

# AdaDiff: Adaptive Step Selection for Fast Diffusion

Hui Zhang<sup>1,2</sup> Zuxuan Wu<sup>1,2</sup> Zhen Xing<sup>1,2</sup> Jie Shao<sup>3</sup> Yu-Gang Jiang<sup>1,2</sup>  
<sup>1</sup>Shanghai Key Lab of Intell. Info. Processing, School of CS, Fudan University  
<sup>2</sup>Shanghai Collaborative Innovation Center of Intelligent Visual Computing  
<sup>3</sup>ByteDance Inc

## Abstract

Diffusion models, as a type of generative models, have achieved impressive results in generating images and videos conditioned on textual conditions. However, the generation process of diffusion models involves denoising for dozens of steps to produce photorealistic images/videos, which is computationally expensive. Unlike previous methods that design “one-size-fits-all” approaches for speed up, we argue denoising steps should be sample-specific conditioned on the richness of input texts. To this end, we introduce AdaDiff, a lightweight framework designed to learn instance-specific step usage policies, which are then used by the diffusion model for generation. AdaDiff is optimized using a policy gradient method to maximize a carefully designed reward function, balancing inference time and generation quality. We conduct experiments on three image generation and two video generation benchmarks and demonstrate that our approach achieves similar results in terms of visual quality compared to the baseline using a fixed 50 denoising steps while reducing inference time by at least 33%, going as high as 40%. Furthermore, our qualitative analysis shows that our method allocates more steps to more informative text conditions and fewer steps to simpler text conditions.

## 1. Introduction

Diffusion models [5, 11, 31, 38, 43], as a class of generative models, have made significant strides. These models have the capability to generate specified visual contents, including images and videos, based on specific input conditions such as text, semantic maps, representations, and images. For example, models like Stable Diffusion [31] and ModelScopeT2V [43] can produce perceptually-convincing or artistic images and videos conditioned on textual descriptions, *a.k.a.* prompts. These diffusion models often contain an iterative denoising process during generation, and more iterations typically indicate better visual quality. However, the improvements come at the cost of increased computational resources, even with the use of state-of-the-art sampling



Figure 1. **A conceptual overview of our approach.** AdaDiff allocates different numbers of generation steps for prompts with varying levels of richness, aiming to minimize inference time while maintaining high image quality. Images with red borders are produced by AdaDiff.

methods [38]. Therefore, to strike a balance between quality and inference speed, the number of denoising steps is often empirically set as a fixed value (typically 50), regardless of the prompts used for generation.

*But do we really need a fixed number of denoising steps for all different prompts?* Intuitively, using more steps may lead to higher quality and more detailed contents [21, 38]. Nonetheless, in real-world applications, the richness in textual prompts, *i.e.*, the number of objects, and how they relate on each other, vary significantly. For certain easy and coarse-grained prompts, involving only one or a few objects, using

fewer steps is already sufficient to generate satisfactory results, and increasing the number of steps may lead to only marginal improvements and does not necessarily produce better results. For complex textual prompts, containing many objects, detailed descriptions, and intricate interactions between objects, a larger number of steps becomes necessary to achieve the desired result. Therefore, the goal of this paper is to develop a dynamic framework for diffusion models by adaptively determining the number of denoising steps needed for generating photorealistic contents conditioned on textual inputs. This is in contrast yet complimentary to existing approaches that speed up diffusion models by designing new schedulers [20, 21, 38] or fast sampling mechanisms [23, 34], all of which adopt the same number of steps regardless of the complexity in textual prompts.

In light of this, we introduce AdaDiff, an end-to-end framework that aims to achieve efficient diffusion models by learning adaptive step selection in the denoising process based on prompts. For each prompt, deriving a dynamic generation strategy involves: I) determining the required number of steps for generation and II) ensuring high-quality generation even with a relatively smaller number of steps. With this, AdaDiff is capable of allocating more computational resources to more descriptive prompts while using fewer resources for simpler ones. While this approach is highly appealing, learning the dynamic step selection is a non-trivial task, as it involves non-differentiable decision-making processes.

To address this challenge, AdaDiff is built upon a reinforcement learning framework [41]. Specifically, given a prompt (text condition), AdaDiff trains a lightweight step selection network to produce a policy for step usage. Subsequently, based on this derived policy, a dynamic sampling process is performed on a pre-trained diffusion model for efficient generation. The step selection network is optimized using a policy gradient method to maximize a meticulously crafted reward function. The primary objective of this reward function is to encourage the generation of high-quality visual contents while minimizing computational resources. It is also worth pointing out that the step selection network conditioned on textual inputs is lightweight with negligible computational overhead.

We conduct extensive experiments to evaluate our proposed method, and the results demonstrate that AdaDiff saves between 33% and 40% of inference time compared to the baseline method while maintaining similar visual quality across various image and video generation benchmarks [19, 36, 44, 45, 51]. In addition, we demonstrate that our approach can be combined with various existing acceleration paradigms. Moreover, the learned policy from one dataset can be successfully transferred to another. Finally, through qualitative and quantitative analysis, we show that AdaDiff flexibly allocates fewer sampling steps for less

informative prompts and more for informative prompts.

## 2. Related Work

**Diffusion Models.** Diffusion models [5, 11, 40] have emerged as a powerful force in the field of deep generative models, achieving top-notch performance in various applications, spanning image generation [18, 25–27, 30–33, 54, 56], video generation [3, 12, 13, 37, 43, 46, 49, 50, 55, 58], image segmentation [1, 52], object detection [4], and image restoration [9, 39, 48], among others. Notably, in image and video generation, these diffusion models [31, 43] have demonstrated the ability to produce desired results based on diverse input conditions, including text, semantic maps, representations, and images. However, the inherent iterative nature of the diffusion process has led to a substantial demand for computational resources and inference time during the generation process.

