FID of 17.2

at 400k iterations and with the help of REPA achieves FID of 7.9 at 400k iterations, while DyA reaches FID of 7.5 at same iterations. The model trained with DyA exhibits even better convergence speed than REPA. The experimental results are averaged over five independent runs. A Wilcoxon rank-sum test reveals that DyA significantly outperforms REPA in both FID and IS ( $p < 0.05$ , Rank-Biserial  $r = 0.6$ ), indicating a statistically meaningful difference.

We provide a quantitative evaluation comparing SiT-XL/2 with DyA to all model variants. As result can be seen from Table 3, DyA shows consistent and significant improvement. At 200 epochs, SiT-XL/2 with DyA achieves FID of 1.71 with a classifier-free guidance scale of  $\omega = 1.5$ , already better than all previous methods. Our method outperforms the original SiT-XL/2 with 7x fewer epochs and out performs REPA with 4x fewer epochs and it is further improved with longer training. As the number of model training epochs increases, the FID shows a continuous decline, further dropping to 1.59 at the 400 epochs. The Rec. value reached 0.65 at 600 epoch which is better than all previous

Table 3: System-level comparison on ImageNet 256x256 without CFG (classifier guidance). DyA use the 9th layer and the final layer as the alignment target and use  $\lambda = 1$ . The  $\downarrow$  and  $\uparrow$  indicate whether lower or higher values are better, respectively.

Model	Epochs	FID $\downarrow$	IS $\uparrow$	Rec. $\uparrow$
<i>Pixel diffusion</i>				
ADM-U	400	3.94	186.7	0.52
VDM++	560	2.40	225.3	-
Simple diffusion	800	2.77	211.8	-
CDM	2160	4.88	158.7	-
<i>Latent Diffusion, Unet</i>				
LDM-4	200	3.60	247.7	0.48
<i>Latent Diffusion, Transformer+U-Net hybrid</i>				
DiffT	-	1.73	276.5	0.62
U+ViT-H/2	240	2.29	263.9	0.57
<i>Latent Diffusion, Transformer</i>				
MaskDiT	1600	2.28	276.6	0.61
SD-DiT	480	3.23	-	-
DiT-XL/2	1400	2.27	278.2	0.57
SiT-XL/2	1400	2.06	270.3	0.59
+REPA	800	1.80	<b>284.0</b>	0.61
<b>+DyA(ours)</b>	200	<b>1.71</b>	254.9	<b>0.63</b>
<b>+DyA(ours)</b>	400	<b>1.59</b>	264.0	<b>0.64</b>
<b>+DyA(ours)</b>	600	<b>1.55</b>	272.4	<b>0.65</b>

Table 4: Comparison on Artbench 256x256 with no CFG. The  $\downarrow$  and  $\uparrow$  indicate whether lower or higher values are better, respectively.

SiT	sFID $\downarrow$	Pre. $\uparrow$	Rec. $\uparrow$
+REPA	29.68	0.55	0.35
+DyA	31.70	0.54	0.35
+DyA with Truncate	27.48	<b>0.61</b>	0.31
+DyA with SDS	<b>26.67</b>	0.60	<b>0.36</b>

Table 5: **Component-wise analysis** on ImageNet  $256 \times 256$ . All models are SiT-L/2 trained for 400K iterations. All metrics are measured with the SDE Euler-Maruyama sampler with NFE=250. We fix  $\lambda = 1$  here.  $\downarrow$  and  $\uparrow$  indicate whether lower or higher values are better, respectively.

