

CFG-Zero*: Improved Classifier-Free Guidance for Flow Matching Models

Weichen Fan¹ Amber Yijia Zheng² Raymond A. Yeh² Ziwei Liu^{1,*}

¹S-Lab, Nanyang Technological University

²Department of Computer Science, Purdue University

weichen.fan@u.nus.edu, {zheng709, rayyeh}@purdue.edu,

zwei.liu@ntu.edu.sg

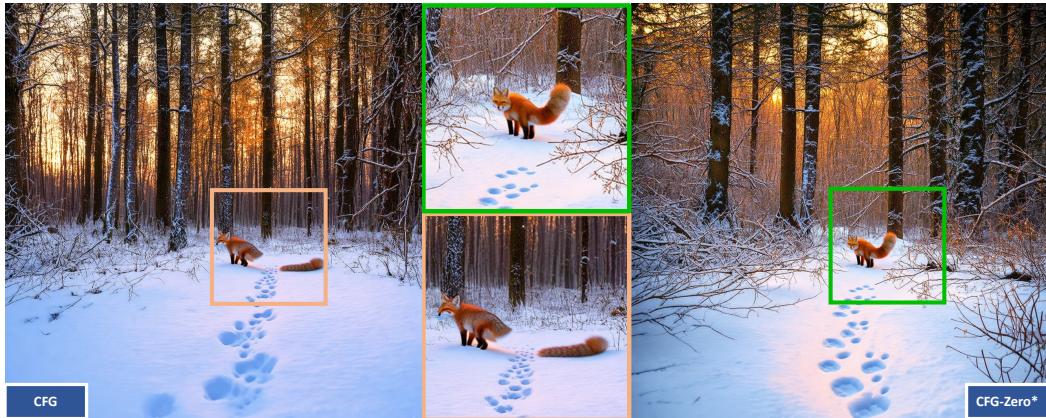


Figure 1. Comparison for the prompt: “A dense winter forest with snow-covered branches, the golden light of dawn filtering through the trees, and a lone fox leaving delicate paw prints in the fresh snow.” Images generated using SD3.5 [5] with CFG and CFG-Zero* (Ours).

Abstract

Classifier-Free Guidance (CFG) is a widely adopted technique in diffusion/flow models to improve image fidelity and controllability. In this work, we first analytically study the effect of CFG on flow matching models trained on Gaussian mixtures where the ground-truth flow can be derived. We observe that in the early stages of training, when the flow estimation is inaccurate, CFG directs samples toward incorrect trajectories. Building on this observation, we propose CFG-Zero, an improved CFG with two contributions: (a) optimized scale, where a scalar is optimized to correct for the inaccuracies in the estimated velocity, hence the * in the name; and (b) zero-init, which involves zeroing out the first few steps of the ODE solver. Experiments on text-to-image (Lumina-Next, Stable Diffusion 3, and Flux) and text-to-video (Wan-2.1) generation demonstrate that CFG-Zero* consistently outperforms CFG, highlighting its effectiveness in guiding Flow Matching models. (Code is available at github.com/WeichenFan/CFG-Zero-star)*

1. Introduction

Diffusion and flow-based models are the state-of-the-art (SOTA) for generating high-quality images and videos, with

recent advancements broadly categorized into Score Matching [10, 27, 30, 34–36] and Flow Matching [5, 6, 20, 24, 45] approaches. Flow matching directly predicts a velocity, enabling a more interpretable transport process and faster convergence compared with score-based diffusion methods. Hence, recent SOTA in text-to-image/video models increasingly adopt flow matching. In this paper, we follow the unifying perspective on diffusion and flow models presented by Lipman et al. [22]. We broadly use the term *flow matching models* to refer to any model trained using flow matching, where samples are generated by solving an ordinary differential equation (ODE).

Next, classifier-free guidance (CFG) [9, 44] is a widely used technique in flow matching models to improve sample quality and controllability during generation. In text-to-image tasks, CFG improves the alignment between generated images and input text prompts. In other words, CFG is used because the conditional distribution induced by the learned conditional velocity does not fully match with the user’s “intended” conditional distribution; see example in Fig. 2. We hypothesize that this mismatch arises from two fundamental factors. First, it may be from dataset limitations, where the user’s interpretation of a text prompt and its corresponding image differs from the dataset distribution. Second, it could result from a learning limitation,



Figure 2. **(Left)** Conditional generation. **(Right)** CFG generation. (Prompt: “A mysterious underwater city with bioluminescent corals and towering glass domes.”)

where the learned velocity fails to accurately capture the dataset’s distribution. In this work, we focus on the latter issue. When the model is underfitted, a mismatch exists between the conditional and unconditional predictions during sampling, causing CFG to guide the sample in a direction that deviates significantly from the optimal trajectory. Specifically, the velocity estimated by CFG in the first step at x_0 may contradict the optimal velocity. This suggests that skipping this prediction could lead to better results.

We empirically analyze the effect of CFG when the learned velocity is underfitted, *i.e.*, inaccurate, using a mixture of Gaussians as the data distribution. In this setting, the ground-truth (optimal) velocity has a closed-form solution, allowing us to compare it with the learned velocity throughout training. Based on the observations, we then propose CFG-Zero*, which introduces two key improvements over the vanilla CFG: *optimized scale* and the *zero init* technique.

In Sec. 4, we provide an empirical analysis to motivate our approach with supporting observations. Furthermore, we validate that the observations go beyond mixture of Gaussians by conducting experiments on ImageNet for the task of class-conditional generation. Finally, in Sec. 5, we apply CFG-Zero* to text-guided generation, showing that our method achieves strong empirical performance in guiding SOTA flow matching models.

Our main contributions are summarized as follows:

- We analyze the sources of error in flow-matching models and propose a novel approach to mitigate inaccuracies in the predicted velocity.
- We empirically show that zeroing out the first ODE solver step for flow matching models improves sample quality when the model is underfitted.
- Extensive experiments validate that CFG-Zero* achieves competitive performance in both discrete and continuous conditional generation tasks, demonstrating its effectiveness as an alternative to CFG.

2. Related Work

Diffusion and Flow-based Models. Unlike generative ad-

versarial methods [7] that rely on one-step generation, diffusion models [4] have demonstrated significantly improved performance in generating high-quality samples. Early diffusion models were primarily score-based generative models, including DDPM [10], DDIM [34], EDM [16], and Stable Diffusion [30], which focused on learning the SDEs governing the diffusion process.

Next, Flow Matching [21] provides an alternative approach by directly modeling sample trajectories using ordinary differential equations (ODEs) instead of SDEs. This enables more stable and efficient generative processes by learning a continuous flow field that smoothly transports samples from a prior distribution to the target distribution. Several works, including Rectified Flow [24], SD3 [5], Lumina-Next [45], Flux [20], Vchitect-2.0 [6], Lumina-Video [23] HunyuanVideo [18], SkyReels-v1 [33], and Wan2.1 [39] have demonstrated that ODE-based methods achieve faster convergence and improved controllability in text-to-image and text-to-video generation. As a result, Flow Matching has become a compelling alternative to stochastic diffusion models, offering better interpretability and training stability. Thus, our analysis is based on Flow Matching models, which aim to provide more accurate classifier-free guidance.

Guidance in Diffusion Models. Achieving better control over diffusion models remains challenging yet essential. Early approaches, such as classifier guidance (CG) [4], introduce control by incorporating classifier gradients into the sampling process. However, this method requires separately trained classifiers, making it less flexible and computationally demanding. To overcome these limitations, classifier-free guidance (CFG) [9] was proposed, enabling guidance without the need for an external classifier. Instead, CFG trains conditional and unconditional models simultaneously and interpolates between their outputs during sampling.

Despite its effectiveness, CFG relies on an unbounded empirical parameter, known as the guidance scale, which determines how strongly the generated output is influenced by the conditional model. Improper tuning of this scale can lead to undesirable artifacts, either over-saturated outputs with excessive conditioning or weakened generation fidelity due to under-conditioning. In fact, previous studies [1] have observed that CFG estimation does not provide an optimal denoising direction.

As a result, several studies have explored adaptive or dynamically scaled guidance to address these issues, including ADG [31], Characteristic-Guidance [43], ReCFG [41], Weight-Scheduler [40], and CFG++ [2]. Additionally, AutoG [17] replaces the unconditional model with a smaller, less-trained version of the model itself. Other researchers [19] have proposed limiting the use of CFG to a specific interval during sampling.

Unlike previous approaches, our work is motivated by an

observation that CFG prediction is inaccurate when a model is underfitted, and specifically in the first step, *i.e.*, the prediction is even worse than zeroing out the first ODE solver step. Thus, we study the error of CFG within the Flow Matching, where we derive the upper bound of the error term and propose to minimize it. From this analysis, we derive a dynamic parameterization technique that adjusts the unconditional output, leading to more stable and effective guidance in Flow Matching based diffusion models.