**Reinforcement Learning in Diffusion Models.** Pre-training objectives of generative models often do not align perfectly with human intent. Therefore, some work focuses on fine-tuning generative models through reinforcement learning [41] to align their outputs with human preferences, using human feedback or carefully designed reward functions. Typically, these models [2, 6, 8, 16, 28, 42, 47, 53] enhance aspects such as text-to-image alignment, aesthetic quality, and human-perceived image quality. Nevertheless, we are the first to leverage reinforcement learning to accelerate image and video generation by learning an instance-specific step usage policy.

**Acceleration of Diffusion Models.** Recently, effort has been made to accelerate the reverse process of diffusion models. These approaches can be broadly categorized into two paradigms. The first paradigm [14, 20, 21, 38, 57] involves designing new samplers that extract a subsequence from the original sampling sequence as a new schedule, typically with fewer sampling steps, to expedite the generation process. The second paradigm is based on knowledge distillation, which uses a student model to learn how to approximate the output of a teacher model with fewer steps [22, 23, 34], or learning more efficient sampling paths to replace the slower sampling process of the teacher model [7, 17].

The proposed AdaDiff method is orthogonal to, yet complementary to, the aforementioned acceleration methods. While previous methods reduce the total number of sampling steps, they still rely on a manually set and “one-size-fits-all” denoising step, ignoring the complexity in textual prompts that are used to condition the generation process. Instead, AdaDiff aims to dynamically decide the optimal step that achieves a balance between visual quality and inference time on a per-input basis. It is also worth noting that AdaDiff can be used together with existing speed-up methods, as will be shown in experiments.

### 3. Methodology

AdaDiff reduces computational cost and inference time for diffusion models by learning step usage policies conditioned on different textual prompts. The intuition is to encourage using fewer denoising steps while generating high-quality results. In the following, we will first review the background knowledge of diffusion models (Sec. 3.1). Subsequently, we will delve into the components of AdaDiff and how it facilitates adaptive generation based on diffusion models, including image and video generation (Sec. 3.2 and Sec. 3.3).

#### 3.1. Background on diffusion models

Diffusion models [11, 31, 43] achieve new state-of-the-art performance in the field of deep generative models inspired by the principles of equilibrium thermodynamics. Specifically, diffusion models involve a forward process where noise is gradually added to the input and a reverse process that learns to recover the desired noise-free data from noisy data. In the forward process, the posterior probability of the diffusion image  $x_t$  at time step  $t$  has a closed form:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad (1)$$

where  $\bar{\alpha}_t = \prod_{i=0}^t \alpha_i = \prod_{i=0}^t (1 - \beta_i)$  and  $\beta_i \in (0, 1)$  represents the noise variance schedule.

Once the diffusion model  $\epsilon_\theta(x_t)$  is trained, during the reverse process, traditional diffusion models like DDPM [11] denoise  $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  step by step for a total of  $T$  steps. One can also use a discrete-time DDIM [38] sampler to speed up the sampling process. Initially, the total number of sampling steps  $S$  is predetermined. Sampling updates are performed at every  $\lceil T/S \rceil$  steps according to the plan, reducing the original  $T$  steps to a new sampling plan that consists of a subset of  $S$  diffusion steps  $\hat{T} = \{\eta_1, \dots, \eta_S\}$ . As a result, the inference process becomes:

$$x_{t-1} = \sqrt{\alpha_{t-1}} \left( \frac{x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t, t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(x_t, t) + \sigma_t \epsilon_t \quad (2)$$

where  $t-1, t \in \hat{T}$  and  $\sigma_t^2 = \eta \cdot \tilde{\beta}_t = \eta \cdot \beta_t \cdot \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}$ . DDIM sampler significantly reduces the total number of sampling steps and is widely employed in various generative tasks using diffusion models.

**Latent Diffusion Models (LDMs).** LDMs [31, 43] employ an approach where the diffusion process operates in the latent space rather than the pixel space. This reduces the training cost and improves the inference speed. It uses the pre-trained encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$  to encode the pixel-space image into a low-dimensional latent and decode the latent back into the image, respectively. In addition, LDMs incorporate flexible conditional information, such as text

conditions and semantic maps, through a cross-attention mechanism to guide the visual generation process.

To date, Stable Diffusion [31] and ModelScopeT2V [43] have been influential approaches for image and video generation conditioned on texts. They utilize a DDIM sampler during generation, with a default of 50 steps for all prompts. In this paper, AdaDiff aims to enhance the inference speed of these LDMs by implementing dynamic step selection on a per-input basis.

#### 3.2. Adaptive step selection for image generation.

In image generation, AdaDiff learns a step usage policy conditioned on the text description to reduce the denoising step of Stable Diffusion such that steps vary for different prompts. To this end, AdaDiff builds upon a lightweight selection network trained to determine the total number of steps  $\mathbf{t}$  of the DDIM sampler used in the reverse process of Stable Diffusion, as shown in Fig. 2. The selection network makes non-differentiable categorical decisions, *i.e.*, steps to be used. We design  $N$  distinct schedulers for the DDIM sampler, each corresponding to a specific total number of sampling steps. In this paper, unless specified otherwise, we set  $N = 5$ , which corresponds to the set of step values  $\mathcal{S} = \{10, 20, 30, 40, 50\}$ . The state space is defined as the input prompts, and actions in the model involve categorizing them in these discrete action spaces. Then, a carefully designed reward function balances the quality of the generated images with the computational cost.