Iter.	Time Modulate	Guidance	Dropout Rate	FID $\downarrow$	sFID $\downarrow$	IS $\uparrow$	Pre. $\uparrow$	Rec. $\uparrow$
400K	-	L11	-	10.0	5.34	111.9	0.68	0.65
400K	$1 - \xi, \xi$	+L7	-	10.4	5.20	106.9	0.68	0.65
400K	$\xi$	+L7	-	9.39	5.09	114.1	0.69	0.65
400K	$\xi$	+L5	-	9.65	5.15	112.2	0.68	0.65
400K	$\xi$	+L7	-	9.39	5.09	114.1	0.69	0.65
400K	$\xi$	+L9	-	9.00	5.17	116.1	0.69	0.65
400K	$\xi$	+L9	5%	9.23	5.19	114.8	0.69	0.65
400K	$\xi$	+L9	10%	9.18	5.28	116.1	0.68	0.66
400K	$\xi$	+L9	15%	9.19	5.31	115.8	0.68	0.65
400K	$\xi$	+L9	20%	9.08	5.17	116.0	0.69	0.65
400K	$\xi$	+L9	25%	8.96	5.56	117.4	0.70	0.65
400K	$\xi$	+L9	Truncation	17.44	6.66	80.7	0.58	0.68

method listed above. We can clearly see from the table that both the FID and IS scores improve significantly with the increase in training epochs.

To investigate the generalization capability of our model and its hyper-parameter choices across diverse datasets, we conduct training and evaluation on ArtBench. Training was performed for 50 000 steps with a batch size of 512, using SiT-B/2 as the backbone and comparing both REPA and the truncation training strategy. Our results demonstrate that DyA combined with the SDS training schedule achieves competitive performance.

### 4.3 ABLATION STUDY

We have experimented with two different ways of using the **Temporal Module**. In our initial attempt, we use  $\xi$ , the output of the temporal module, to adjust the regularization intensity of both semantic and texture information, where  $\xi \in [0, 1]$ . We set  $\xi$  and  $1 - \xi$  as the regularization intensities for the two types of information, respectively. As shown in Table 5, although the model’s performance improved in terms of sFID, its performance in IS and FID was lower than that of REPA. This indicates that the regularization intensity demands for different types of information during the denoising process are not simply inversely proportional. In the subsequent experiments, we adopt the latter approach.

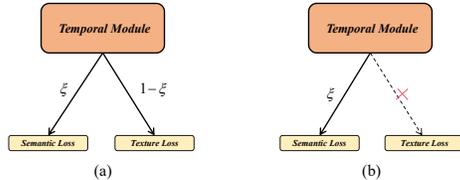


Figure 6: (a) represents the method in which temporal module regulate both semantic loss and texture loss. (b) shows the method that temporal module regulate texture loss only.

Then we conducted ablation studies on different **Guidance**. As result can be seen from Table 5, DyA with the modification outperforms REPA in all five metrics, with the improvement being particularly evident in the sFID score. From the data, we can observe that, at a fixed number of training steps, increasing the richness of guidance can significantly enhance the training efficiency of generative models. Additionally, as visualized in Figure 4, features that emerge too early are not suitable for use as generative guidance. A plausible explanation is that the image information has not been processed through a sufficient number of layers to produce convergent guiding representations.

From Table 5 it can be discovered that DyA achieves better evaluation score regardless of the layer choice of alignment target, and we attribute this phenomenon to the presence of temporal module. Additionally, the performance gains contributed by different layers are non-uniform. Therefore To preserve methodological simplicity and uphold the algorithmic essence—enhancing guidance complexity in the most straightforward manner—we ultimately employ double target layers.

We have evaluated the model with the **Stochastic Dropout Strategy** in different dropout rate. As shown in Table 5, we experiment with dropout rates of 5%, 10%, 15%, 20% and 25% for the L-sized model. Overall, we find that DyA still outperforms baseline. Moreover, compared with truncation strategy which ceases the alignment from a pre-set step, SDS also have a better performance. Com-