3. Preliminaries

We briefly recap flow matching following the unifying perspective presented in Lipman et al. [22].

Conditional flow matching. Given a source distribution $p(x|y)$ and an unknown target distribution $q(x|y)$, conditional flow matching (CFM) defines a probability path $p_t(x|y)$, where $t \in [0, 1]$ is a continuous time variable that interpolates between p and q , such that $p_0(x|y) = p(x|y)$ and $p_1(x|y) = q(x|y)$. An effective choice for $p_t(x|y)$ is the linear probability path

$$p_t(x|y) \triangleq (1-t) \cdot p(x|y) + t \cdot q(x|y), \quad (1)$$

where a sample $x_t = (1-t)x_0 + tx_1$, with $x_0 \sim p(x|y)$ and $x_1 \sim q(x|y)$. Next, a continuous flow model is trained by learning a time-dependent velocity field $\frac{d}{dt}x_t = \mathbf{v}_t^\theta(x|y)$ that governs the trajectory of x over t . Here, the velocity is represented using a deep-net with trainable parameters θ . A flow matching model is trained by minimizing the CFM loss expressed as

$$L_{\text{CFM}}(\theta) = \mathbb{E}_{t,x_0,x_1} \|\mathbf{v}_t^\theta(x_t|y) - (x_1 - x_0)\|_2^2. \quad (2)$$

At generation time, a new sample can be obtained by using any ODESolver, *e.g.*, the midpoint method [37].

Classifier free guidance (CFG) [9, 44] improves the quality of conditional generation by steering a sample toward the given input condition, *e.g.*, a class label or a text prompt. In CFG, a single flow model $\mathbf{v}_t^\theta(x|y)$ is trained to output both conditional and unconditional velocity fields. This is done by introducing $y = \emptyset$, which does not contain any conditioning information.

At inference, the guided velocity field is formed by

$$\hat{\mathbf{v}}_t^\theta(x|y) \triangleq (1-w) \cdot \mathbf{v}_t^\theta(x|y = \emptyset) + w \cdot \mathbf{v}_t^\theta(x|y), \quad (3)$$

where w is the guidance scale and \emptyset denotes the null condition. When $w = 1$, it is equivalent sampling only using the conditional velocity $\mathbf{v}_t^\theta(x|y)$, *i.e.*, no guidance.

Closed form velocity for Gaussian distributions. When both the source and target consist of Gaussian mixtures, the optimal velocity has a closed form.

For example, let the source distribution $p = \mathcal{N}(0, \mathbf{I})$ and target distribution $q = \mathcal{N}(\mu, \mathbf{I})$ both be a single Gaussian.

Algorithm 1 CFG-Zero*

```

1: Input: Trained velocity  $\mathbf{v}^\theta$ , noise sample  $x_0$ , guidance weight  $\omega$ , and number of steps to zero out  $K$ .
   #  $K$  equals to 1 by default
2:  $s_t^* \leftarrow \frac{\mathbf{v}_t^\theta(\cdot|y)^\top \mathbf{v}_t^\theta(\cdot|\emptyset)}{\|\mathbf{v}_t^\theta(\cdot|\emptyset)\|^2}$  # Optimized scale
3:  $\tilde{\mathbf{v}}_t(\cdot) \leftarrow (1 - \omega) \cdot s_t^* \cdot \mathbf{v}_t^\theta(\cdot|\emptyset) + \omega \cdot \mathbf{v}_t^\theta(\cdot|y)$ 
4: # Solve ODE
5: for  $t = 0$  to  $T$  do
6:   if  $t < K$  then
7:      $x_{t+1} \leftarrow x_t$  # Zero-init
8:   else
9:      $x_{t+1} \leftarrow \text{ODEStep}(\tilde{\mathbf{v}}_t(\cdot), x_t)$ 
10:  end if
11: end for
12: Return generated sample  $x_T$ 

```

Then, the corresponding optimal velocity $\mathbf{v}_t^*(x) =$

$$[(2t-1)\mathbf{I}] \left[(1-t)^2 \mathbf{I} + t^2 \mathbf{I} \right]^{-1} (x - t\mu) + \mu. \quad (4)$$

With a closed-form optimal \mathbf{v}_t^* , we can now empirically study the gap between the optimal and learned velocity $\|\mathbf{v}_t^*(\cdot) - \mathbf{v}_t^\theta(\cdot)\|$ throughout training, specifically in the underfitting regime, which motivated our proposed CFG-Zero*.

4. Methodology

We propose a guidance algorithm CFG-Zero* with two improvements from the standard CFG: **(a) optimized scale**, where a scalar parameter is optimized to compensate for inaccuracies in the learned velocity (Sec. 4.1); **(b) zero-init**, where we zero out the first step of the ODE solver (Sec. 4.2). These modifications can be easily integrated into the existing CFG code base and introduce minimal additional computational cost. The overall algorithm is summarized in Alg. 1.

4.1. Optimizing an additional scaler in CFG

As motivated in the introduction, we aim to study the use of CFG in the setting where the velocity is underfitted, *i.e.*, CFG is used in the hope that the guided velocity $\tilde{\mathbf{v}}^\theta$ approximates the ground-truth flow $\tilde{\mathbf{v}}^*$, *i.e.*,

$$\tilde{\mathbf{v}}_t^\theta(x|y) \approx \mathbf{v}_t^*(x|y). \quad (5)$$

To further improve this approximation, we introduce an optimizable scaler $s \in \mathbb{R}_{>0}$ to CFG,

$$\tilde{\mathbf{v}}_t^\theta(x|y) \triangleq (1-\omega) \cdot s \cdot \mathbf{v}_t^\theta(x) + \omega \cdot \mathbf{v}_t^\theta(x|y) \quad (6)$$

$$= -\omega' \cdot s \cdot \mathbf{v}_t^\theta(x) + (1+\omega')\mathbf{v}_t^\theta(x|y), \quad (7)$$

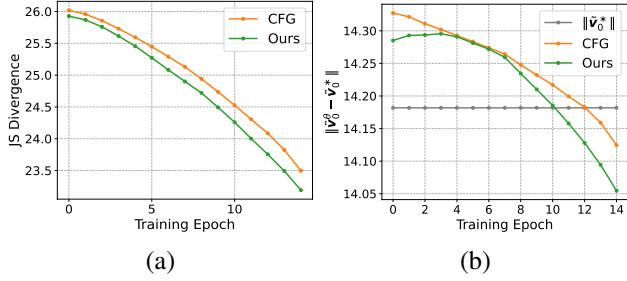


Figure 3. Results on mixture of Gaussians in \mathbb{R}^2 . **Left:** The Jensen–Shannon divergence between the model’s final flow sample distribution and the target distribution *v.s.* training epoch. **Right:** The velocity error norm $\|\tilde{v}_t^\theta - \tilde{v}_t^*\|$, with the ground truth norm shown in gray *v.s.* training epoch.

where $\omega' \triangleq 1 + \omega$. The choice of learning a scalar s is inspired by classifier guidance [4], where they introduce a scaling factor to balance between the gradient and the unconditioned direction. The remaining challenge is how to optimize s .

In a hypothetical case where we do have ground-truth flow \mathbf{v}_t^* , then one can formulate the approximation in Eq. (5) as a least-squares, *i.e.*, minimizing s over the loss

$$\mathcal{L}(s) \triangleq \|\tilde{\mathbf{v}}_t^\theta(x|y) - \mathbf{v}_t^*(x|y)\|_2^2. \quad (8)$$

However, as \mathbf{v}_t^* is unknown, we instead minimize over an upperbound of $\mathcal{L}(s)$ in Eq. (8) established using triangle inequality as follows:

$$\begin{aligned} \mathcal{L}(s) &= \|-\omega' s \mathbf{v}_t^\theta(x) + (1 + \omega') \mathbf{v}_t^\theta(x|y) - \mathbf{v}_t^*(x|y)\|_2^2 \\ &\leq \|\mathbf{v}_t^\theta(x|y)\|_2^2 + \|\mathbf{v}_t^*(x|y)\|_2^2 \\ &\quad + \omega' \|\mathbf{v}_t^\theta(x|y) - s \cdot \mathbf{v}_t^\theta(x)\|_2^2. \end{aligned} \quad (9)$$

Observe that only the last term has dependencies on s , *i.e.*, optimizing Eq. (9) is equivalent to

$$\min_s \|\mathbf{v}_t^\theta(x|y) - s \cdot \mathbf{v}_t^\theta(x)\|_2^2, \quad (10)$$

where the solution s^* is a projection of the condition velocity onto the unconditional velocity, *i.e.*,

$$s^* = (\mathbf{v}_t^\theta(x|y)^\top \mathbf{v}_t^\theta(x)) / \|\mathbf{v}_t^\theta(x)\|^2. \quad (11)$$

We empirically validate the approach on a toy example consisting of a Gaussian mixture, with results shown in Fig. 3(a), where we observe that samples generated with CFG-Zero* more closely match the target distribution than those with CFG.