Formally, given a prompt  $\mathbf{p}$ , Stable Diffusion first uses a text encoder  $\tau$  to extract the text features, denoted as  $\mathbf{c} = \tau(\mathbf{p})$ . Following this, the step selection network  $f_s$ , parameterized by  $\mathbf{w}$ , learns the informativeness of  $\mathbf{c}$  using self-attentional mechanisms and then further maps it to  $\mathbf{s} \in \mathbb{R}^N$  through a Multi-Layer Perceptron (MLP):

$$\mathbf{s} = f_s(\mathbf{c}; \mathbf{w}), \quad (3)$$

where each entity in  $\mathbf{s}$  indicates the probability score of choosing this step. We then define a step selection policy  $\pi^f(\mathbf{u} | \mathbf{p})$  with a  $N$ -dimensional Categorical Distribution. Here,  $\mathbf{u}$  is a one-hot vector of length  $N$ , denoted as  $\mathbf{u} \in \{0, 1\}^N$ , and  $\mathbf{u}_j = 1$  indicates that the step  $\mathbf{t}$  with index  $j$  in  $\mathcal{S}$  is selected. In the training phase,  $\mathbf{u}$  is generated by sampling from the corresponding policy, and in the testing phase, a greedy approach is employed.

So far, the total number of steps  $\mathbf{t}$  for DDIM sampling is determined on a per-prompt  $\mathbf{p}$  basis. Then, Stable Diffusion starts with a latent  $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , fuses the text-condition features  $\mathbf{c}$  through cross-attention, and further generates the desired latent  $\mathbf{z}_0$  through  $\mathbf{t}$  denoising steps. Finally, the decoder  $\mathcal{D}$  decodes the latent into a pixel-space generated image  $\mathbf{x}$ . This process can be formalized as:

$$\mathbf{x} = \mathcal{D}(\text{LDM}_{\text{sampler}}^{\mathbf{p} \rightarrow \mathbf{t}}(\mathbf{z}_T, \mathbf{c})). \quad (4)$$

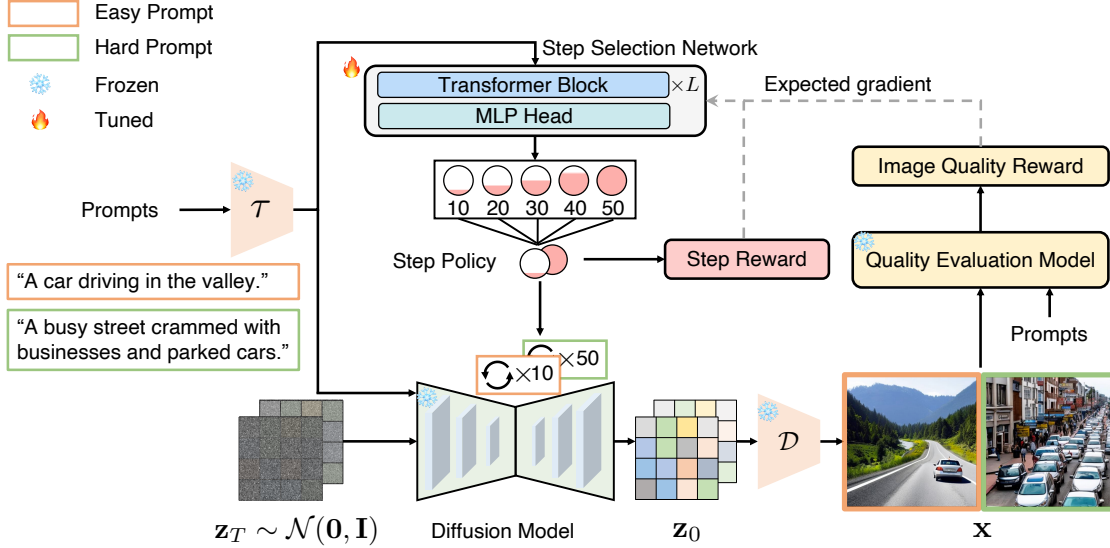


Figure 2. **An overview of the proposed pipeline AdaDiff.** Given the input prompts, the step selection network learns the information richness of each prompt and derives the corresponding step usage policy. These policies determine the number of steps required for the diffusion model to generate images. Subsequently, a carefully designed reward function balances the trade-off between inference time and image quality. See texts for more details.

Recall that the primary objective of AdaDiff is to generate high-quality images in a smaller number of sampling steps, hence, it is crucial to design a suitable reward function to evaluate these actions efficiently. The reward function consists primarily of two parts: an image quality reward and a step reward, balancing quality and inference time. Firstly, to assess image quality, we leverage a quality evaluation model  $f_q$  specifically designed to evaluate the quality of generated images [53], abbreviated as the IQS model. This model assesses image quality in two dimensions: image-text alignment and perceptual fidelity. Image-text alignment means that the generated image should match the user-provided text conditions, while perceptual fidelity means that the generated image should be faithful to the shape and characteristics of the object rather than being generated chaotically. Typically, the higher the IQS score  $f_q(\mathbf{x})$ , the higher the quality of the generated image  $\mathbf{x}$ . Thus, in this paper, we design the image quality reward as  $Q(\mathbf{u}) = f_q(\mathbf{x})$ .