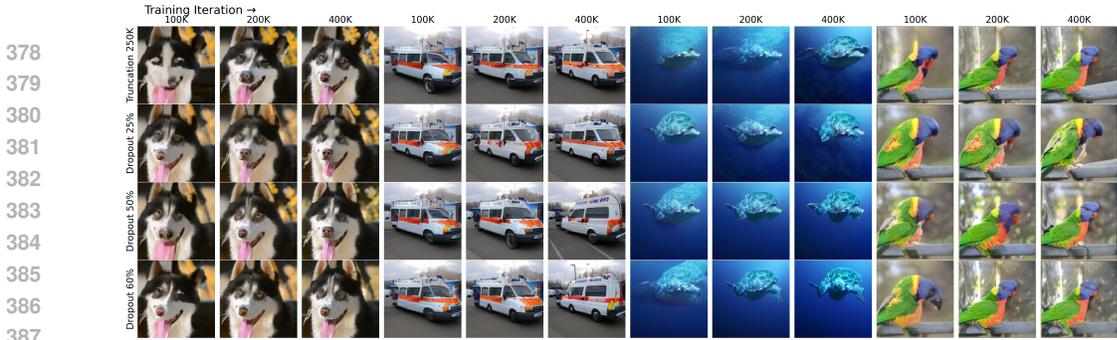


Figure 7: The impact of the SDS strategy on the generation effect. A high dropout rate can lead to insufficient pre-training guidance strength, thereby slowing down the convergence speed. Furthermore, the truncation method generally yields inferior generation quality before the truncation point compared to the SDS approach.

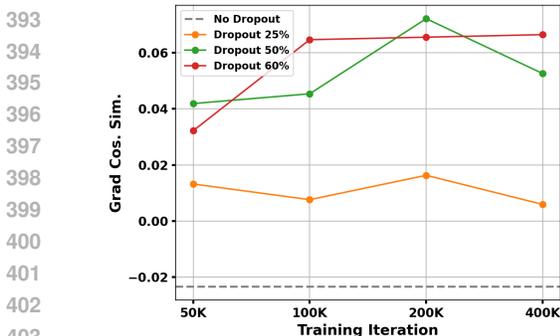


Figure 8: Gradient angle between Denoising Loss and Auxiliary Loss across training iteration.

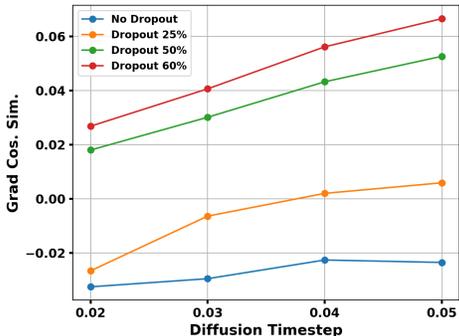


Figure 9: Gradient angle between Denoising Loss and Auxiliary Loss across denoising timesteps.

pared to the truncation approach, the guidance provided by SDS spans the entire training process. It also does not significantly affect the generation quality due to the setting of a too early or too late truncation point.

During training, we observed that the diffusion loss of the SDS-augmented model is consistently higher than its REPA counterpart in the early phase. Once the training step exceeds 30 000, the auxiliary loss decreases at a significantly faster rate. Such a situation indicates that, from a training dynamic perspective, the impact of SDS on the generative-related denoising loss is superior to full-process alignment, effectively alleviating the issue of over-alignment.

Furthermore, from Figure 7 the truncation method generally yields inferior generation quality before the truncation point compared to the SDS approach. This indicates that issues of over-alignment and capacity mismatch do not solely emerge in the later stages of training.

Thus, it can be concluded that the underlying assumption of truncation, which posits that over-alignment and capacity mismatch occur only in the later stages of training, is erroneous. In contrast, the SDS strategy does not rely on such a premise. Instead, the SDS approach attributes the problem to the inherent differences in the tasks of various models, adjusting guidance throughout the entire training process. This effectively enhances both the quality of generation and the efficiency of training.

From a gradient perspective, as can be seen from Figure 9 and Figure 8, A higher dropout rate has a more pronounced effect in mitigating the issue of over-alignment. At dropout rates of 50% and 60%, the gradients of the auxiliary loss and the denoising loss remain positive throughout the entire training process, indicating that the guidance consistently contributes to positive optimization of the generation effect. However, from the perspective of the generation effect, the positive cosine similarity is actually due to "insufficient alignment." As can be seen from Figure 7, the images with a 25% dropout rate clearly show more mature image generation at earlier training steps.