4.2. Zero-init for ODE Solver

During the empirical validation of s^* , we also studied how well the guided velocity from CFG-Zero* matches the

ground-truth. As shown in Fig. 3(b), during the early stages of training, the difference between both types of estimated velocities and the ground-truth velocity at $t = 0$ is greater than when just using an all-zero velocity (*do nothing*), *i.e.*,

$$\|\tilde{\mathbf{v}}_0^\theta(x|y) - \mathbf{v}_0^*(x|y)\|_2^2 \geq \|\mathbf{0} - \mathbf{v}_0^*(x|y)\|_2^2. \quad (12)$$

Based on this observation, we propose zero-init, which zeros out the velocity for the first few steps of the ODESolver. As an estimation of all zeros for the velocity would be more accurate. Next, we investigate whether this behavior is specific to the small scale mixture of Gaussians dataset or could be generalized to real data. Specifically, we consider the experiment of ImageNet-256.

4.3. Validation beyond Mixture of Gaussians

We experiment on the ImageNet-256 [3] benchmark where we train a class-conditioned flow model, *i.e.*, given an input class label, generate an image from the class. We report results across multiple evaluation metrics, including *Inception Score* (IS) [32], *Fréchet Inception Distance* (FID) [8], *sFID* [26], *Precision*, and *Recall*.

Validating Zero-Init. To study the impact of zero-init on a real dataset with flow matching models, we compare samples generated by the standard CFG with and without the zero-init throughout the training epochs, as shown in Tab. 1. Theoretically, if a model is well-trained, *zero-init* could degrade its performance. For quick validation, we use a smaller version of DiT [27] with 400M parameters as our base model and train it from scratch on the ImageNet-256 dataset using the Flow Matching loss [21].

The results show that CFG with *Zero-Init* consistently outperforms standard CFG up to 160 epochs. This indicates that zeroing out the first step of the ODESolver is suitable during early training. Beyond 160 epochs, standard CFG surpasses Zero-Init, which is consistent with our hypothesis that as the velocity is trained to be more accurate, the benefit of zero-init lessens.

Validation of CFG-Zero*. To highlight the differences between our method and other classifier-free guidance approaches, we conduct experiments using a pre-trained SiT-XL model [25] (700M parameters) on the ImageNet-256 benchmark. The model has not yet reached the "turning point" mentioned earlier. All experiments are performed using the standard guidance scale.

In Tab. 2, our results demonstrate that CFG-Zero* achieves the best overall performance, outperforming both CFG++ [2] and ADG [31] across key metrics. Specifically, CFG-Zero* attains the highest *Inception Score* of 258.87, highlighting its ability to generate diverse and high-quality images. Furthermore, CFG-Zero* achieves the best *FID Score* of 2.10 and *sFID Score* of 4.59, indicating improved perceptual quality and stronger alignment with the target

Epochs	Methods	Metrics				
		IS↑	FID↓	sFID↓	Precision↑	Recall↑
10	CFG	53.27	28.57	18.52	0.61	0.36
	Zero-Init	52.78	28.55	17.32	0.62	0.37
20	CFG	257.23	11.00	11.64	0.92	0.24
	Zero-Init	255.79	10.65	10.95	0.92	0.25
40	CFG	339.39	12.61	11.17	0.94	0.23
	Zero-Init	338.40	12.29	10.47	0.94	0.24
80	CFG	383.06	13.53	10.99	0.94	0.24
	Zero-Init	383.45	12.18	10.39	0.94	0.26
160	CFG	222.13	2.84	4.56	0.81	0.56
	Zero-Init	218.90	2.85	4.97	0.80	0.56

Table 1. Validation on ImageNet-256. We evaluate a model at different training stages and observe a turning point at 160 epochs, where zero-init results in poorer performance when the model converges. This experiment validates that high-dimensional models also suffer from inaccuracies in initial sampling.

Method	IS↑	FID↓	sFID↓	Precision↑	Recall↑
Baseline	125.13	9.41	6.40	0.67	0.67
w/ CFG	257.03	2.23	4.61	0.81	0.59
w/ ADG [31]	257.92	2.37	5.51	0.80	0.58
w/ CFG++ [2]	257.04	2.25	4.65	0.79	0.57
w/ CFG-Zero*	258.87	2.10	4.59	0.80	0.61

Table 2. Comparison of different guidance strategy on ImageNet-256 benchmark. Lower FID is better (\downarrow) and higher IS is better (\uparrow). Baseline here denotes using the conditional prediction only.

distribution. In terms of fidelity metrics, our method maintains a competitive *Precision* of 0.80—comparable to both CFG and ADG—while achieving the highest *Recall* of 0.61. This suggests that our approach better captures the underlying data distribution, leading to more representative and well-balanced samples. (Both ADG [31] and CFG++ [2] are not designed for classifier-free guidance in Flow Matching, and therefore may not perform well.)

5. Experiments

In this section, we evaluate CFG-Zero* on large-scale models for text-to-image (Lumina-Next [45], SD3 [5], SD3.5 [5], Flux [20]) and text-to-video (Wan2.1 [39]) generation. Note: Flux is CFG-distilled, so directly applying classifier-free guidance may yield different results.

5.1. Text-to-Image Generation

To evaluate the effectiveness of our proposed method, CFG-Zero*, in continuous class-conditional image generation, we conducted experiments using four state-of-the-art flow matching models: Lumina-Next [45], SD3 [5], SD3.5 [5], and Flux [20]. These models were selected for their strong performance in class-conditional image synthesis. We applied both CFG-Zero* and the standard

Model	Method	Aesthetic Score↑	CLIP Score↑
Lumina-Next [45]	CFG	6.85	34.09
	CFG-Zero*	7.03	34.37
SD3 [5]	CFG	6.73	34.00
	CFG-Zero*	6.80	34.11
SD3.5 [5]	CFG	6.96	34.60
	CFG-Zero*	7.10	34.68
Flux [20]	CFG	7.06	34.60
	CFG-Zero*	7.12	34.69

Table 3. Quantitative evaluation of Text-to-Image generation, using Lumina-Next, Stable Diffusion 3, Stable Diffusion 3.5, and Flux. The evaluation is based on *Aesthetic Score* and *CLIP Score* as key metrics. Results indicate that CFG-Zero* consistently enhances image quality and improves alignment with textual prompts across different models.

CFG under default settings on a diverse set of self-curated prompts, allowing us to provide fair comparisons.

Quantitative Evaluation. Tab. 3 presents the quantitative comparison of CFG-Zero* and CFG across all tested models. The results indicate that CFG-Zero* consistently achieves superior performance, as evidenced by higher *Aesthetic Score* [38] and *CLIP Score* [28, 29]. The improvement in *Aesthetic Score* suggests that CFG-Zero* enhances the visual appeal of generated images, producing outputs with more coherent textures, lighting, and structure. Additionally, the increase in *CLIP Score* demonstrates that CFG-Zero* improves text-image alignment, ensuring that generated images better capture the semantics of the given prompts. These results validate the effectiveness of our proposed modifications in refining the quality of diffusion-based generation.

Qualitative Evaluation. Fig. 4 provides a comparisons of the images generated using CFG-Zero* and vanilla CFG. Our method produces high-fidelity outputs that exhibit richer details, sharper textures, and better preservation of object structures compared to the baseline CFG. Notably, CFG-Zero* mitigates common artifacts observed in CFG-generated images, particularly those that introduce unintended distortions or elements unrelated to the given prompt. This reduction in artifacts highlights the robustness of CFG-Zero* in preserving semantic consistency, ensuring that generated images adhere more closely to the given prompts.

Additionally, we observe that CFG-Zero* have better color consistency and level of detail, reducing blurry artifacts that are sometimes present in CFG-based outputs. These improvements are particularly evident in complex prompts that require the precise rendering of intricate textures or fine-grained semantic attributes. Further visual comparisons and additional generated samples can be found in the Appendix.