For the step reward, we define it as  $\mathcal{O}(\mathbf{u}) = 1 - \frac{\mathbf{t}}{S_{max}}$ , which represents the normalized steps saved relative to the maximum steps in  $\mathcal{S}$ . Finally, the overall reward function is formalized as follows:

$$R(\mathbf{u}) = \begin{cases} \mathcal{O}(\mathbf{u}) + \lambda Q(\mathbf{u}) & \text{for high quality image} \\ -\gamma & \text{else} \end{cases} \quad (5)$$

where  $\lambda$  is a hyperparameter that controls the effect of image quality reward  $Q(\mathbf{u})$  and  $\gamma$  is the penalty imposed on the reward function when the generated image quality is low. Instead of a straightforward comparison between the image quality score  $f_q(\mathbf{x})$  and a predefined threshold (e.g., 0) to

discern whether  $\mathbf{x}$  is a high-quality image, we design the determination to be a relative manner. Specifically, for a given prompt, we individually generate an image for each step in the step set  $\mathcal{S} = \{10, 20, 30, 40, 50\}$  and we consider the image quality score to be high if it ranks top  $k$  among these five images (we empirically set  $k$  to 3). At this point, the step selection network can be optimized to maximize the expected reward:

$$\max_{\mathbf{w}} \mathcal{L} = \mathbb{E}_{\mathbf{u} \sim \pi_f} R(\mathbf{u}). \quad (6)$$

In this paper, we use the policy gradient method [41] to learn the parameters  $\mathbf{w}$  for the step selection network. The expected gradient can be derived as follows:

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbb{E} [R(\mathbf{u}) \nabla_{\mathbf{w}} \log \pi^f(\mathbf{u} | \mathbf{p})], \quad (7)$$

which is further approximated with Monte-Carlo sampling using mini-batches of samples Eq. (7) as:

$$\nabla_{\mathbf{w}} \mathcal{L} \approx \frac{1}{B} \sum_{i=1}^B [R(\mathbf{u}_i) \nabla_{\mathbf{w}} \log \pi^f(\mathbf{u}_i | \mathbf{p}_i)]. \quad (8)$$

where  $B$  is the total number of prompts in the mini-batch. The gradient is then propagated back to train the step selection network using the Adam optimizer.

Following the aforementioned training process, the selection network learns the step usage policy that strikes a balance between inference time and generation quality. During the inference phase, for different prompts, the maximum probability score in  $\mathbf{s}$  is used to determine the number of generation steps, enabling dynamic inference.



### 3.3. Adaptive step selection for video generation.

In addition to image generation, the proposed step selection strategy can also be applied to video diffusion models, such as ModelScopeT2V [43], as the prompts used to guide video generation also have varying levels of richness. The overall implementation paradigm is similar to Fig. 2. To assess the quality of the generated videos  $\mathcal{V}$ , we individually score each frame using the IQS model and subsequently calculate the average of all frames to obtain the video quality reward. This procedure can be formalized as follows:

$$\mathcal{Q}(u) = \frac{1}{F} \sum_{i=1}^F f_q(\mathcal{V}_i). \quad (9)$$

Here,  $F$  is the number of frames in the generated video. Then, we calculate the step reward in the same way as in Sec. 3.2, and the overall reward is obtained from Eq. (5).

## 4. Experiments

### 4.1. Experimental Details

**Datasets.** To evaluate the effectiveness and generalizability of our approach, we conduct extensive experiments on three image datasets: MS COCO 2017 [19], Laion-COCO [36], DiffusionDB [45], and two video datasets: MSR-VTT [51] and InternVid [44]. In MS COCO 2017, our training set consists of 118, 287 textual descriptions, and all 25, 014 text-image pairs from the validation set are employed for testing. Regarding Laion-COCO, we randomly select 200K textual descriptions for training and 20K text-image pairs for testing. The partitioning of the training and testing sets for DiffusionDB follows the same paradigm as Laion-COCO. The training sets for MSR-VTT and InternVid consist of 6, 651 and 24, 911 text descriptions, respectively. The test set for MSR-VTT comprises 2, 870 text-video pairs.

**Evaluation Metrics.** Following [23, 43, 49], we assess the quality of generated images or videos using several metrics, including the Fréchet Inception Distance Score (FID) [10], Inception Score (IS) [35], text-to-image similarity (CLIP Score) [29], Natural Image Quality Evaluator score (NIQE) [24], and the recently introduced Image Quality Score (IQS) [53]. Additionally, we use the average number of denoising steps and time per image or video generation to measure the inference speed.

**Implementation Details.** We design the step selection network as a lightweight architecture consisting of three self-attention layers and a multi-layer perceptron. For image generation, we employ the Stable Diffusion v2.1-base model to generate  $512 \times 512$  images. For video generation, we use ModelScopeT2V to generate 16-frame videos with a resolution of  $256 \times 256$ . The computational cost of determining the total number of sampler steps per prompt is only 1.93

	Speed		Image Quality				
	Step↓	Time↓	IQS↑	CLIP↑	IS↑	FID↓	NIQE↓
<b>COCO 2017</b>							
	10	0.58	0.215	0.311	35.18	22.08	4.17
	20	0.97	0.357	0.312	37.15	22.12	3.85
SD	30	1.44	0.376	0.313	37.56	22.48	3.87
	40	1.81	0.401	0.314	37.21	22.05	3.76
	50	2.24	<b>0.419</b>	0.314	37.48	22.13	<b>3.75</b>
Random	30.03	1.41	0.354	0.313	36.85	22.50	3.88
AdaDiff	28.61	1.35	0.412	<b>0.314</b>	<b>37.60</b>	<b>21.92</b>	3.76
<b>Laion-COCO</b>							
	10	0.54	0.116	0.314	28.45	22.14	4.67
	20	0.94	0.304	0.318	30.29	22.64	4.58
SD	30	1.42	0.316	0.320	30.26	21.97	4.67
	40	1.83	0.336	0.319	30.39	22.08	4.57
	50	2.27	<b>0.350</b>	0.319	<b>30.71</b>	22.08	4.58
Random	30.07	1.46	0.286	0.318	29.51	22.45	4.62
AdaDiff	31.34	1.50	0.345	<b>0.324</b>	30.51	<b>22.07</b>	<b>4.58</b>
<b>DiffusionDB</b>							
	10	0.57	0.065	0.325	15.32	9.29	4.76
	20	0.95	0.225	0.327	15.43	8.57	4.42
SD	30	1.43	0.249	0.327	15.48	8.67	4.46
	40	1.84	0.274	0.326	15.14	8.62	4.30
	50	2.28	<b>0.281</b>	0.327	15.41	8.57	4.30
Random	30.01	1.42	0.221	0.326	15.34	8.74	4.48
AdaDiff	32.38	1.52	0.273	<b>0.328</b>	<b>15.52</b>	<b>8.56</b>	4.33

Table 1. Comparison of AdaDiff with other baselines on three image generation benchmarks. SD represents the Stable Diffusion v2-1-base Model.