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446

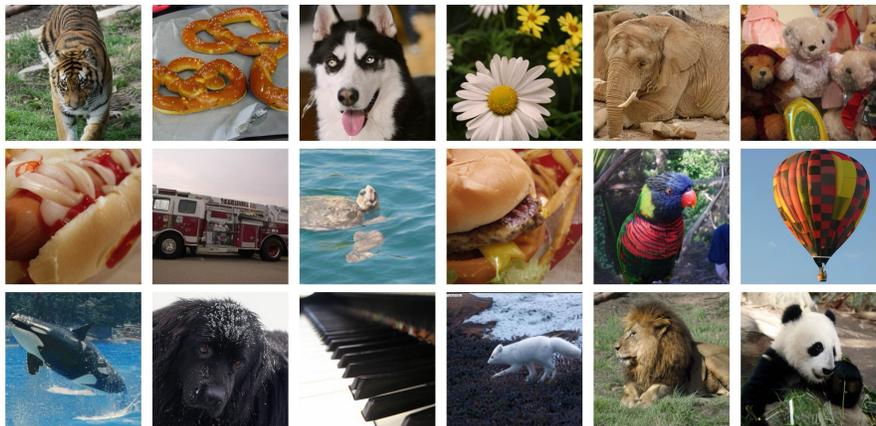


Figure 10: Selected samples from models trained on ImageNet 256 ×256 resolution with  $cfg = 4.0$ .

447  
448  
449  
450  
451  
452  
453  
454

Using the SDS strategy requires a balance between the positive optimization brought by alignment and the actual generation results. Through experimentation, we found that a dropout rate of 25% provides the greatest enhancement to generation quality, while also ensuring that the alignment strategy continues to provide positive guidance for the generation capability training throughout the long-term training process.

455  
456  
457  
458  
459  
460

From Figure 9, we can see SDS actually alleviate the over-alignment issue. From denoising time 0.02 to 0.05, the cosine similarity between gradient of the auxiliary loss and the gradient of the denoising loss are negative without the use of the SDS strategy. With the implementation of the SDS strategy, the cosine similarity during this time period is almost entirely positive. As mentioned above, considering the dual aspects of balancing guidance for positive optimization and generation quality, we ultimately adopted a dropout rate of 25% as the final solution.

461  
462  
463  
464  
465  
466

Table 6: Ablation study for Stu. Layer .

Table 7: Ablation study for  $\lambda$ .

Stu. Layer	Layer6	Layer8	Layer10	$\lambda$	0.25	0.5	0.75	1
FID↓	22.6	26.4	34.4	FID↓	10.21	10.10	9.79	<b>9.39</b>
IS↑	65.44	57.91	46.30	IS↑	107.1	108.8	111.5	<b>114.1</b>

467  
468  
469

We investigate the impact of different **Student Layers** on model performance. As shown in Table 6, we find that the model’s performance actually shows a decline as the relative depth of the alignment layer increases, consistent with the findings in the REPA.

470  
471  
472  
473

We examine the effect of the regularization coefficient  $\lambda$  by training DyA with different coefficients 0.5 to 1.0 and comparing the FID and IS. As shown in Table 7, the performance have stability over different values and peak at  $\lambda = 1$ . Although performance exhibits a monotonic rise with increasing  $\lambda$ , we ultimately cap  $\lambda$  at 1 to keep the optimization focus unambiguously on generation quality.

474  
475  
476

## 5 CONCLUSION

477  
478  
479  
480  
481  
482

In this paper, we presented Dynamic Alignment. We introduced multi-level information coupled with temporal module to address the capacity mismatch issue in conventional representation-guided generation. Building upon this, Stochastic Dropout Strategy was further employed to mitigate over-alignment issue. Systematic comparisons on ImageNet and ArtBench demonstrated that DyA significantly enhances the generation quality of model and accelerate training convergence.

