

ControlNeXt: Powerful and Efficient Control for Image and Video Generation

Bohao Peng¹ Jian Wang¹ Yuechen Zhang¹ Wenbo Li¹
 Ming-Chang Yang¹ Jiaya Jia^{1,2}

¹CUHK ²SmartMore

<https://github.com/dvlab-research/ControlNeXt>

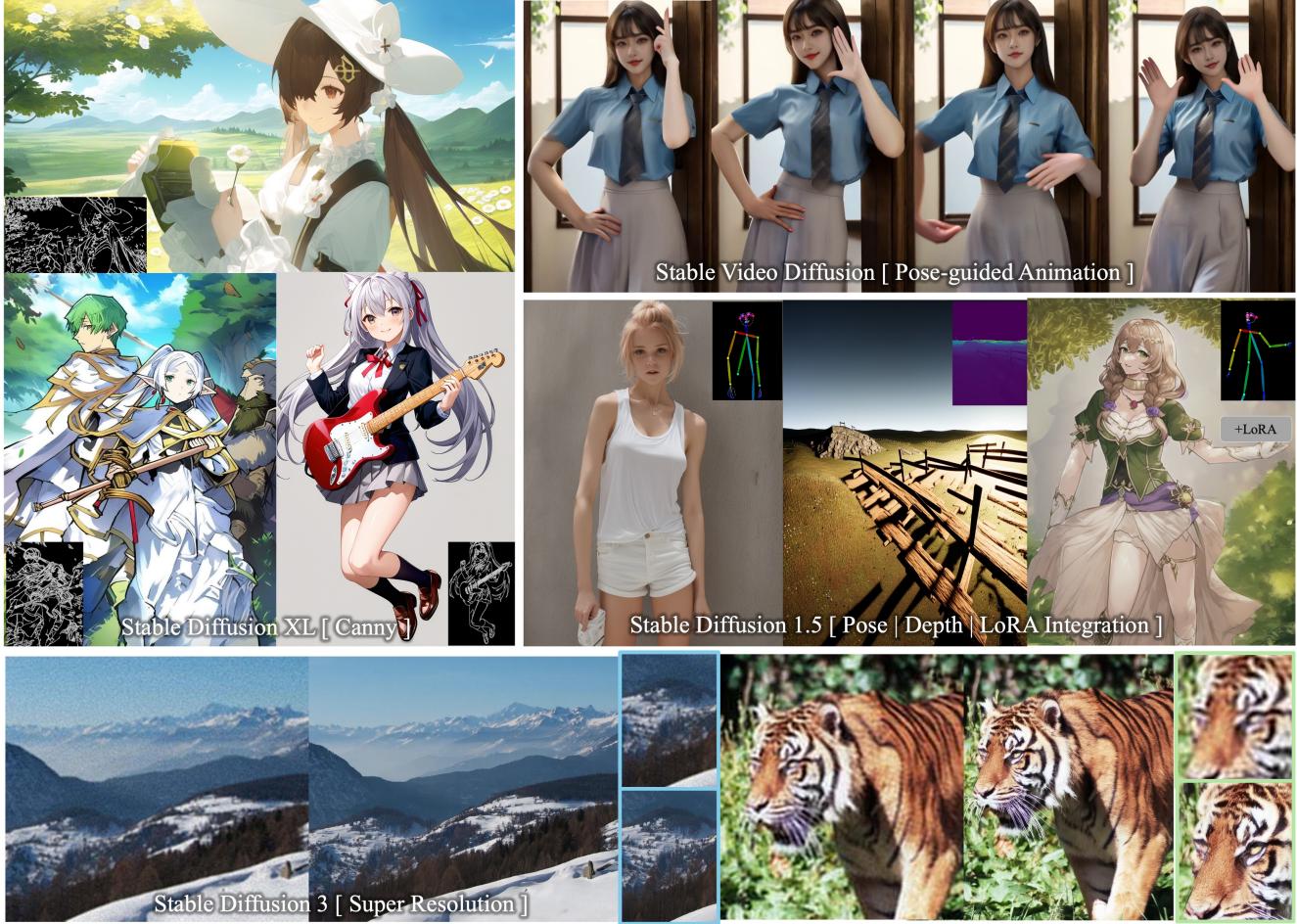


Figure 1. ControlNeXt is a powerful yet efficient method for controllable generation, emphasizing improved efficiency and robustness. For more examples, please refer to our project page: <https://pbihao.github.io/projects/controlnext/index.html>

Abstract

Diffusion models have demonstrated remarkable and robust abilities in both image and video generation. To achieve greater control over generated results, researchers

introduce additional architectures, such as ControlNet, Adapters and ReferenceNet, to integrate conditioning controls. However, current controllable generation methods often require substantial additional computational resources, especially for video generation, and face challenges in

training or exhibit weak control. In this paper, we propose ControlNeXt: a powerful and efficient method for controllable image and video generation. We first design a more straightforward and efficient architecture, replacing heavy additional branches with minimal additional cost compared to the base model. Such a concise structure also allows our method to seamlessly integrate with other LoRA weights, enabling style alteration without the need for additional training. As for training, we reduce up to 90% of learnable parameters compared to the alternatives. Furthermore, we propose another method called Cross Normalization (CN) as a replacement for “zero-convolution” to achieve fast and stable training convergence. We have conducted various experiments with different base models across images and videos, demonstrating the robustness of our method.

1. Introduction

Diffusion models generate data by iteratively transforming an initial simple distribution into a complex, high-quality distribution. This process allows diffusion models to produce realistic and high-fidelity results across various generative tasks, including both the images and videos [4, 8, 10, 23, 34]. However, relying solely on base generation models makes it challenging to achieve precise results and necessitates iterative adjustments of prompts and seeds, incurring significant costs. To achieve finer-grained controllable generation, pioneering works [5, 19, 44] recently introduce additional controls, such as depth, human pose skeletons, and edge maps, *etc.*, as the guidance.