GFLOPs, which is negligible compared to the substantial computational requirements of diffusion models. For example, when Stable Diffusion generates an image in 50 steps, its computational cost is approximately 35,140 GFLOPs. In the reward function, we set the parameter  $\lambda$  to 2, and we define images with IQS scores within the top 3 as high-quality images. We train the step selection network for 200 epochs with a batch size of 256. Optimization is performed using the Adam optimizer with an initial learning rate of  $10^{-5}$ .

### 4.2. Main Results

**Performance on Image Generation.** To validate the effectiveness and general applicability of AdaDiff, we compare it with different baselines across three image generation benchmarks. One of these baselines is Stable Diffusion (SD), which employs a fixed number of sampling steps for various

MSR-VTT	Speed		Image Quality				
	Step↓	Time↓	IQS↑	CLIP↑	IS↑	FID↓	NIQE↓
ModelScope	50	21.2	-0.518	0.293	<b>18.79</b>	44.85	6.57
Random	29.98	13.5	-0.723	0.293	18.22	47.41	6.75
AdaDiff	31.14	13.6	<b>-0.517</b>	<b>0.299</b>	18.74	<b>44.49</b>	<b>6.36</b>

Table 2. Comparison of AdaDiff on video generation.

prompts, *i.e.*  $S = \{10, 20, 30, 40, 50\}$ . Another baseline involves using a random step selection strategy for different prompts, and we report mean steps used of 5 runs.

Tab. 1 provides a comprehensive breakdown of the performance of AdaDiff across different benchmarks. On COCO 2017, Laion-COCO, and DiffusionDB, AdaDiff assigns average sampling step counts of 28.61, 31.34, and 32.38 per prompt, respectively. When compared to Stable Diffusion, which uses a fixed 50-step sampling process, AdaDiff demonstrates comparable or superior performance across five image quality metrics while saving 40% (1.35 vs. 2.24), 34% (1.50 vs. 2.27), and 33% (1.52 vs. 2.28) of the time required to generate each image on the three respective datasets. This confirms that AdaDiff learns efficient step usage policies, thus conserving inference time while maintaining competitive performance across various datasets.

Compared to baselines with similar computational resources, such as a fixed step count of 30 and a random policy, AdaDiff outperforms them by a significant margin regarding image quality on all three datasets. Notably, AdaDiff improves IQS Score by 16.4% (0.412 vs. 0.354), 20.6% (0.345 vs. 0.286), and 23.5% (0.273 vs. 0.221) on the three datasets, respectively, compared to the random strategy. When compared to a fixed step count of 30, AdaDiff improves IQS Score by 9.6%, 9.2%, and 9.6% on the same datasets. These results confirm that AdaDiff generates adaptive policies and allocates different sampling step counts for each prompt to maintain image quality. In addition, AdaDiff consistently demonstrates superior performance over the random policy across the CLIP, IS, FID, and NIQE metrics. However, the improvement in some metrics is not significant, as they typically diverge from human judgments of the quality of the generated images [15, 47, 53]. For example, metrics such as IS, FID, and NIQE focus on the perceptual quality of images but overlook alignment with prompts. On the other hand, CLIP Score emphasizes text-image alignment but does not assign equal importance to perceptual quality. Consequently, the fluctuations in these metrics are not highly pronounced when a diffusion model generates results for a given prompt across different step numbers. In contrast, the IQS Score demonstrates greater sensitivity and a more consistent alignment with human judgment regarding instance quality [53].

**Extension to other diffusion models.** To verify the applicability of AdaDiff to other diffusion models, we evaluate

COCO 2017	Speed		Image Quality				
	Step↓	Time↓	IQS↑	CLIP↑	IS↑	FID↓	NIQE↓
SD	50	2.42	<b>0.436</b>	0.312	<b>38.76</b>	21.72	3.78
Random	30.04	1.48	0.416	0.311	38.37	21.69	3.86
AdaDiff	24.09	1.19	0.434	<b>0.313</b>	38.57	<b>21.35</b>	<b>3.76</b>

Table 3. Extension to another acceleration method, DPM-Solver.

its performance in video generation, as depicted in Tab. 2. Compared to the results generated by ModelScope using 50 steps, AdaDiff exhibits superior video quality and reduces the generation time per video by 35.8% (13.6 vs. 21.2). When compared to a random step usage strategy, AdaDiff utilizes similar computational resources but significantly improves video quality across all aspects, particularly with a 28.5% improvement in IQS Score. These results demonstrate the versatility of our proposed method, as it can be deployed on a variety of text-conditioned diffusion models, providing more efficiency.

**Extension to other acceleration methods.** To validate the compatibility of AdaDiff with different acceleration paradigms in the diffusion models, we evaluate its performance in image generation when using the DPM-Solver [21] sampler in Stable Diffusion, as illustrated in Tab. 3. Compared to 50-step generation, AdaDiff allocates an average of 24.09 steps per prompt, further saving 50.8% (1.19 vs. 2.42) of generation time while maintaining comparable image quality. Compared to a random policy, the learned adaptive policy not only speeds up inference time by 19.6% but also improves image quality across five metrics. These results demonstrate that AdaDiff can be used as a plug-and-play component in combination with other acceleration methods for dynamic generation.