483  
484  
485

We empirically observed that augmenting representation richness together with the SDS strategy consistently elevates generation quality. And semantic structures are able to emerge markedly earlier in the target layer. These findings corroborate the validity of our methodological framework. We will endeavor to transcend the constraints of pre-trained models and further refine this paradigm.

486 ETHICS STATEMENT  
487

488 Our work investigates diffusion models for generative modeling. We believe this research has the  
489 potential to contribute positively by enabling creative applications, advancing understanding of prob-  
490 abilistic modeling, and improving efficiency in downstream scientific and industrial tasks.

491 However, we acknowledge possible negative impacts. The proposed method could be misused for  
492 generating misleading, harmful, or otherwise inappropriate content. Our experiments are conducted  
493 exclusively on publicly available benchmark datasets, which do not contain personally identifiable  
494 information. Nevertheless, any biases present in the datasets may be propagated or amplified by our  
495 models. We encourage further work on bias detection and mitigation in generative modeling.

496 Regarding environmental impact, our experiments were performed on  $4 \times A100$  GPUs. We recognize  
497 the importance of efficient model design and responsible use of computational resources in order to  
498 reduce the carbon footprint of large-scale model training.

499 Overall, we emphasize that our research should be used for beneficial purposes only, and we dis-  
500 courage applications that may cause harm to individuals or society.  
501

502 REPRODUCIBILITY STATEMENT  
503

504 We are committed to ensuring the reproducibility of our work. All datasets used in our experiments  
505 are publicly available. We will release our code, along with detailed instructions for training and  
506 evaluation.

507 We describe all necessary implementation details in the main papersupplementary material, and  
508 appendix, including model architectures, optimization settings (learning rate, batch size, optimizer,  
509 scheduler), and data preprocessing steps. We set fixed random seeds for all experiments to ensure  
510 consistent results across runs.

511 Our experiments were conducted on  $4 \times A100$  GPUs.

512 All reported results in tables and figures can be reproduced using the released code and configuration  
513 files. We also provide scripts to regenerate the main figures and evaluation metrics directly from  
514 trained checkpoints.  
515

516 REFERENCES  
517

518 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic inter-  
519 polants. In *11th International Conference on Learning Representations, 2023*, 2023.

520 Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A uni-  
521 fying framework for flows and diffusions, 2023. URL <https://arxiv.org/abs/2303.08797>.

522 Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth  
523 words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on*  
524 *computer vision and pattern recognition*, pp. 22669–22679, 2023.

525 Junsong Chen, YU Jincheng, GE Chongjian, Lewei Yao, Enze Xie, Zhongdao Wang, James Kwok,  
526 Ping Luo, Huchuan Lu, and Zhenguo Li. Fast training of diffusion transformer for photorealistic  
527 text-to-image synthesis. In *The Twelfth International Conference on Learning Representations*,  
528 2023.

529 Xinlei Chen, Zhuang Liu, Saining Xie, and Kaiming He. Deconstructing denoising diffusion models  
530 for self-supervised learning, 2024.  
531

532 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-  
533 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
534 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.

535 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
536 *in neural information processing systems*, 34:8780–8794, 2021.  
537  
538  
539