Current controllable generation methods typically incorporate control information by adding a parallel branch or adapter to process and inject extra conditions, as exemplified by ControlNet [41, 44], T2I-Adapter [27] and ReferenceNet [16]. These architectures process auxiliary controls in parallel with the denoising main branch to extract fine-grained features. It then incorporates conditional controls to guide the denoising process, utilizing techniques such as *zero convolution* and *cross attention*.

Unfortunately, such operation typically leads to a significant increase in computational cost and training challenges. It can at most double the GPU memory consumption and require to introduce a massive number of new parameters for training. This issue is particularly serious for video generation, which repeatedly processes each individual frame. The T2I-Adapter attempts to reduce costs but sacrifices control abilities, making it unsuitable for a wide range of tasks, such as video generation or low-level tasks. Furthermore, the introduction of *zero convolution* also increases training challenges, slows convergence, and results in the “sudden convergence phenomenon” [44].

This paper presents ControlNeXt, a powerful yet efficient method for controllable visual generation, emphasizing

improved efficiency and robustness. We present its powerful capabilities in Fig. 1, demonstrating that it can be applied to various files and tasks across different backbones. We first delve into refining the architecture design of the ControlNet to provide a more compatible and straightforward structure. Typically, controllable generation involves fine-tuning on small-scale datasets to enable pre-trained large diffusion models [4, 31, 34, 47] to effectively capture and utilize conditional controls. Fine-tuning datasets are usually much smaller compared to the original pre-training datasets, such as LAION-5B [36]. Therefore, we propose that it is unnecessary to adopt additional components of comparable size to the base model, as the base model itself can effectively adjust to conditional controls. Instead, we replace the parallel control branch with a lightweight convolutional network to extract the conditional control features. As for the training, we fine-tune the base model by freezing most of its modules and selectively training a much smaller subset of the pretrained parameters. This approach avoids overfitting and catastrophic forgetting [6, 11, 15, 26] potentially caused by training. It significantly reduces the number of trainable parameters while incurring almost no additional inference latency overhead.

Furthermore, we introduce *Cross Normalization* as a replacement for *Zero Convolution* in fine-tuning large pre-trained models. *Zero Convolution*, which initializes weights to zeros allowing them to progressively grow during training, functions as a “bridge layer” that connects the control branch to the main branch. This operation is widely used in fine-tuning large pre-trained models and other parameter efficient fine-tune (PEFT) methods, as directly introducing new components and parameters to the pre-trained model can lead to training collapse [20, 43, 46]. *Zero Convolution* plays a role in gradually introducing the impact of new parameters. However, it also leads to slow convergence and training challenges because the learnable parameters initially struggle to receive the correct and instance-specific gradients. In this paper, we propose that the key reason for training collapse is the new initialized parameters sharing a different data distribution in terms of mean and standard deviation compared to the pre-trained parameters. Such distribution dissimilarity renders the two groups of parameters incompatible. This paper introduces *Cross Normalization* to align data distributions, leading to a more efficient and stable training process.

We conduct a series experiments on various generation backbones across the image and video generation [4, 31, 34, 38]. It demonstrates that ControlNeXt is robust and compatible with various types of conditional controls and network architectures. ControlNeXt retains the original base model’s architecture to a great extent, introducing only a few auxiliary components. This lightweight design allows it to function as a plug-and-play plugin compatible

with other methods. As shown in Figs. 1, 8, ControlNeXt can be integrated with other LoRA weights [15, 35] to modify styles without additional training. In summary:

- We present ControlNeXt, a powerful yet efficient method for image and video generation that significantly reduces latency overhead and parameters.
- We introduce *Cross Normalization* for fine-tuning pre-trained large models, which facilitates efficient and stable training convergence.
- ControlNeXt can serve as a lightweight plug-and-play module. It can be integrated with other LoRA weights to alter generation styles without additional training.

2. Related Work

Image and video diffusion models. Diffusion probability models [8, 14, 37] are advanced generative models that restore original data from pure Gaussian noise by learning the distribution of noisy data at various levels of noise. With their powerful capability to fit complex data distributions, diffusion models have excelled in several domains, including image and video generation.

In the field of image generation, diffusion models have surpassed previous Generative Adversarial Networks (GANs) in terms of image fidelity and diversity [8]. Moreover, the quality and consistency of videos generated by diffusion models continue to break new ground as the models are explored further. The most widely used neural network architectures in diffusion models are UNet [8, 14, 28, 29] and DiT [30], with others including U-ViT [3].

Specifically, diffusion models gradually remove noise from a Gaussian noise sample to obtain an image sample that follows the original data distribution. The forward and reverse diffusion processes are Markovian. In forward diffusion, the image at a certain time step depends only on the image from the previous time step. Diffusion models learn reverse diffusion from forward diffusion to gradually restore information from noise.

In recent years, latent diffusion models have incorporated variational autoencoders (VAEs) to transfer the diffusion process to latent space, significantly accelerating the model’s training and inference efficiency. Leading image generation models like Stable Diffusion [9, 31, 34, 38] have been widely adopted, used, and modified by the community, becoming a significant force in the advancement of AI-generated content (AIGC). This success is attributed to streamlined and efficient model architecture designs, including the network structures of diffusion models and various additional and extended components.

Controllable generation. Most recent models are guided by textual information, utilizing NLP models such as CLIP [32], Bert [7], or T5 [33] to extract textual features that guide the generated content. There are two main methods for introducing controllable conditions into image

or video generation models: (i) training a large diffusion model from scratch to achieve controllability under multiple conditions [17], (ii) fine-tuning a lightweight Adapter on a pretrained large model while keeping the original model parameters frozen [44]. Compared to the former’s substantial training costs, the latter’s lightweight nature makes it affordable for the community and individual users, facilitating easier dissemination and use.