**Extension to other datasets.** We also evaluate AdaDiff’s ability to generalize its learned step selection policy from one dataset to another, which we refer to as zero-shot gen-

	Speed		Image / Video Quality				
	Step↓	Time↓	IQS↑	CLIP↑	IS↑	FID↓	NIQE↓
<b>COCO 2017 → Laion-COCO</b>							
SD	50	2.27	<b>0.350</b>	0.319	<b>30.71</b>	22.08	4.58
Random	30.07	1.46	0.286	0.318	29.51	22.45	4.62
AdaDiff	30.50	1.48	0.341	<b>0.320</b>	30.68	<b>21.94</b>	<b>4.57</b>
<b>InternVid → MSR-VTT</b>							
ModelScope	50	21.20	<b>-0.518</b>	0.293	18.79	44.85	6.57
Random	29.98	13.51	-0.723	0.293	18.22	47.41	6.75
AdaDiff	32.23	14.03	-0.521	<b>0.299</b>	<b>18.96</b>	<b>44.58</b>	<b>6.25</b>

Table 4. Validation of AdaDiff’s zero-shot adaptive generation.

eration performance. For image generation, we evaluate the step usage policy derived from COCO 2017 on Laion-COCO, while for video generation, we use the policy of InternVid for validation on MSR-VTT. As shown in Tab. 4, the dynamic strategy saves 33.6% and 33.8% of generation time on the two datasets separately compared to the fixed 50-step generation while maintaining comparable generation quality. Besides, AdaDiff exhibits superior generation quality when compared to random strategies. These results demonstrate the transferability of step selection strategies trained on large-scale data.

**Analyses of learned policies.** To better understand the learned policy of AdaDiff, we investigate the relationship between step selection and prompt richness through Fig. 3. We evaluate prompt richness in terms of the number of words and objects, assuming an overall increase in information richness as both grow. We observe that as the number of words in the prompt increases, the generated results may include more detailed descriptions, spatial relationships, attribute definitions, *etc.* Consequently, AdaDiff allocates more denoising steps for each prompt. Besides, as the number of objects in the prompt increases, the results involve a larger amount of details and interactions among the objects. As a result, AdaDiff also assigns more steps per prompt.

We further qualitatively analyze our approach as shown in Fig. 4. We specifically examine five distinct scenarios {indoor, food, animal, outdoor, and sports}, and explore the impact of the richness of the prompts on the step usage policy within the same scenario. Our observations reveal that in instances where the prompt is straightforward (easy, Fig. 4), typically involving only one or a few objects, AdaDiff assigns 10 or 20 steps to obtain satisfactory generated results. For prompts characterized by an increased number of objects or the inclusion of detailed descriptions (medium, Fig. 4), such as “multiple pizzas” or “ankle-deep”, AdaDiff allocates a higher number of steps, typically around 30. In the case of challenging prompts (hard, Fig. 4), which often involve numerous objects, intricate interactions between them, and diverse detailed descriptions, AdaDiff allocates 40 or 50 steps to achieve satisfactory generation.

### 4.3. Discussion

**Reward function.** Our reward function is designed to generate images of comparable quality to the baseline while utilizing fewer computational resources. Tab. 5 illustrates a comparison of different reward function designs from three main perspectives: 1) utilizing only step savings as a reward; 2) integrating step savings with various image quality metrics; 3) employing different criteria for assessing high-quality images. When the reward function is solely based on step savings, it becomes relatively opaque in terms of image quality, posing challenges in selecting higher-quality images. The introduction of image quality metrics such as

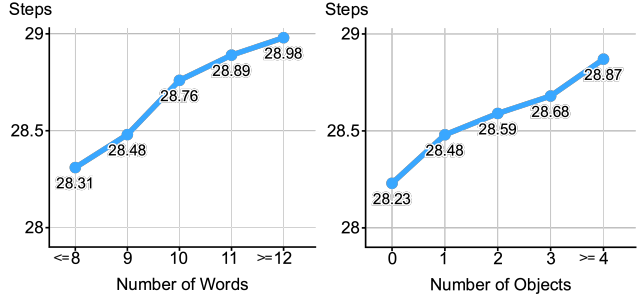


Figure 3. Learned step usage policies for different prompts, which were categorized based on the number of words (left) and objects (right). As the richness of prompt information increases, AdaDiff allocates more steps on average.

CLIP score, NIQE score, and IQS score into the reward function results in varying effectiveness of the step selection strategy. IQS score, which considers both perceptual fidelity and text-image alignment, notably enhances the strategy’s effectiveness. However, when focusing solely on text-image alignment (CLIP score) or perceptual fidelity (NIQE score), despite achieving a lower average step count, there is a more pronounced decrease in image quality. In addition, we compare the impact of employing different methods to determine high-quality images on the step selection strategy. The relative approach ranks results generated by the different numbers of steps according to IQS, considering the top- $k$  as high quality. In contrast, the absolute approach defines high quality when the IQS value exceeds a manually set threshold (*e.g.*, 0). The findings in Tab. 5 indicate that the relative approach learns a more efficient step selection policy.

Step	Reward function				Performance				
	CLIP	NIQE	IQS	Top- $k$	Threshold	Step↓	IQS↑	IS↑	NIQE↓
✓				✓		18.12	0.335	36.58	3.967
✓	✓			✓		18.86	0.334	36.53	3.967
✓		✓		✓		29.62	0.392	37.11	3.799
✓			✓	✓		28.61	<b>0.412</b>	<b>37.60</b>	<b>3.761</b>
✓			✓		✓	25.23	0.377	37.25	3.907

Table 5. Comparisons of different reward functions.