- 540 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam  
541 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers  
542 for high-resolution image synthesis. In *Forty-first international conference on machine learning*,  
543 2024a.
- 544 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam  
545 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English,  
546 Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow trans-  
547 formers for high-resolution image synthesis. In *Forty-first international conference on machine*  
548 *learning*, 2024b.
- 549 Wen-Shu Fan, Xin-Chun Li, and De-Chuan Zhan. Exploring dark knowledge under various teacher  
550 capacities and addressing capacity mismatch. *arXiv preprint arXiv:2405.13078*, 2024.
- 551 Yingqing Guo, Hui Yuan, Yukang Yang, Minshuo Chen, and Mengdi Wang. Gradient guidance  
552 for diffusion models: An optimization perspective. *Advances in Neural Information Processing*  
553 *Systems*, 37:90736–90770, 2024.
- 554 Yongxin Guo, Jingyu Liu, Mingda Li, Dingxin Cheng, Xiaoying Tang, Dianbo Sui, Qingbin Liu,  
555 Xi Chen, and Kevin Zhao. Vtg-llm: Integrating timestamp knowledge into video llms for en-  
556 hanced video temporal grounding. In *Proceedings of the AAAI Conference on Artificial Intelli-*  
557 *gence*, volume 39, pp. 3302–3310, 2025.
- 558 Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Fei-Fei Li, Irfan Essa, Lu Jiang, and  
559 José Lezama. Photorealistic video generation with diffusion models. In *European Conference on*  
560 *Computer Vision*, pp. 393–411. Springer, 2024.
- 561 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
562 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*  
563 *neural information processing systems*, 30, 2017.
- 564 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on*  
565 *Deep Generative Models and Downstream Applications*, 2022.
- 566 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
567 *neural information processing systems*, 33:6840–6851, 2020.
- 568 Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Sali-  
569 mans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning*  
570 *Research*, 23(47):1–33, 2022.
- 571 Tao Huang, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge distillation from a stronger  
572 teacher. *Advances in Neural Information Processing Systems*, 35:33716–33727, 2022.
- 573 Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved  
574 precision and recall metric for assessing generative models. *Advances in neural information*  
575 *processing systems*, 32, 2019.
- 576 Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng.  
577 Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint*  
578 *arXiv:2504.10483*, 2025.
- 579 Alexander C. Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your dif-  
580 fusion model is secretly a zero-shot classifier. In *Proceedings of the IEEE/CVF International*  
581 *Conference on Computer Vision*, pp. 2206–2217, 2023.
- 582 Peiyuan Liao, Xiuyu Li, Xihui Liu, and Kurt Keutzer. The artbench dataset: Benchmarking genera-  
583 tive models with artworks. *arXiv preprint arXiv:2206.11404*, 2022.
- 584 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
585 for generative modeling. In *11th International Conference on Learning Representations*, 2023,  
586 2023.

- 594 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and  
595 transfer data with rectified flow. In *The Eleventh International Conference on Learning Repre-*  
596 *sentations*, 2023.
- 597 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast  
598 ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural*  
599 *Information Processing Systems*, 35:5775–5787, 2022.
- 600 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast  
601 solver for guided sampling of diffusion probabilistic models, 2023. URL <https://arxiv.org/abs/2211.01095>.
- 602 Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Sain-  
603 ing Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant  
604 transformers. In *European Conference on Computer Vision*, pp. 23–40. Springer, 2024.
- 605 Lejun Min, Junyan Jiang, Gus Xia, and Jingwei Zhao. Polyffusion: A diffusion model for poly-  
606 phonic score generation with internal and external controls. In *ISMIR 2023 Hybrid Conference*,  
607 2023.
- 608 Soumik Mukhopadhyay, Matthew Gwilliam, Vatsal Agarwal, Namitha Padmanabhan, Archana  
609 Swaminathan, Srinidhi Hegde, Tianyi Zhou, and Abhinav Shrivastava. Diffusion models beat  
610 gans on image classification, 2023. URL <https://arxiv.org/abs/2307.08702>.
- 611 Charlie Nash, Jacob Menick, Sander Dieleman, and Peter Battaglia. Generating images with sparse  
612 representations. In *International Conference on Machine Learning*, pp. 7958–7968. PMLR, 2021.
- 613 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
614 Mcgrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and  
615 editing with text-guided diffusion models. In *International Conference on Machine Learning*, pp.  
616 16784–16804. PMLR, 2022.
- 617 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
618 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning  
619 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 620 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
621 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 622 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
623 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
624 synthesis. In *The Twelfth International Conference on Learning Representations*, 2023.
- 625 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
626 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
627 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 628 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-  
629 yar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Sal-  
630 imans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image dif-  
631 fusion models with deep language understanding. *Advances in neural information processing*  
632 *systems*, 35:36479–36494, 2022a.
- 633 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
634 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
635 text-to-image diffusion models with deep language understanding. *Advances in neural informa-*  
636 *tion processing systems*, 35:36479–36494, 2022b.
- 637 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
638 Improved techniques for training gans. *Advances in neural information processing systems*, 29,  
639 2016.
- 640 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Interna-*  
641 *tional Conference on Learning Representations*, 2022.