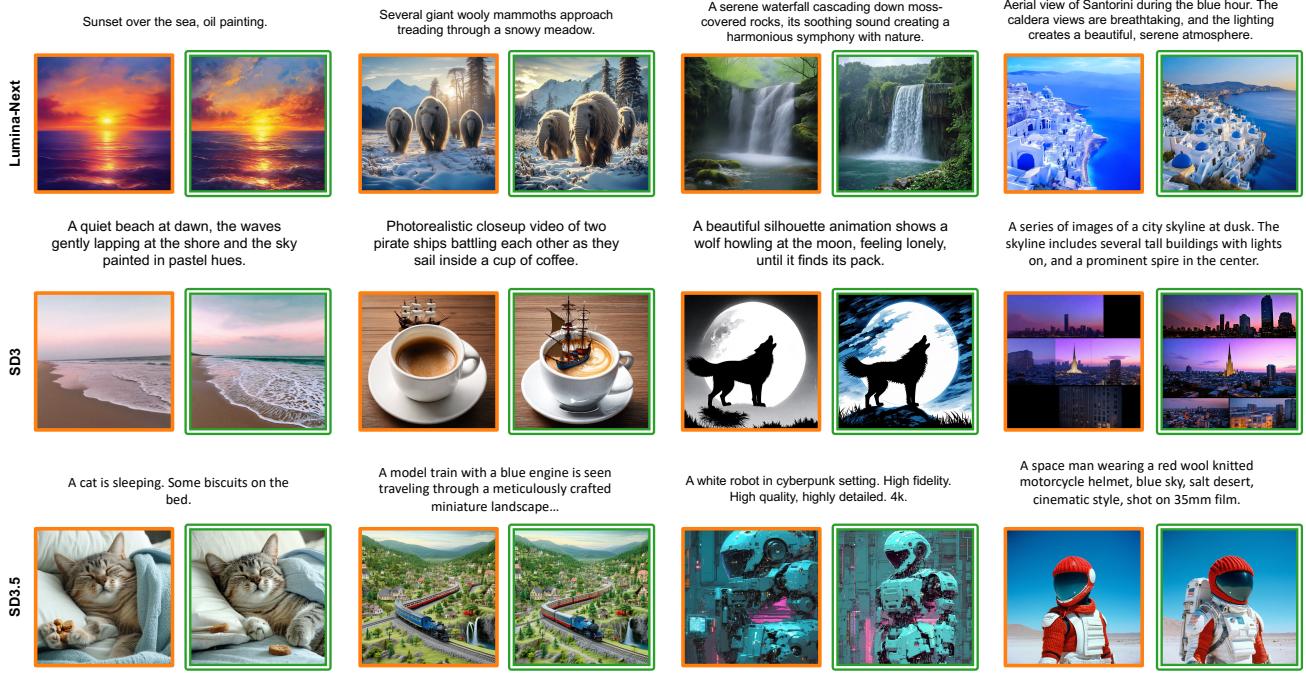


Figure 4. **Qualitative comparisons between CFG and CFG-Zero^{*}.** Experiments are conducted using Lumina-Next, Stable Diffusion 3, and Stable Diffusion 3.5, with each model evaluated under its recommended optimal sampling steps and guidance scale settings. CFG results are shown in orange and Ours are highlighted in green boxes.

Method	Color↑	Shape↑	Texture↑	Spatial↑
Lumina-Next [45]	0.51	0.34	0.41	0.19
+ CFG-Zero [*]	0.52	0.36	0.45	0.29
SD3 [5]	0.81	0.57	0.71	0.31
+ CFG-Zero [*]	0.83	0.58	0.72	0.31
SD3.5 [5]	0.76	0.59	0.70	0.27
+ CFG-Zero [*]	0.78	0.60	0.71	0.28

Table 4. **Quantitative evaluation on T2I-CompBench [12],** using Lumina-Next, Stable Diffusion 3, and Stable Diffusion 3.5. Compared to CFG, CFG-Zero^{*} demonstrates consistent improvements across all evaluated aspects.

Benchmark Results. We compare our method against standard CFG across three different Flow Matching models using T2I-CompBench [11, 12]. As shown in Tab. 4, integrating CFG-Zero^{*} leads to notable improvements in *Color*, *Shape*, and *Texture* quality in generated images. Meanwhile, the *Spatial* dimension remains comparable to the baseline, indicating that CFG-Zero^{*} improves image fidelity without compromising structural coherence.

User Study. To further assess CFG-Zero^{*}, we conduct a user study to compare its performance against standard CFG across various flow matching models. Participants were presented with image pairs generated using both CFG-Zero^{*} and CFG and were asked to evaluate them based on three key aspects: *detail preservation*, *color con-*

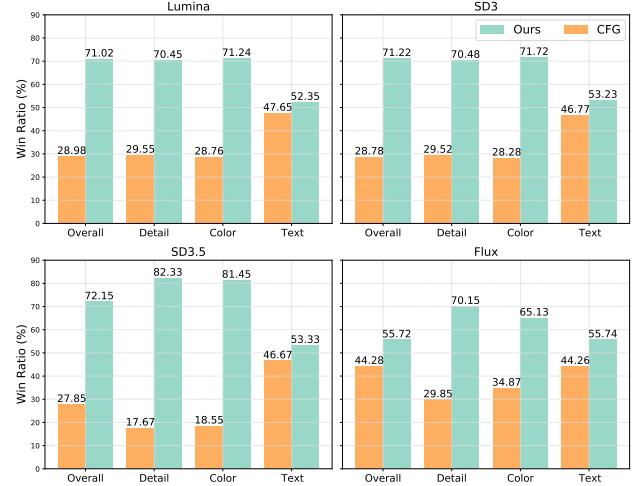


Figure 5. **User study on Lumina-Next, Stable Diffusion 3, Stable Diffusion 3.5, and Flux.** The win rate of our method compared to CFG is presented.

sistency, and *image-text alignment*. The overall preference score was then computed as the percentage of times CFG-Zero^{*} was favored over CFG.

Fig. 5 summarizes the results. CFG-Zero^{*} demonstrates a clear advantage over CFG across all tested models. Notably, our method achieves the highest win ratio on SD3.5, with an overall win ratio score of **72.15%**, primar-

Method	Total Score	subject consistency	aesthetic quality	imaging quality	color	spatial relationship	temporal style	motion smoothness
Vchitect-2.0 [2B] [6]	81.57	61.47	65.60	86.87	86.87	54.64	25.56	97.76
CogVideoX-1.5 [5B] [42]	82.17	96.87	62.79	65.02	87.55	80.25	25.19	98.31
Wan2.1 [14B] [39]	83.99	93.33	69.13	67.48	83.43	80.46	25.90	98.05
w/ CFG-Zero*	84.06	93.34	69.22	67.55	85.39	79.28	25.98	98.00
Wan2.1 [1B] [39]	80.52	93.89	61.67	65.40	87.57	72.75	24.13	97.24
w/ CFG-Zero*	80.91	94.93	64.24	68.13	89.36	73.84	23.36	98.16

Table 5. **Qualitative evaluation on VBench [13].** We use the Wan-2.1 [39] model as our base model. Compared to vanilla CFG, CFG-Zero* improves both frame quality and overall video smoothness.



Figure 6. **Qualitative comparisons between CFG-Zero* and CFG.** Experiments are conducted using Wan-2.1 [1B] [39], under its recommended optimal sampling steps and guidance scale settings.

ily driven by significant improvements in *detail preservation* (**82.33%**) and *color fidelity* (**81.45%**). This suggests that CFG-Zero* effectively enhances fine-grained structures and maintains consistent color distributions compared to its baseline. In terms of *image-text alignment*, CFG-Zero* also performs favorably, surpassing CFG in most cases. Specifically, with CFG++, SD3.5 exhibits strong text-image consistency (**53.33%** win ratio), indicating that our method improves coherence between generated images and textual prompts across different architectures.

5.2. Text-to-Video Generation

To further evaluate the effectiveness of our proposed method, we further conduct experiments on the text-to-video generation task, using the most recent state-of-the-art model, Wan-2.1 [39].

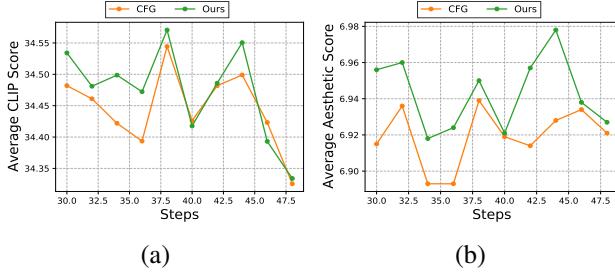
Benchmark Results. As shown in Tab. 5, we evaluate our method using all metrics from VBench [13, 14]. Compared to CFG, Wan-2.1 [1.3B] equipped with CFG-Zero* achieves a higher *Total Score*. Specifically, our method

improves *Aesthetic Quality* by **2.57** and *Imaging Quality* by **2.73**, indicating that CFG-Zero* enhances video fidelity. Additionally, CFG-Zero* improves *Motion Smoothness* (**+0.92**) and *Spatial Relationship* (**+1.09**), demonstrating superior temporal coherence and spatial understanding. However, we also observe a decrease in *Temporal Style* (**-0.77**), which can be attributed to the base model's poor capability in generating stylized videos.

Qualitative Evaluation. Fig. 6 presents a visual comparison between CFG-Zero* and CFG on Wan2.1-14B [39], demonstrating that the refined velocity produced by our method leads to more plausible content with natural motion.