**Different trade-offs between speed and quality.** As mentioned earlier, the parameter top- $k$  is employed to determine whether an image is of high quality. In cases of low quality, a penalty of  $\gamma$  is applied to the strategy. Consequently, the choice of  $k$  controls the trade-off between speed and quality. Here, we present the image quality of AdaDiff at different number of steps. Fig. 5 illustrates that when  $k$  is smaller, imposing stricter requirements for high-quality images leads to an increase in the average IQS, accompanied by a rise in the average number of steps.

It is worth noting that AdaDiff outperforms the random policy in terms of image quality by 15.5% to 17.2%, with



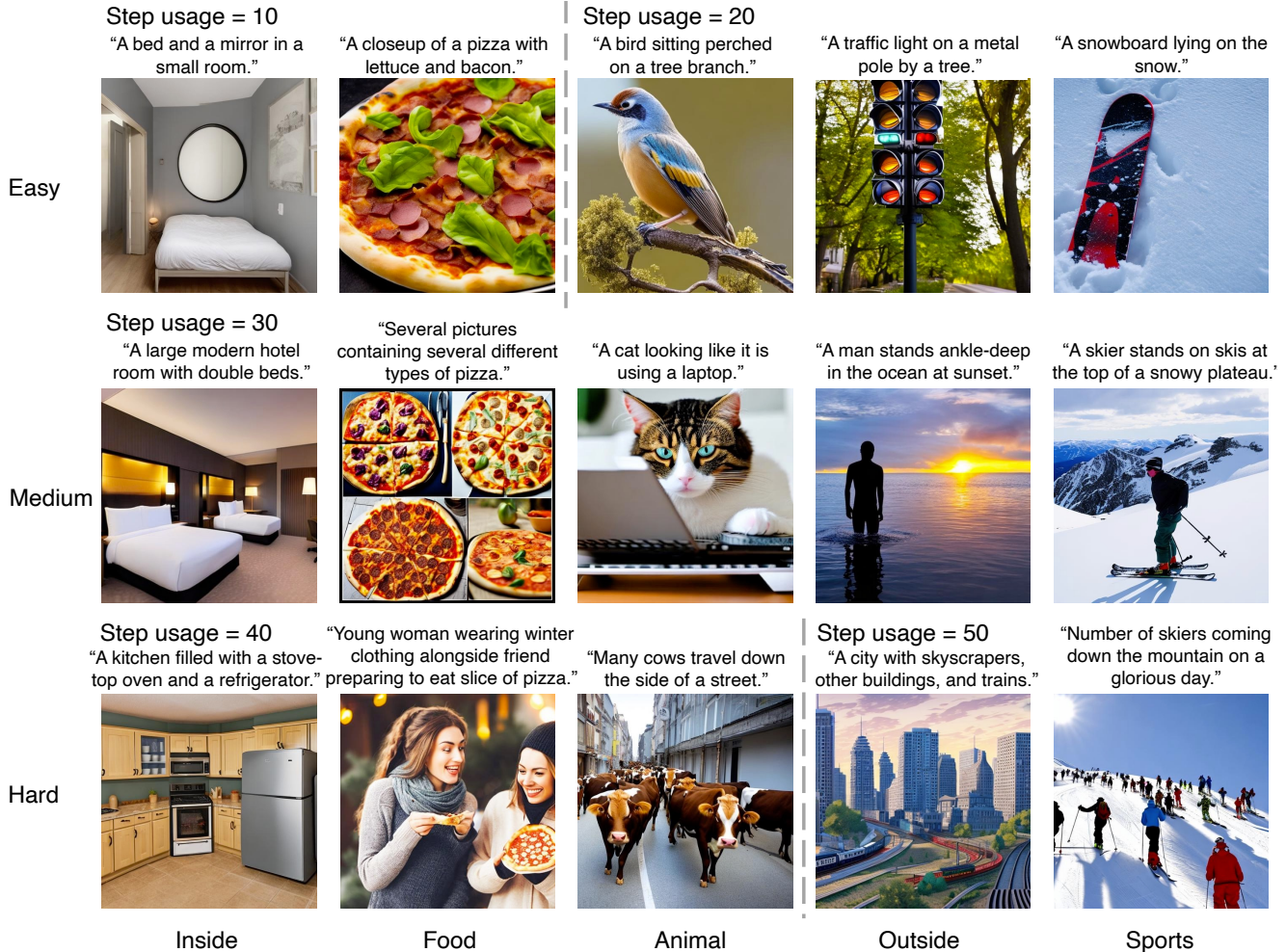


Figure 4. **Qualitative Results.** We categorize prompts into three classes (easy, medium, and hard) based on their information richness, and AdaDiff implements an instance-specific dynamic sampling step policy in five different scenarios (inside, food, animal, outside, and sports).

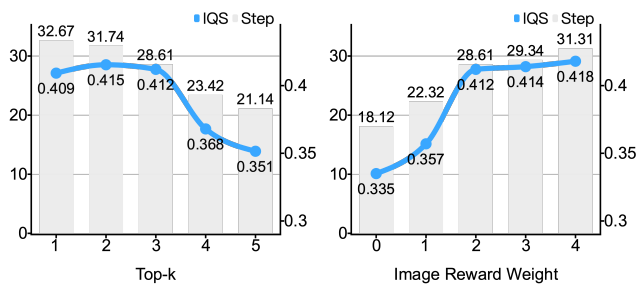


Figure 5. The different trade-offs between speed and quality controlled by top- $k$  (left) and the image reward weight (right).

close average number of steps. Meanwhile, when achieving similar image quality, AdaDiff demonstrates savings of 21.9% and 29.5% in the average number of steps. Additionally, the trade-off between speed and quality can be further influenced by the image reward weight  $\lambda$  in the reward function. A larger  $\lambda$  tends to result in the generation

of higher-quality images, leading to a higher average image quality along with an increase in the average number of steps. Compared to the random policy, with a similar number of steps, AdaDiff outperforms in terms of image quality by 16.4% to 18.1%, while achieving similar image quality results in terms of average step savings by 25.7%.