- 648 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.  
649 *Advances in neural information processing systems*, 32, 2019.  
650
- 651 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
652 Poole. Score-based generative modeling through stochastic differential equations. In *International  
653 Conference on Learning Representations*, 2020.
- 654 Yuchuan Tian, Hanting Chen, Mengyu Zheng, Yuchen Liang, Chao Xu, and Yunhe Wang. U-repa:  
655 Aligning diffusion u-nets to vits. *arXiv preprint arXiv:2503.18414*, 2025.  
656
- 657 Ziqiao Wang, Wangbo Zhao, Yuhao Zhou, Zekai Li, Zhiyuan Liang, Mingjia Shi, Xuanlei Zhao,  
658 Pengfei Zhou, Kaipeng Zhang, Zhangyang Wang, et al. Repa works until it doesn't: Early-  
659 stopped, holistic alignment supercharges diffusion training. *arXiv preprint arXiv:2505.16792*,  
660 2025.
- 661 Weilai Xiang, Hongyu Yang, Di Huang, and Yunhong Wang. Denoising diffusion autoencoders are  
662 unified self-supervised learners. In *Proceedings of the IEEE/CVF International Conference on  
663 Computer Vision*, pp. 15802–15812, 2023.
- 664 Zeyue Xue, Guanglu Song, Qiushan Guo, Boxiao Liu, Zhuofan Zong, Yu Liu, and Ping Luo.  
665 Raphael: Text-to-image generation via large mixture of diffusion paths. *Advances in Neural  
666 Information Processing Systems*, 36:41693–41706, 2023.  
667
- 668 Jinhui Ye, Zihan Wang, Haosen Sun, Keshigeyan Chandrasegaran, Zane Durante, Cristobal Eyza-  
669 guirre, Yonatan Bisk, Juan Carlos Niebles, Ehsan Adeli, Li Fei-Fei, et al. Re-thinking temporal  
670 search for long-form video understanding. In *Proceedings of the Computer Vision and Pattern  
671 Recognition Conference*, pp. 8579–8591, 2025.
- 672 Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and  
673 Saining Xie. Representation alignment for generation: Training diffusion transformers is easier  
674 than you think. In *The Twelfth International Conference on Learning Representations*, 2025.  
675
- 676 Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation.  
677 In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pp.  
678 11953–11962, 2022.  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A ALGORITHM

---

### Algorithm 1 DyA with SDS.

---

**Input:** probability, batchsize, representations, timestep and target;

**Output:**  $L_{DyA}$ ;

```

semantic, texture = target;
Dropout = random() > probability;
LDyA = 0;
if not Dropout then
   $\xi = \text{TemporalModule}(\text{timestep})$ 
  for representation in representations: do
     $\text{loss}_{sem} = 0, \text{loss}_{tex} = 0;$ 
    for batch in representation do
       $\text{loss}_{sem} += \text{sim}(\text{batch}, \text{semantic});$ 
       $\text{loss}_{tex} += \text{sim}(\text{batch}, \text{texture});$ 
       $L_{DyA} += (\text{loss}_{sem} + \xi * \text{loss}_{tex}) / \text{batchsize};$ 
    end for
  end if
return  $L_{DyA}$ 

```

---

## B ACKNOWLEDGEMENTS

We used large language models (LLMs) such as ChatGPT for minor editing support, including grammar checking and language polishing. No LLMs were used for generating technical content, experimental results, or ideas in this work.