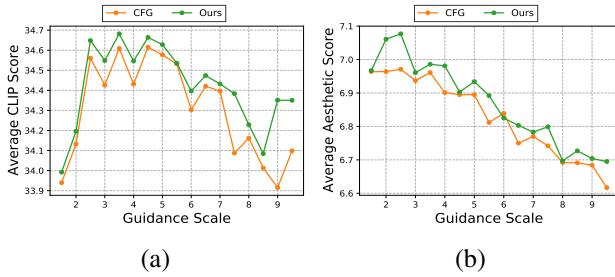
5.3. Ablation Studies

Different Sampling Steps. To further investigate the impact of sampling steps, we conduct experiments using a strong baseline, SD3.5. As shown in Fig. 7, our method consistently enhances both image quality and text-matching accuracy across different sampling steps, demonstrating its



(a) (b)

Figure 7. **Ablation study on different sampling steps.** Comparison of *CLIP Score* and *Aesthetic Score* between our method and CFG across different sampling steps.



(a) (b)

Figure 8. **Ablation study on different guidance scale.** Comparison of *CLIP Score* and *Aesthetic Score* between our method and CFG across different guidance scale.

robustness in varied sampling settings.

Different Guidance Scale. As shown in Fig. 8, we conduct an experiment using SD3.5 to analyze the impact of different guidance scales on our method. The results in sub-figure (a) demonstrate that CFG-Zero* achieves higher CLIP scores than CFG across all guidance scales, validating its effectiveness in improving image-text alignment. Similarly, sub-figure (b) shows that our method consistently attains higher aesthetic scores, highlighting its ability to generate visually appealing images.

Effectiveness of CFG-Zero*. We conduct an ablation study using SD3.5 as the baseline to assess the impact of each component in CFG-Zero*. As shown in Tab. 6, we compare four variants: vanilla CFG, CFG with zero-init, pure optimized scaler, and CFG-Zero*. The optimized scaler reduces the gap between predicted and true velocity, enhancing stability, while zero-init improves performance by skipping the first step. Combining both, CFG-Zero* achieves the highest Aesthetic Score (7.10) and CLIP Score (34.68), outperforming all variants. These results confirm that both modifications help to enhance image quality and text alignment.

Effect of Zeroing Out Initial Steps. To assess whether extending the zero-out strategy beyond the first step can further improve performance, we conduct an ablation study using Lumina-Next, SD3, and SD3.5, as shown in Tab. 7.

Metrics	w/ CFG	w/ CFG-Zero	w/ Scaler	w/ CFG-Zero*
Aesthetic Score	6.96	7.00	6.96	7.10
CLIP Score	34.60	34.65	34.64	34.68

Table 6. **Effectiveness of CFG-Zero*.** Comparison of vanilla CFG, CFG with zero-init, dynamic scaling, and CFG-Zero*, highlighting the impact of *zero-init* and *dynamic scaling* in improving performance.

Model	Zero-out / Total (steps)	Aesthetic Score↑	Clip Score↑
		First 3 / 30	6.78
Lumina-Next [45]	First 2 / 30	7.06	34.73
	First 1 / 30	7.03	34.37
	First 3 / 28	6.95	34.01
SD3 [5]	First 2 / 28	6.98	34.33
	First 1 / 28	6.80	34.11
	First 3 / 28	6.78	34.02
SD3.5 [5]	First 2 / 28	6.99	34.54
	First 1 / 28	7.10	34.68

Table 7. **Ablation study on zero-out steps.** For SD3.5 [5], more initial zero-out steps lead to worse performance, while Lumina-Next [45] and SD3 [5] achieve the highest *Aesthetic Score* and *Clip Score* with first 7% zero out.

Resolution	81x1280x720	81x832x480	1024x1024	512x512
FLOPs	19.4 M	0.84 M	$1.6e^{-3}$ M	$4.1e^{-4}$ M
Memory	18.46 MB	8.00 MB	64 KB	16 KB

Table 8. **Computational costs.** FLOPs [15] and GPU memory usage of our method for 5-second video generation at 720p/480p using Wan2.1 [39], and at 1024/512 resolution using SD3 [5].

The results indicate that for some models, zeroing out can be beneficial not only in the first step but also in the initial few steps. Specifically, both Lumina-Next and SD3 achieve optimal performance when the first 2 steps are zeroed out. However, SD3.5 exhibits a decline in performance when a higher proportion of initial steps are zeroed out, suggesting that SD3.5 is well-trained and approaching convergence.

Computational Cost. Tab. 8 summarizes the FLOPs and GPU memory usage of our method. For 720p and 480p resolutions using Wan2.1, our method requires 19.4M and 0.84M FLOPs, with memory usage of 18.46MB and 8.00MB. For SD3, the computational cost is significantly lower at 1024×1024 and 512×512, requiring only $1.6e^{-3}$ M and $4.1e^{-4}$ M FLOPs, with memory usage of 64KB and 16KB. These results confirm that CFG-Zero* introduces minimal computational costs.

6. Conclusion

We introduce CFG-Zero*, an improved classifier-free guidance method for flow-matching diffusion models. CFG-Zero* addresses CFG’s limitations with two key techniques: (1) an optimized scale factor for accurate velocity

estimation and (2) a zero-init technique to stabilize early sampling. Theoretical analysis and extensive experiments on text-to-image (SD3.5, Lumina-Next, Flux) and text-to-video (Wan-2.1) models show that CFG-Zero* outperforms standard CFG, achieving higher aesthetic scores, better text alignment, and fewer artifacts. Ablation studies further validate the effectiveness of both components in improving sample quality without significant computational cost.

References

- [1] Arwen Bradley and Preetum Nakkiran. Classifier-free guidance is a predictor-corrector. *arXiv preprint arXiv:2408.09000*, 2024. 2
- [2] Hyungjin Chung, Jeongsol Kim, Geon Yeong Park, Hyelin Nam, and Jong Chul Ye. CFG++: Manifold-constrained classifier free guidance for diffusion models. In *Int. Conf. Learn. Represent.*, 2025. 2, 4, 5
- [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2009. 4
- [4] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *Adv. Neural Inform. Process. Syst.*, 2021. 2, 4
- [5] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Int. Conf. Mach. Learn.*, 2024. 1, 2, 5, 6, 8
- [6] Weichen Fan, Chenyang Si, Junhao Song, Zhenyu Yang, Yinan He, Long Zhuo, Ziqi Huang, Ziyue Dong, Jingwen He, Dongwei Pan, et al. Vchitect-2.0: Parallel transformer for scaling up video diffusion models. *arXiv preprint arXiv:2501.08453*, 2025. 1, 2, 7
- [7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 2020. 2
- [8] 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 *Adv. Neural Inform. Process. Syst.*, 2017. 4
- [9] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 1, 2, 3
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Adv. Neural Inform. Process. Syst.*, 2020. 1, 2
- [11] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2I-CompBench: A comprehensive benchmark for open-world compositional text-to-image generation. In *Adv. Neural Inform. Process. Syst.*, 2023. 6, 13
- [12] Kaiyi Huang, Chengqi Duan, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2I-CompBench++: An Enhanced and Comprehensive Benchmark for Compositional Text-to-Image Generation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2025. 6, 12, 13
- [13] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2024. 7
- [14] Ziqi Huang, Fan Zhang, Xiaojie Xu, Yinan He, Jiashuo Yu, Ziyue Dong, Qianli Ma, Nattapol Chanpaisit, Chenyang Si, Yuming Jiang, Yaohui Wang, Xinyuan Chen, Ying-Cong Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench++: Comprehensive and versatile benchmark suite for video generative models. *arXiv preprint arXiv:2411.13503*, 2024. 7
- [15] Vivek Jain, Charles Schmidt, Paritosh Goyal, Cliff Young, et al. Benchmarking deep learning inference: Characteristics of AI workloads on edge vs. cloud. In *International Symposium on Workload Characterization (IISWC)*. IEEE, 2020. 8
- [16] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Adv. Neural Inform. Process. Syst.*, 2022. 2
- [17] Tero Karras, Miika Aittala, Tuomas Kynkänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine. Guiding a diffusion model with a bad version of itself. In *Adv. Neural Inform. Process. Syst.*, 2024. 2
- [18] Weijie Kong et al. Hunyanvideo: A systematic framework for large video generative models, 2024. 2
- [19] Tuomas Kynkänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. In *Adv. Neural Inform. Process. Syst.*, 2025. 2
- [20] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 1, 2, 5
- [21] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *Int. Conf. Learn. Represent.*, 2023. 2, 4
- [22] Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint arXiv:2412.06264*, 2024. 1, 3
- [23] Dongyang Liu, Shicheng Li, Yutong Liu, Zhen Li, Kai Wang, Xinyue Li, Qi Qin, Yufei Liu, Yi Xin, Zhongyu Li, Bin Fu, Chenyang Si, Yuwen Cao, Conghui He, Ziwei Liu, Yu Qiao, Qibin Hou, Hongsheng Li, and Peng Gao. Lumina-video: Efficient and flexible video generation with multi-scale Next-DiT, 2025. 2
- [24] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *Int. Conf. Learn. Represent.*, 2023. 1, 2
- [25] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. SiT: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In *Eur. Conf. Comput. Vis.*, 2024. 4
- [26] Pierre Nguyen, Arthur Leclaire, Léo Gautron, Philippe Robert, and Nicolas Papadakis. Sliced Wasserstein generative models. In *Int. Conf. Learn. Represent.*, 2022. 4
- [27] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Int. Conf. Comput. Vis.*, 2023. 1, 4