## 5. Conclusion

In this paper, we introduced AdaDiff, a method that derives adaptive step usage policies tailored to each prompt, facilitating efficient image and video generation. More precisely, a step selection network is trained using policy gradient methods to generate these policies, striking a balance between generation quality and the reduction of overall computational costs. Extensive experiments validated the capability of AdaDiff to generate strong step usage policies on per-input bias, providing compelling qualitative and quantitative evidence.



## References

- [1] Tomer Amit, Tal Shaharbany, Eliya Nachmani, and Lior Wolf. Segdiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*, 2021. 2
- [2] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023. 2
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *CVPR*, 2023. 2
- [4] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. Diffusiondet: Diffusion model for object detection. In *ICCV*, 2023. 2
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021. 1, 2
- [6] Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023. 2
- [7] Ying Fan and Kangwook Lee. Optimizing ddp sampling with shortcut fine-tuning. *arXiv preprint arXiv:2301.13362*, 2023. 2
- [8] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models. *arXiv preprint arXiv:2305.16381*, 2023. 2
- [9] Sicheng Gao, Xuhui Liu, Bohan Zeng, Sheng Xu, Yanjing Li, Xiaoyan Luo, Jianzhuang Liu, Xiantong Zhen, and Baochang Zhang. Implicit diffusion models for continuous super-resolution. In *CVPR*, 2023. 2
- [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 5
- [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1, 2, 3
- [12] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- [13] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video diffusion models. In *NeurIPS*, 2022. 2
- [14] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022. 2
- [15] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *arXiv preprint arXiv:2305.01569*, 2023. 6
- [16] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023. 2
- [17] Sangyun Lee, Beomsu Kim, and Jong Chul Ye. Minimizing trajectory curvature of ode-based generative models. *arXiv preprint arXiv:2301.12003*, 2023. 2
- [18] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *CVPR*, 2023. 2
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 2, 5
- [20] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *ICLR*, 2021. 2
- [21] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. In *NeurIPS*, 2022. 1, 2, 6
- [22] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021. 2
- [23] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *CVPR*, 2023. 2, 5
- [24] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a “completely blind” image quality analyzer. *IEEE Signal processing letters*, 2012. 5
- [25] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaoju Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2
- [26] Alexander Quinn 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 *ICML*, 2022.
- [27] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *CVPR*, 2023. 2
- [28] Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023. 2
- [29] 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 *ICML*, 2021. 5
- [30] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2
- [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 2, 3

- [32] 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 *CVPR*, 2023.
- [33] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 2
- [34] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 2
- [35] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *NeurIPS*, 2016. 5
- [36] Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont. Laioncoco. <https://laion.ai/blog/laion-coco/>, 2022. 2, 5
- [37] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. In *ICLR*, 2023. 2
- [38] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 1, 2, 3
- [39] Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion models for inverse problems. In *ICLR*, 2022. 2
- [40] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. In *NeurIPS*, 2020. 2
- [41] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018. 2, 4
- [42] Bram Wallace, Akash Gokul, Stefano Ermon, and Nikhil Naik. End-to-end diffusion latent optimization improves classifier guidance. In *ICCV*, 2023. 2
- [43] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. 1, 2, 3, 5
- [44] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, et al. Internvid: A large-scale video-text dataset for multimodal understanding and generation. *arXiv preprint arXiv:2307.06942*, 2023. 2, 5
- [45] Zijie J Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. Diffusiondb: A large-scale prompt gallery dataset for text-to-image generative models. *arXiv preprint arXiv:2210.14896*, 2022. 2, 5
- [46] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *ICCV*, 2023. 2
- [47] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference. In *ICCV*, 2023. 2, 6
- [48] Bin Xia, Yulun Zhang, Shiyin Wang, Yitong Wang, Xinglong Wu, Yapeng Tian, Wenming Yang, and Luc Van Gool. Diffir: Efficient diffusion model for image restoration. In *ICCV*, 2023. 2
- [49] Zhen Xing, Qi Dai, Han Hu, Zuxuan Wu, and Yu-Gang Jiang. Simda: Simple diffusion adapter for efficient video generation. *arXiv preprint arXiv:2308.09710*, 2023. 2, 5
- [50] Zhen Xing, Qijun Feng, Haoran Chen, Qi Dai, Han Hu, Hang Xu, Zuxuan Wu, and Yu-Gang Jiang. A survey on video diffusion models. *arXiv preprint arXiv:2310.10647*, 2023. 2
- [51] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, 2016. 2, 5
- [52] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *ICCV*, 2023. 2
- [53] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. In *NeurIPS*, 2023. 2, 4, 5, 6
- [54] Xingqian Xu, Zhangyang Wang, Gong Zhang, Kai Wang, and Humphrey Shi. Versatile diffusion: Text, images and variations all in one diffusion model. In *ICCV*, 2023. 2
- [55] Shengming Yin, Chenfei Wu, Huan Yang, Jianfeng Wang, Xiaodong Wang, Minheng Ni, Zhengyuan Yang, Linjie Li, Shuguang Liu, Fan Yang, et al. Nuwa-xl: Diffusion over diffusion for extremely long video generation. *arXiv preprint arXiv:2303.12346*, 2023. 2
- [56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 2
- [57] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. *arXiv preprint arXiv:2302.04867*, 2023. 2
- [58] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022. 2