Recent studies have attempted to control the outcomes of generative models by integrating additional neural networks into the foundation of diffusion models. ControlNet guides the generation of images that align with control information by duplicating certain layers from pre-trained large models and connecting them with zero convolutions to the original layers [41, 44]. In image and video generation tasks, ControlNet learns the relationship with the associated additional control conditions during training, thus allowing for the generation of related images based on the extra control conditions input during inference.

Low-Rank adaptation (LoRA). In recent years, as the AIGC community has grown, there has been a significant demand for personalized content customization. LoRA [15] reduces the costs of fine-tuning towards specific branches and disseminating models in the community by compressing the shift of the original model’s parameters into low-rank matrices, notably reducing both the amount of model parameters and the size of model files.

Information in diffusion. From an information theory perspective, the reverse diffusion in diffusion models is a process of gradually reconstructing information from Gaussian noise. Samples at a specific time step exhibit Markovian properties in their diffusion trajectory. Some studies alter the generated outcomes by perturbing the image information at certain time steps. For example, disturbing the initial noise distribution results the model’s incorrect learning of data distributions during training, thus mitigating the issue of average grayness in generation [45]. Image-to-image and inpainting guide the model to generate images similar to reference images by replacing the original Gaussian noise with a noisy image that has been added noise for several time steps and denoising from the middle of the diffusion trajectory.

3. Method

In this section, we provide a detailed technical overview of ControlNeXt. We first introduce the necessary preliminaries regarding diffusion models for controllable generation in Sec. 3.1. In Sec. 3.2, we delve into the analysis of the architecture design and prune it in order to make a concise and straightforward structure. Next, we introduce *Cross Normalization* in Sec 3.3, which is designed for the fine-tuning of large pre-trained models with additional introduced parameters.

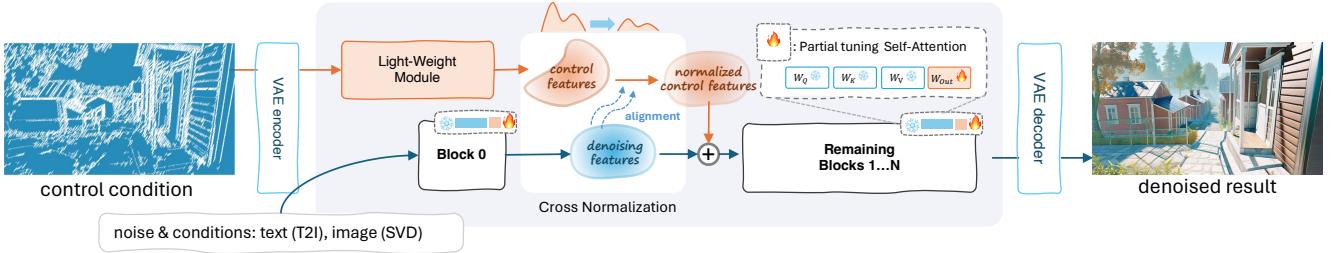


Figure 2. Training pipeline of ControlNeXt. We explore a more remarkable parameter-efficient framework than directly adopting a trainable copy.

3.1. Preliminaries

Diffusion model (DM) is a type of generative model that generates data by reversing a gradual noise-adding process, transforming random noise into coherent data samples. The model’s prediction for x_t at time step t depends only on x_{t+1} and t :

$$p_\theta(x_t|x_{t+1}) = \mathcal{N}(x_t; \tilde{\mu}_t, \tilde{\beta}_t I), \quad (1)$$

where θ represents the pre-trained model, $\tilde{\mu}_t$ is the model’s predicted target, and the variance $\tilde{\beta}_t$ is computed from the posterior of forward diffusion:

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t. \quad (2)$$

The loss function of diffusion models is the MSE loss function:

$$\mathcal{L} = w \cdot \mathbb{E}_{x_0, t, c_t, \epsilon \sim \mathcal{N}(0, 1)} [\|x - \hat{x}_\theta(x_t, t, c_t)\|^2], \quad (3)$$

where c_t represents text prompts, and w denotes the weight of the loss function. When $w = 1$, \mathcal{L} represents the x -prediction loss function; $w = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t}$, \mathcal{L} represents the noise-prediction loss function; $w = 1 + \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t}$, \mathcal{L} represents the v -prediction loss function. ControlNet [44] introduces controllable generation by integrating conditional controls. It calculates the loss function as:

$$\mathcal{L} = w \cdot \mathbb{E}_{x_0, t, c_t, c_f, \epsilon \sim \mathcal{N}(0, 1)} [\|x - \hat{x}_\theta(x_t, t, c_t, c_f)\|^2]. \quad (4)$$

3.2. Architecture Pruning

Motivation. The key innovation of ControlNet [44] is the addition of a control branch, which extracts conditional controls and integrates them into the main network. This branch shares trainable parameters initialized as a copy of the original half branch and operates in parallel, using a *zero convolution* as a bridge to integrate the conditional controls. Specifically:

$$y_c = \mathcal{F}_m(x) + \mathcal{Z}(\mathcal{F}_{cn}(x, c; \Theta_{cn}); \Theta_z), \quad (5)$$

where $\mathcal{F}(\cdot; \Theta)$ denotes a neural model with learnable parameters Θ , $\mathcal{Z}(\cdot; \Theta_z)$ indicates the *zero convolution* layer, and $x, y_c \in \mathbb{R}^{h \times w \times c}$ and c are the 2D feature maps and conditional controls, respectively. The pretrained large generation model \mathcal{F}_m , with pretrained parameters Θ_m , is completely frozen, while the ControlNet branch \mathcal{F}_{cn} is initialized by copying from the main branch.

However, such design, while introducing control capabilities, also incurs significant costs. The additional branch increases latency by at most 50%, which is particularly significant for video generation as each frame needs to be processed. Furthermore, the trainable parameters are substantial and fixed, equating to almost half of all pretrained parameters. Beyond the redundancy, optimizing only the ControlNet also limits the upper bound of the entire model, as it does not influence the pretrained generation model. To improve efficiency and clarity, we first simplify the original design of ControlNet by removing the additional branch. We then train a selected subset of the pretrained model, resulting in a more effective and efficient architecture.