- [28] William Peebles, Shuang Li, Michal Lukasiewicz, Richard Zhang, Alexei A Efros, and Eli Shechtman. Image reward: Learning and evaluating human preferences for text-to-image generation. In *Adv. Neural Inform. Process. Syst.*, 2023. 5
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Int. Conf. Mach. Learn.*, 2021. 5
- [30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 1, 2
- [31] Seyedmorteza Sadat, Otmar Hilliges, and Romann M Weber. Eliminating oversaturation and artifacts of high guidance scales in diffusion models. In *Int. Conf. Learn. Represent.*, 2024. 2, 4, 5
- [32] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *Adv. Neural Inform. Process. Syst.*, 2016. 4
- [33] SkyReels-AI. SkyReels v1: Human-centric video foundation model. <https://github.com/SkyworkAI/SkyReels-V1>, 2025. 2
- [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Int. Conf. Learn. Represent.*, 2021. 1, 2
- [35] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Adv. Neural Inform. Process. Syst.*, 2019.
- [36] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Int. Conf. Learn. Represent.*, 2021. 1
- [37] Endre Süli and David F Mayers. *An introduction to numerical analysis*. Cambridge University Press, 2003. 3
- [38] Hossein Talebi and Peyman Milanfar. Nima: Neural image assessment. *IEEE Trans. Image Process.*, 2018. 5
- [39] Wan Team. Wan: Open and advanced large-scale video generative models. 2025. 2, 5, 7, 8, 13, 14, 15
- [40] WANG Xi, Nicolas Dufour, Nefeli Andreou, CANI Marie-Paule, Victoria Fernandez Abrevaya, David Picard, and Vicky Kalogeiton. Analysis of classifier-free guidance weight schedulers. *Trans. Mach. Learn. Res.*, 2024. 2
- [41] Mengfei Xia, Nan Xue, Yujun Shen, Ran Yi, Tieliang Gong, and Yong-Jin Liu. Rectified diffusion guidance for conditional generation. *arXiv preprint arXiv:2410.18737*, 2024. 2
- [42] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. CogVideoX: Text-to-video diffusion models with an expert transformer. In *Int. Conf. Learn. Represent.*, 2025. 7
- [43] Candi Zheng and Yuan Lan. Characteristic guidance: Non-linear correction for diffusion model at large guidance scale. *arXiv preprint arXiv:2312.07586*, 2023. 2
- [44] Qinqing Zheng, Matt Le, Neta Shaul, Yaron Lipman, Aditya Grover, and Ricky TQ Chen. Guided flows for generative modeling and decision making. *arXiv preprint arXiv:2311.13443*, 2023. 1, 3
- [45] Le Zhuo, Ruoyi Du, Han Xiao, Yangguang Li, Dongyang Liu, Rongjie Huang, Wenze Liu, Lirui Zhao, Fu-Yun Wang, Zhanyu Ma, et al. Lumina-next: Making Lumina-T2X stronger and faster with Next-DiT. *arXiv preprint arXiv:2406.18583*, 2024. 1, 2, 5, 6, 8

Appendix

A1. Additional Experiments

A1.1. Experiments on Mixed Gaussian

Comparison of flow trajectory. We present the flow sampling trajectories with 10 steps of different methods in Fig. A1. As shown in the last column, samples guided by CFG move across the target distribution, while using only Cond leads to high variance in the sampled distribution. In contrast, CFG-Zero* effectively guides the samples toward the target distribution without excessive variance.

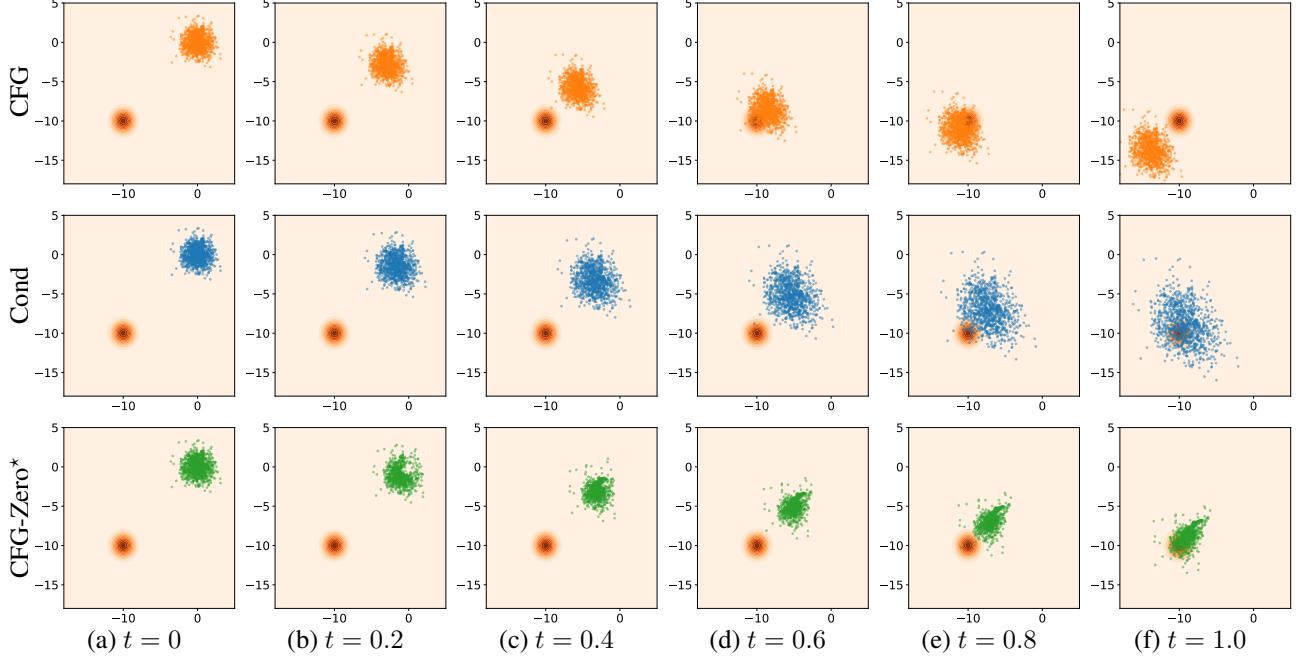


Figure A1. **Flow sampling trajectory.** Each panel shows the sample trajectories at a different time step.

Ablation study on the number of zero-Initialization steps is presented in Fig. A2. Specifically, we initialize the first 1, 2, 3, or 4 steps with zeros and observe that during the early training epochs, increasing the number of zero-initialized steps can be beneficial. However, as the number of zero-init steps increases, the learned model achieves better velocity compared to using zero initialization. At this stage, it becomes better to avoid zeroing out additional steps.

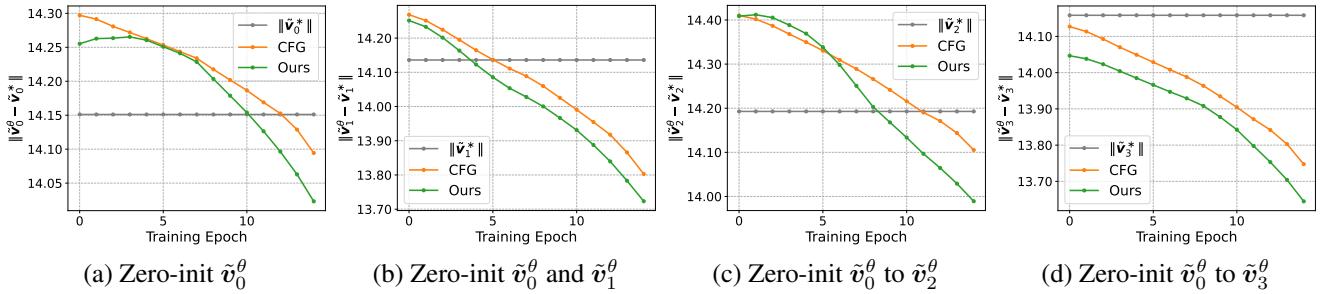


Figure A2. Ablation study of zero-init steps.

A1.2. Experiments on Text-to-Video Generation

A2. Additional Visual Results

In this section, we provide additional visual comparisons between our method and CFG. Fig. A3 presents videos generated by Wan2.1-1B, while Fig. A4 showcases those produced by the larger Wan2.1-14B model. Compared to CFG, our videos exhibit finer details, more vibrant colors, and smoother motion. Figs. A5 to A8 provide qualitative comparisons between CFG-Zero* and CFG. All results are shown without cherry-picking.

A3. Implementation Details

We first present our code, which can be easily integrated with any Flow-Matching-based model.

```
def optimized_scale(positive_flat, negative_flat):

    # Calculate dot production
    dot_product = torch.sum(positive_flat * negative_flat, dim=1, keepdim=True)

    # Squared norm of uncondition
    squared_norm = torch.sum(negative_flat ** 2, dim=1, keepdim=True) + 1e-8

    # st_star = v_cond^T * v_uncond / ||v_uncond||^2
    st_star = dot_product / squared_norm

    return st_star

# Get the velocity prediction
noise_pred_uncond, noise_pred_text = model(...)

positive = noise_pred_text.view(Batchsize,-1)
negative = noise_pred_uncond.view(Batchsize,-1)

# Calculate the optimized scale
st_star = optimized_scale(positive,negative)
# Reshape for broadcasting
st_star = st_star.view(Batchsize, 1, 1, 1)

# Perform CFG-Zero* sampling
if sample_step == 0:
    # Perform zero init
    noise_pred = noise_pred_uncond * 0.
else:
    # Perform optimized scale
    noise_pred = noise_pred_uncond * st_star + \
        guidance_scale * (noise_pred_text - noise_pred_uncond * st_star)
```

A3.1. Text-to-Image Details

Quantitative Evaluation. We evaluate Lumina-Next, Stable Diffusion 3, Stable Diffusion 3.5, and De-distill Flux on our self-curated text prompt benchmark, which consists of 200 short and long prompts covering a diverse range of objects, animals, and humans. Each model is assessed using its default optimal settings. For a fair comparison, we generate 10 images per prompt for each model.

Benchmark Results. We compare our method with CFG using Lumina-Next, SD3, and SD3.5 on T2I-CompBench [12], available at <https://github.com/Karine-Huang/T2I-CompBench/tree/main>. Each image is generated 10 times with different random seeds to ensure a fair comparison, and all models are evaluated using their optimal settings.

User Study. Our user study includes 76 participants, all familiar with text-to-image generation. Each participant an-

swers a questionnaire consisting of 25 questions, with each question randomly sampled from our generated images in T2I-CompBench [11, 12] to ensure fair comparisons.

A3.2. Text-to-Video Details

We evaluate Wan2.1 [39] both quantitatively and qualitatively, with all videos generated using the default settings specified in the official repository [39] (<https://github.com/Wan-Video/Wan2.1>). The VBench evaluation strictly follows the official guidelines (<https://github.com/Vchitect/VBench/tree/master>).



Text prompt: A bear on the right of a zebra, front view.



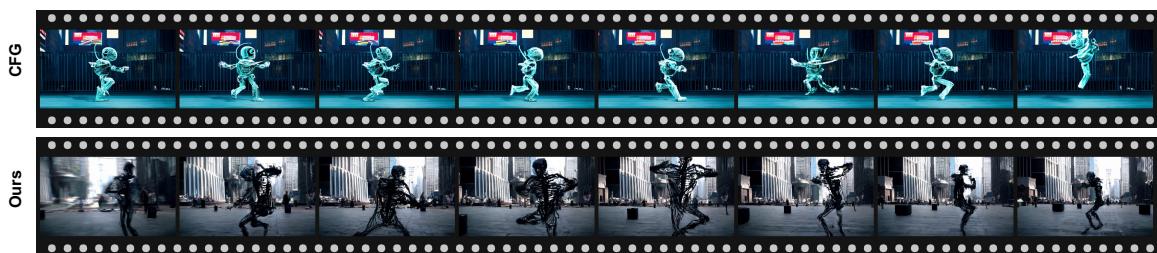
Text prompt: An epic tornado attacking above a glowing city at night, the tornado is made of smoke.



Text prompt: A horse running to join a herd of its kind.

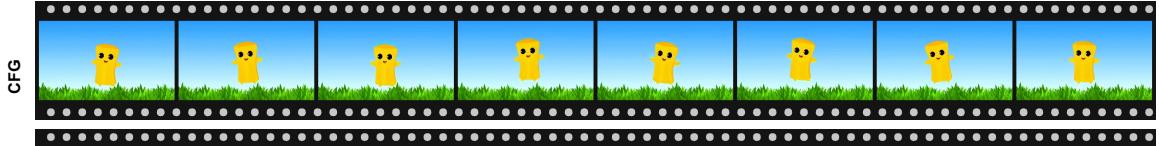


Text prompt: Origami dancers in white paper, 3D render.



Text prompt: Robot dancing in Times Square.

Figure A3. **Additional visual results.** Videos generated by Wan2.1-1B [39]



Text prompt: A cartoon character disco dances.



Text prompt: An extreme close-up of an gray-haired man with a beard in his 60s, he is deep in thought pondering the history of the universe as he sits at a cafe in Paris, ..., cinematic 35mm film.



Text prompt: A cat waking up its sleeping owner demanding breakfast. The owner tries to ignore the cat, but the cat tries new tactics and finally the owner pulls out a secret stash of treats from under the pillow to hold the cat off a little longer.



Text prompt: A old man, in a hat, walking towards the camera, talking to the camera, beautiful sunny day, at the beach, cinematic film shot in 35mm.



Text prompt: A happy dog wearing a green jacket and a pair of sun glasses walking towards the camera, talking, beautiful sunny day, at the beach, cinematic film shot in 35mm.

Figure A4. **Additional visual results.** Videos generated by Wan2.1-14B [39]

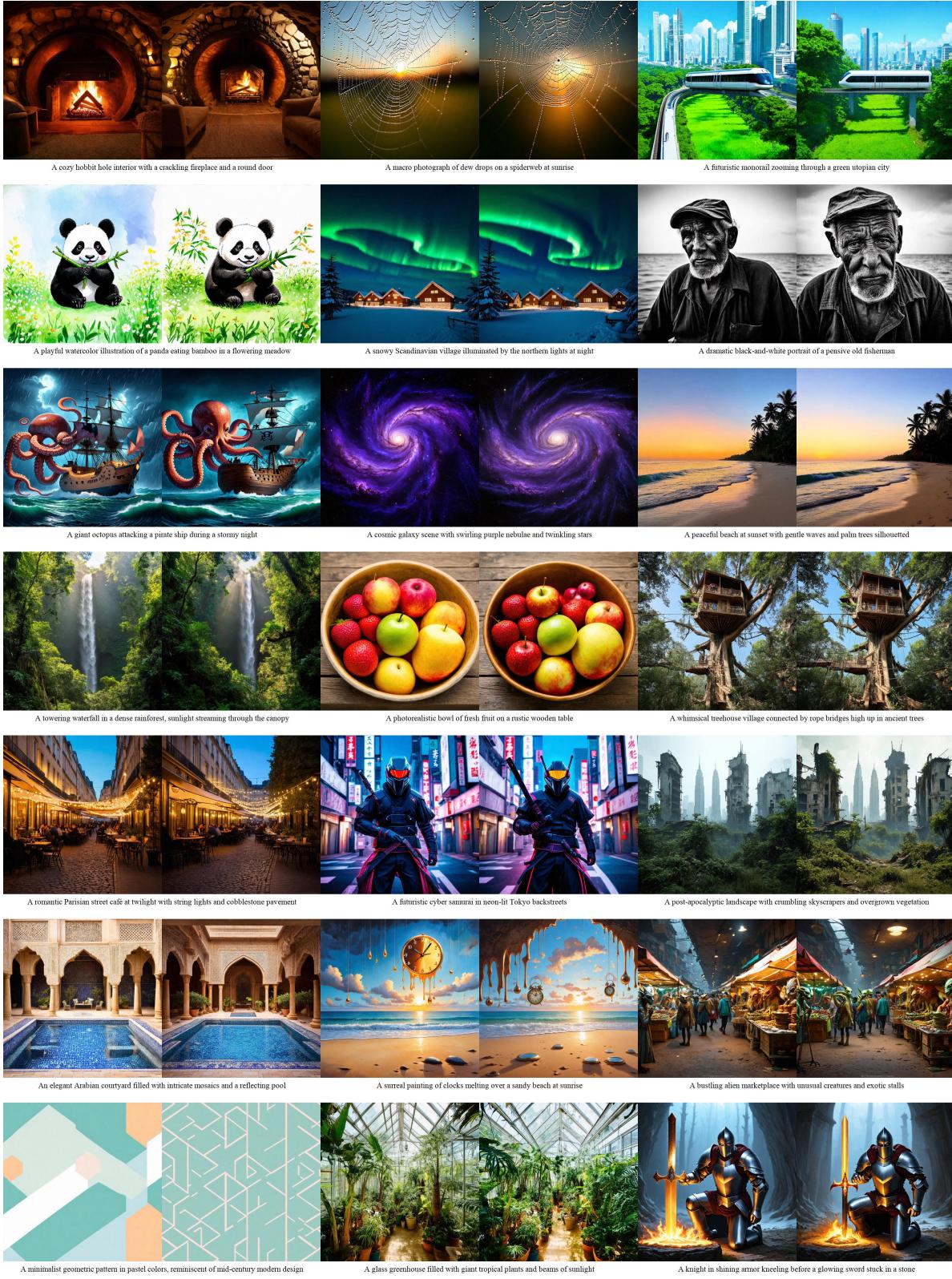


Figure A5. **Additional visual results.** Qualitative comparison between CFG (**left**) and CFG-Zero* (**right**). (Images generated by SD3.)