Architecture pruning. It is important to note that the pretrained model is typically trained on a large-scale dataset, such as LAION-5B [36], whereas fine-tuning is always conducted on a much smaller dataset, often thousands of times smaller. Based on this, we maintain that the pretrained large generation model is sufficiently powerful, and there is no need to introduce such a large number of additional parameters to achieve control generation capabilities.

Specifically, we remove the control branch and replace it with a lightweight convolution module composed solely of multiple ResNet blocks [12]. This module, significantly smaller than the pretrained model, is designed to extract guidance information from the conditional controls and align it with the denoising features. Due to its small size, we rely more on the generation model itself to process the control signals. During training, we freeze most of the pretrained modules and selectively optimize a small subset of the trainable parameters from the pretrained generation model. This approach minimizes the risk of forgetting that could potentially be caused by the training process. It

can also be combined with parameter-efficient fine-tuning methods, such as LoRA [15, 24, 40]. However, we strive to maintain model structure consistency by avoiding significant modifications to the original architecture. Directly training the models also leads to greater effectiveness and efficiency, and can adaptively adjust the scale of the learnable parameters to fit various tasks. Mathematically,

$$\mathbf{y}_c = \mathcal{F}_m(\mathbf{x}, \mathcal{F}_c(\mathbf{c}; \Theta_c); \Theta'_m), \quad (6)$$

where $\Theta'_m \subseteq \Theta_m$ represents a trainable subset of the pre-trained parameters, and \mathcal{F}_c is the lightweight convolution module to extract the conditional controls. Based on the above processes, we strive to maintain the model’s consistency while minimizing additional expenses and latency as much as possible.

Regarding the injection of conditional controls, we observe that for most controllable generation tasks, the controls typically have a simple form or maintain a high level of consistency with the denoising features, eliminating the need to insert controls at multiple stages. We integrate the controls into the denoising branch at a single selected middle block by directly adding them to the denoising features after normalization through Cross Normalization. It can serve as a plug-and-play module constructed with a lightweight convolution module and learnable parameters, which are a subset of the pretrained models, represented as follows:

$$\mathcal{M}_c = \{\mathcal{F}_c(\cdot; \Theta_c), \Theta'_d\}, \quad (7)$$

where $\Theta'_d \subseteq \Theta_d$, and $\Theta_c \ll \Theta_d$.

3.3. Cross Normalization

Motivation. A typical problem in the continual training of pretrained large models is how to appropriately introduce additional parameters and modules. Since directly combining new parameters often leads to training collapse and poor convergence, recent works widely adopt zero initialization [44, 46], initializing the bridge layer that connects the pretrained model and the added module to zeros. Such an operation ensures that there is no influence from the newly introduced modules at the beginning of training, allowing for a stable warm-up phase. However, *zero initialization* can cause slow convergence and increases training challenges because it prevents the modules from receiving accurate gradients from the loss function. This results in a phenomenon known as “sudden convergence,” where the model doesn’t gradually learn the conditions but abruptly starts to follow them after an extended training period.

In this subsection, we analyze the reasons for training collapse when adding new parameters, and propose *cross normalization* to replace *zero convolution*, ensuring stable and efficient training.

Cross normalization. We find that the key reason for training collapse is the unaligned and incompatible data distribution between the introduced modules and pretrained models. After training on large-scale data, the pretrained generation model typically exhibits stable feature and data distributions, characterized by consistent mean and standard deviation. However, the newly introduced neural modules are typically only initialized using random methods [13, 21, 22], such as Gaussian initialization. This leads to the newly introduced neural modules producing feature outputs with significantly different means and standard deviations. Directly adding or combining these features results in model instability.

Normalization methods [2, 18, 42], like batch normalization and layer normalization, standardize layer inputs to improve training stability and speed. They achieve this by normalizing inputs to have zero mean and unit variance, which is widely used in neural network training. Inspired by their insights, we propose *cross normalization* to align the processed conditional controls and main branch features, ensuring training stability and speed.

We represent the feature maps processed from the main denoising branch and the control transferring branch as \mathbf{x}_d and \mathbf{x}_c , respectively, where $\mathbf{x}_m, \mathbf{x}_c \in \mathbb{R}^{h \times w \times c}$. The key of Cross Normalization is to use the mean and variance calculated from the main branch \mathbf{x}_m to normalize the control features \mathbf{x}_c , ensuring their alignment. First, calculate the mean and variance of the denoising features,

$$\boldsymbol{\mu}_m = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{m,i}, \quad (8)$$

$$\boldsymbol{\sigma}_m^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_{m,i} - \boldsymbol{\mu}_m)^2. \quad (9)$$

Then, we normalize the control features using the mean and variance of the denoising features,

$$\hat{\mathbf{x}}_c = \frac{\mathbf{x}_c - \boldsymbol{\mu}_m}{\sqrt{\boldsymbol{\sigma}_m^2 + \epsilon}} * \gamma, \quad (10)$$

where ϵ is a small constant added for numerical stability and γ is a parameter that allows the model to scale the normalized value.

Cross Normalization aligns the distributions of the denoising and control features, serving as a bridge to connect the diffusion and control branches. It accelerates the training process, ensures the effectiveness of the control on generation even at the beginning of training, and reduces sensitivity to the initialization of network weights.

4. Experiments

In this section, we conduct a series of experiments across various tasks and backbones. Our method demonstrates ex-

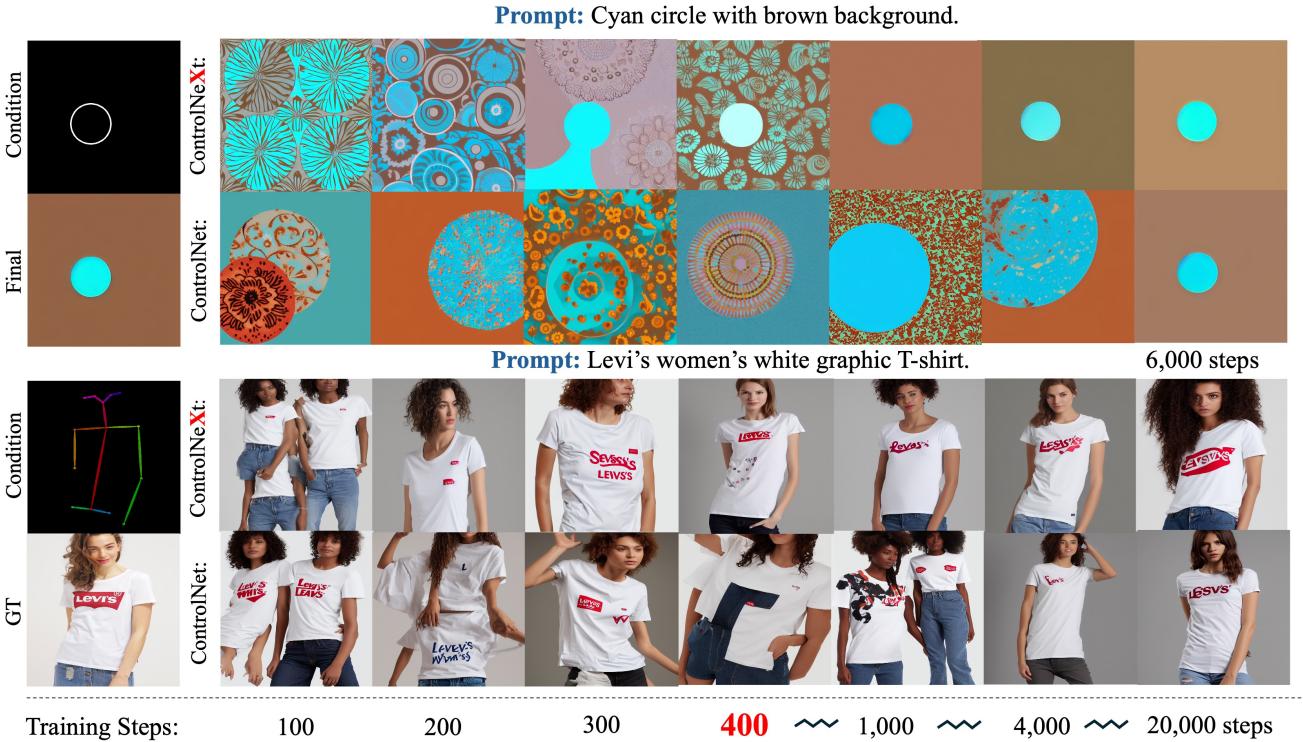


Figure 3. ControlNeXt achieves significantly faster training convergence and data fitting. It can learn to fit the conditional controls with fewer training steps, which also significantly alleviates the *sudden convergence* problem.

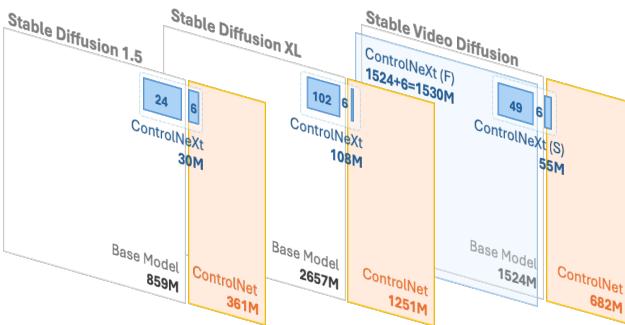


Figure 4. Parameter efficiency of ControlNeXt. We present the number of learnable parameters with various base models.

ceptional efficiency and generality in both image and video generation.

4.1. Generality

To demonstrate the robustness and generality of our methods, we first implement our method on various diffusion-based backbones, such as Stable Diffusion 1.5 [14, 38], Stable Diffusion XL [31], Stable Diffusion 3 [9] and Stable Video Diffusion [4]. It covers a wide range of tasks, including image generation, high-resolution generation, and

video generation, with various types of conditional controls, and the qualitative results are shown in Fig. 1. The results show that our method is robust and generalizable, effectively adapting to various architectures and meeting the requirements of different tasks.

Various conditional controls. ControlNeXt also supports various types of conditional controls. In this subsection, we choose “mask”, “depth”, “canny”, and “pose” as the conditional controls, shown in Fig. 5 from top to bottom, respectively. All the experiments are constructed based on the Stable Diffusion 1.5 architecture [38]. More generation results of the stable video generation, where we utilize the pose sequence as guidance for character animation, are presented in Fig. 6. The results of the SDXL are shown in Fig. 7, where we implement style transfer by extracting the Canny edges from the input images and generating the output with our SDXL model.

4.2. Training Convergence

A typical problem for the controllable geneartion is the hard training convergence, which means that it requires thousands or more than ten thousands steps training to learn the conditional controls. This phenomenon, known as the *sudden convergence* problem [44], occurs when the model initially fails to learn the control ability and then suddenly

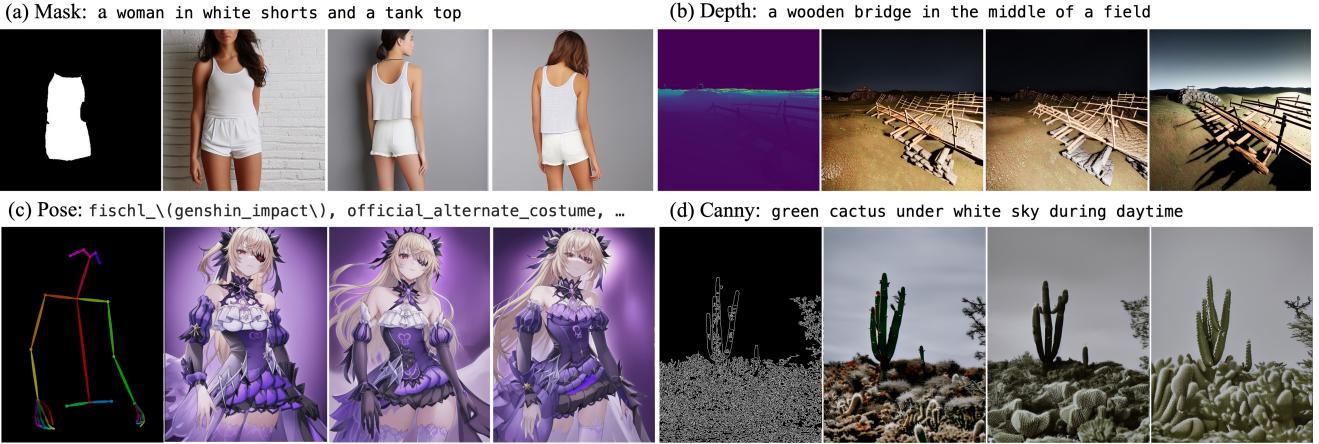


Figure 5. ControlNeXt also supports various conditional controls types. In this subsection, we select “mask”, “depth”, “canny”, and “pose” as the conditional controls, shown from top to bottom, respectively.



Figure 6. Detailed generation results of the stable video diffusion are provided. We utilize the pose sequence as guidance for character animation.

acquires this skill. This is caused from such two aspects:

1. *Zero convolution* inhibits the influence of the loss function, resulting in a prolonged warm-up phase where the model struggles to start learning effectively.
2. The pretrained generation model is completely frozen, and ControlNet functions as an adapter that cannot immediately affect the model.

In ControlNeXt, we eliminate these two limitations, resulting in significantly faster training convergence. We conducted experiments using two types of controls, and the results and comparisons are shown in Fig. 3. It can be seen that ControlNeXt starts to converge after only a few hundred training steps, while ControlNet requires thousands of

steps. This significantly alleviates the *sudden convergence* problem.

4.3. Efficiency

Our method only adds a lightweight module to the original pretrained model, ensuring it remains efficient and does not introduce significant latency. In this section, we provide more details and conduct additional experiments to demonstrate the efficiency of our method.

Parameters. First, we present statistics on the parameters, including the total and learnable parameters, calculated only for the UNet model (excluding the VAE and encoder

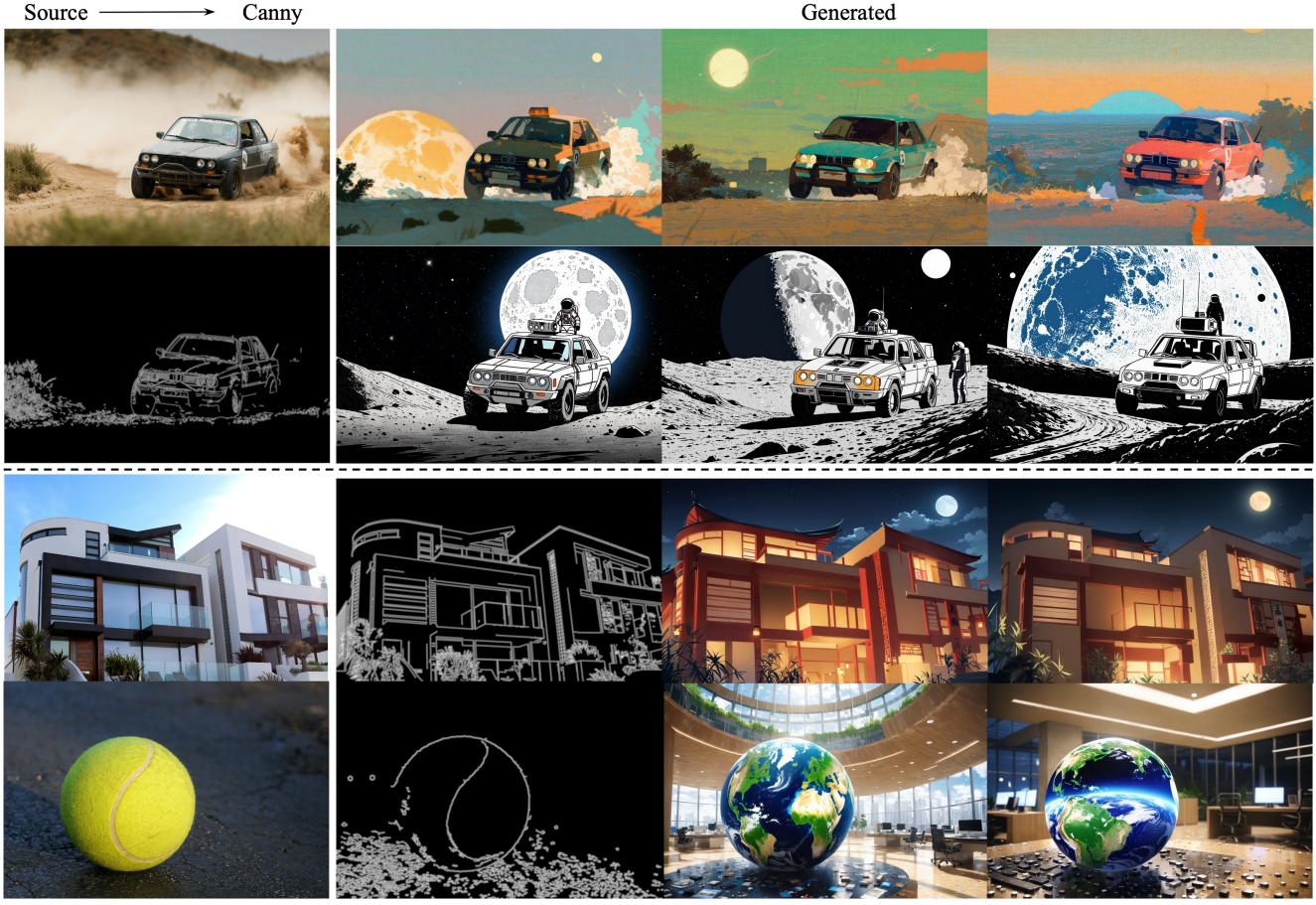


Figure 7. Detailed generation results of the stable diffusion xl are provided. We extract the Canny edges from the input natural image and implement the style transfer using our SDXL model.

Model	Method	Parameters (M)	
		Total	Learnable
SD1.5	ControlNet	1,220	361
	ControlNeXt(Our)	865	30
	Base model	859	-
SDXL	ControlNet	3,818	1,251
	ControlNeXt(Our)	2,573	108
	Base model	2,567	-
SVD	ControlNet	2,206	682
	ControlNeXt-S(Our)	1,530	55
	ControlNeXt-F(Our)	1,530	1,530
	Base model	1,524	-

Table 1. Comparison of the total and learnable parameters of different methods with various backbones. We select three base models: 1. stable diffusion 1.5 (SD1.5), 2. stable diffusion xl (SDXL), 3. stable video diffusion (SVD).

Method	Inference Time (s)			Δ
	SD1.5	SDXL	SVD	
ControlNet	0.31	1.01	1.73	+ 41.9%
ControlNeXt(Our)	0.24	0.82	1.29	+ 10.4%
Base model	0.22	0.70	1.23	-

Table 2. Comparison of the inference time with various backbones. Our method adds only minimal latency compared to the pretrained base model.

parts). And the results are shown in Tab. 1 and Tab. 4. It can be seen that our method only adds a lightweight module with minimal additional parameters, maintaining consistency with the original pretrained model. As for training, our method requires at most less than 10% of the learnable parameters, making it a very lightweight and plug-and-play module. You can also adaptively adjust the amount of learnable parameters for various tasks and



Figure 8. Our method can serve as a plug-and-play module that adapts to various generation models and LoRA weights, enabling changes in generation style without the need for training.



Figure 9. ControlNeXt serving as a plugin-unit to ensure a stable generaion with minimal costs.

performance requirements. More details on the influence of the number of parameters will be discussed later.

Inference time. We also compare the inference time of different methods with various base models. The results are shown in Tab. 2, which presents the computational time of one inference step, considering only the UNet and ControlNet parts and excluding the encoding and decoding processes. It can be seen that, since our method only adds a lightweight module, it increases latency minimally compared to the pretrained base generation model. This ensures outstanding efficiency advantages for our method.

4.4. Plug-and-Play

ControlNeXt is designed to preserve the consistency of the generation model’s original architecture, ensuring its compatibility and effectiveness. It can serve as a plug-and-play, training-free module that seamlessly integrates with various backbones and open-source LoR [15], enabling the alteration of generation styles.

Training free intergration. In this subsection, we first collected various LoRA weights downloaded from Civitai [1], encompassing diverse generation styles. We then construct experiments on various backbones based on the SD1.5 architecture, including SD1.5 [38], AnythingV3 [39] and DreamShaper [25]. The results are shown in Fig. 8. It can be observed that ControlNeXt can integrate with various backbones and LoRA weights in a training-free manner, effectively altering the quality and styles of generated images. This is mainly attributed to our method’s lightweight design, which primarily maintains the consistency of the pretrained base models and minimally adds additional modules. Such advantages enable it to serve as a plug-and-play module with general compatibility.

Stable generation. To generate satisfactory results using generation models, iterative prompt adjustments are often required. ControlNeXt, serving as a plug-in unit, facilitates stable generation with minimal effort and cost. We provide a simple prompt, "A woman," for the generation. The comparison of the generated results with and without our method is shown in Fig 9.

5. Conclusion

This paper presents ControlNeXt, an advanced and efficient method for controllable image and video generation. ControlNeXt employs a streamlined and concise architecture, eliminating heavy auxiliary components to minimize latency overhead and reduce trainable parameters. This lightweight design enables it to act as a plug-and-play module with strong robustness and compatibility, further allowing integration with other LoRA weights to alter generation styles without additional training. We propose *Cross Normalization* for finetuning pre-trained large models with newly introduced parameters, facilitating faster and more stable training convergence. Extensive experiments across various image and video generation backbones demonstrate the effectiveness and robustness of our methods.

Acknowledgements. This work was supported in part by the Research Grants Council under the Areas of Excellence scheme grant: AoE/E-601/22-R.



Figure 10. Stable Diffusion XL.



Figure 11. Stable video diffusion.

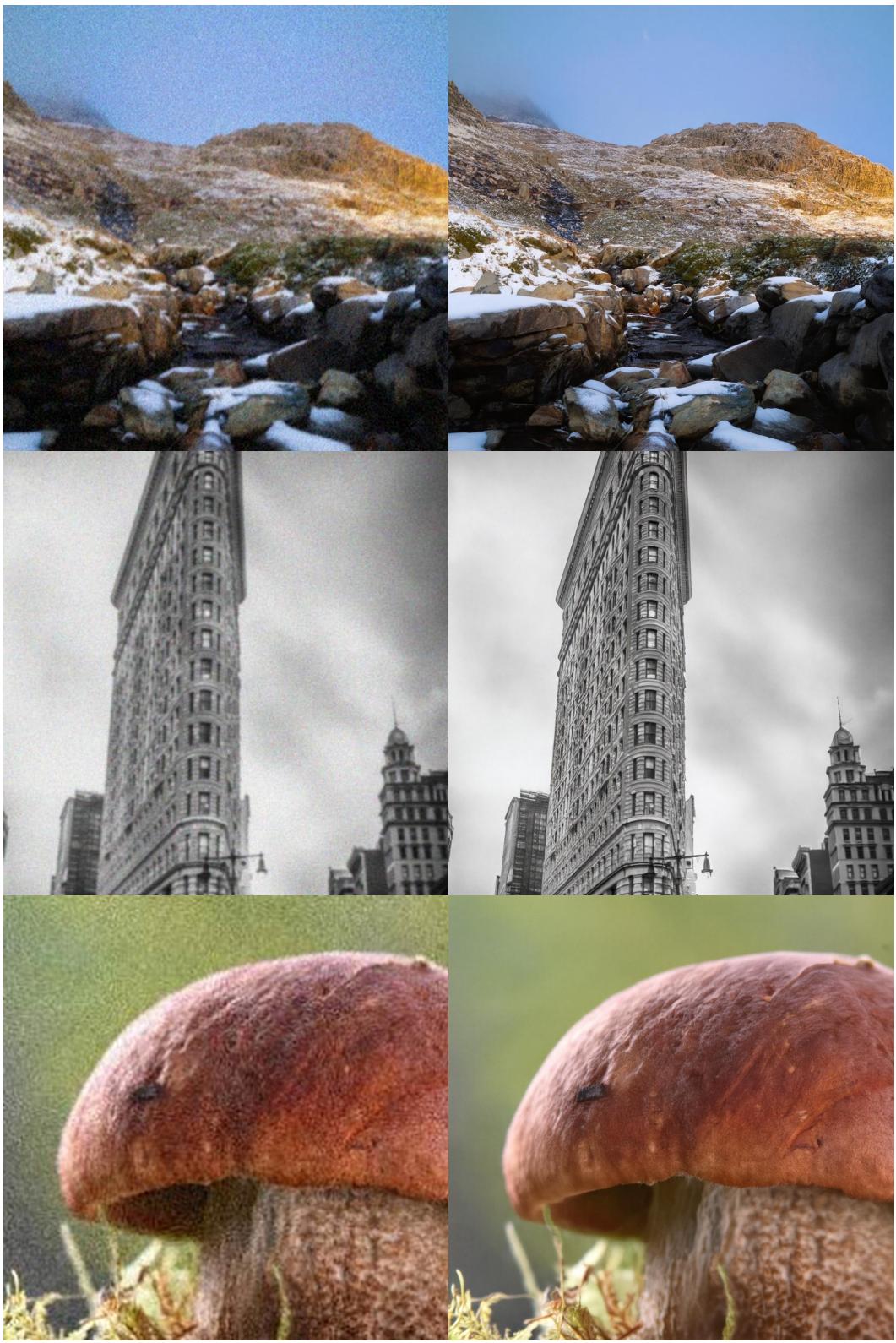


Figure 12. Stable Diffusion 3.

References

- [1] Civitai 2024. Civitai. <https://civitai.com>, 2024. 9
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 5
- [3] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. *arXiv preprint arXiv:2209.12152*, 2022. 3
- [4] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 2, 6
- [5] Pu Cao, Feng Zhou, Qing Song, and Lu Yang. Controllable generation with text-to-image diffusion models: A survey. *arXiv preprint arXiv:2403.04279*, 2024. 2
- [6] Zhe Chen, Yuchen Duan, Wenhui Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 2
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. 3
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2, 3
- [9] Patrick Esser et al. Scaling rectified flow transformers for high-resolution image synthesis. In *International Conference on Machine Learning*, 2024. 3, 6
- [10] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 2
- [11] Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024. 2
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4
- [13] Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4918–4927, 2019. 5
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NIPS*, pages 6840–6851, 2020. 3, 6
- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 3, 5, 9
- [16] Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8153–8163, 2024. 2
- [17] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023. 3
- [18] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015. 5
- [19] Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dreampose: Fashion image-to-video synthesis via stable diffusion. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22623–22633. IEEE, 2023. 2
- [20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2
- [21] Siddharth Krishna Kumar. On weight initialization in deep neural networks. *arXiv preprint arXiv:1704.08863*, 2017. 5
- [22] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015. 5
- [23] Zuzeng Lin, Ailin Huang, and Zhewei Huang. Collaborative neural rendering using anime character sheets. *arXiv preprint arXiv:2207.05378*, 2022. 2
- [24] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022. 5
- [25] Lykon. Dreamshaper. <https://huggingface.co/Lykon/DreamShaper>, 2022. 9
- [26] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022. 2
- [27] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4296–4304, 2024. 2
- [28] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, 2021. 3
- [29] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 3
- [30] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 3

- [31] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. [2](#), [3](#), [6](#)
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [3](#)
- [33] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. [3](#)
- [34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. [2](#), [3](#)
- [35] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. [3](#)
- [36] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. [2](#), [4](#)
- [37] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NIPS*, 2019. [3](#)
- [38] Stability. Stable diffusion v1.5. <https://huggingface.co/runwayml/stable-diffusion-v1-5>, 2022. [2](#), [3](#), [6](#), [9](#)
- [39] Furqanil Taqwa. Anything v3. <https://huggingface.co/Linaqruf/anything-v3.0>, 2022. [9](#)
- [40] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Guihong Cao, Dixin Jiang, Ming Zhou, et al. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*, 2020. [5](#)
- [41] Tan Wang, Linjie Li, Kevin Lin, Yuanhao Zhai, Chung-Ching Lin, Zhengyuan Yang, Hanwang Zhang, Zicheng Liu, and Lijuan Wang. Disco: Disentangled control for realistic human dance generation. *arXiv preprint arXiv:2307.00040*, 2023. [2](#), [3](#)
- [42] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. [5](#)
- [43] Yi Xin, Siqi Luo, Haodi Zhou, Junlong Du, Xiaohong Liu, Yue Fan, Qing Li, and Yuntao Du. Parameter-efficient fine-tuning for pre-trained vision models: A survey. *arXiv preprint arXiv:2402.02242*, 2024. [2](#)
- [44] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. [2](#), [3](#), [4](#), [5](#), [6](#)
- [45] Pengze Zhang et al. Tackling the singularities at the endpoints of time intervals in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. [3](#)
- [46] Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Ao-jun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. [2](#), [5](#)
- [47] Yuechen Zhang, Shengju Qian, Bohao Peng, Shu Liu, and Jiaya Jia. Prompt highlighter: Interactive control for multi-modal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13215–13224, 2024. [2](#)