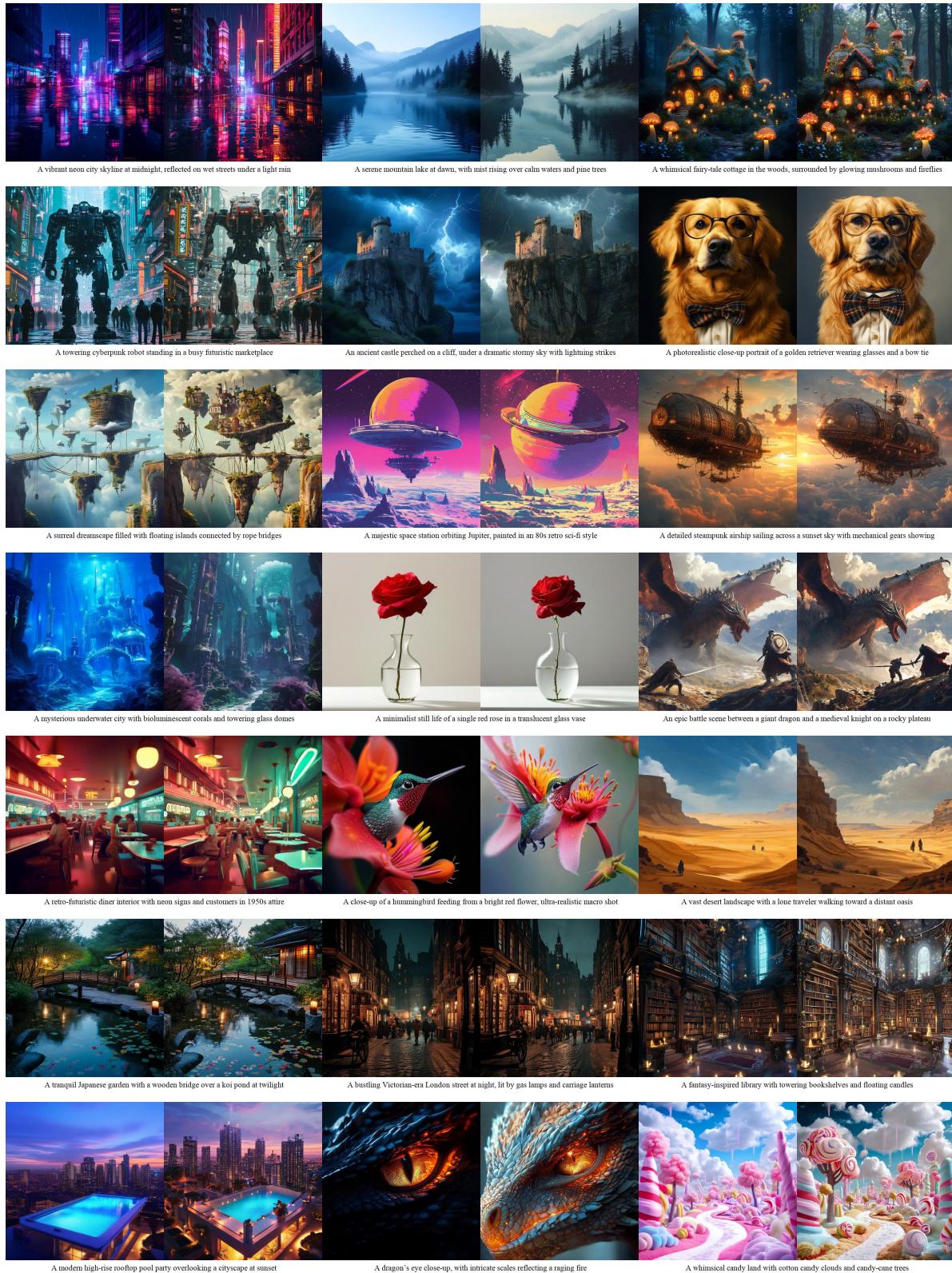


Figure A6. Additional visual results. Qualitative comparison between CFG (**left**) and CFG-Zero* (**right**). (Images generated by Lumina-Next.)

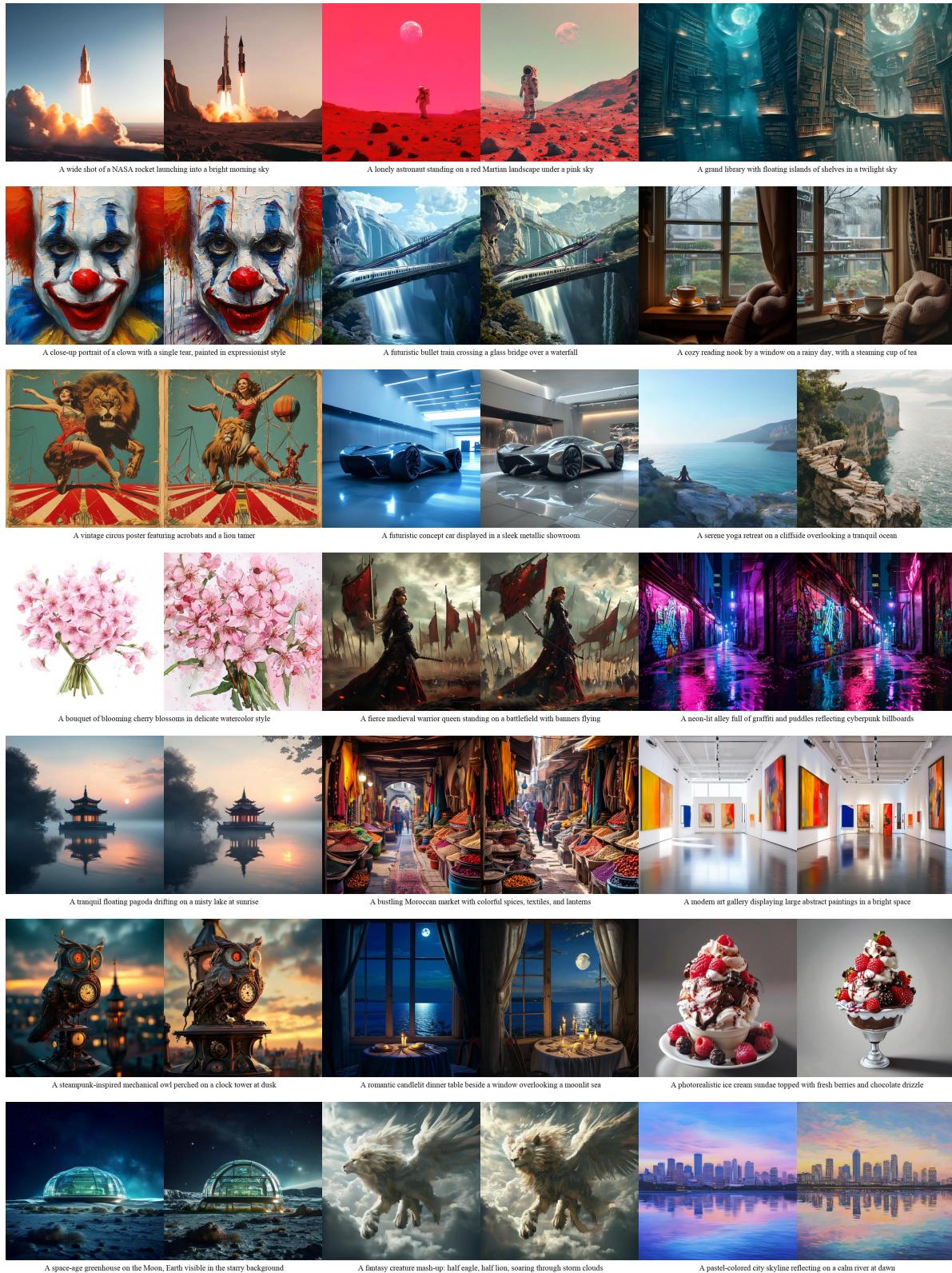


Figure A7. Additional visual results. Qualitative comparison between CFG (**left**) and CFG-Zero* (**right**). (Images generated by Lumina-Next.)

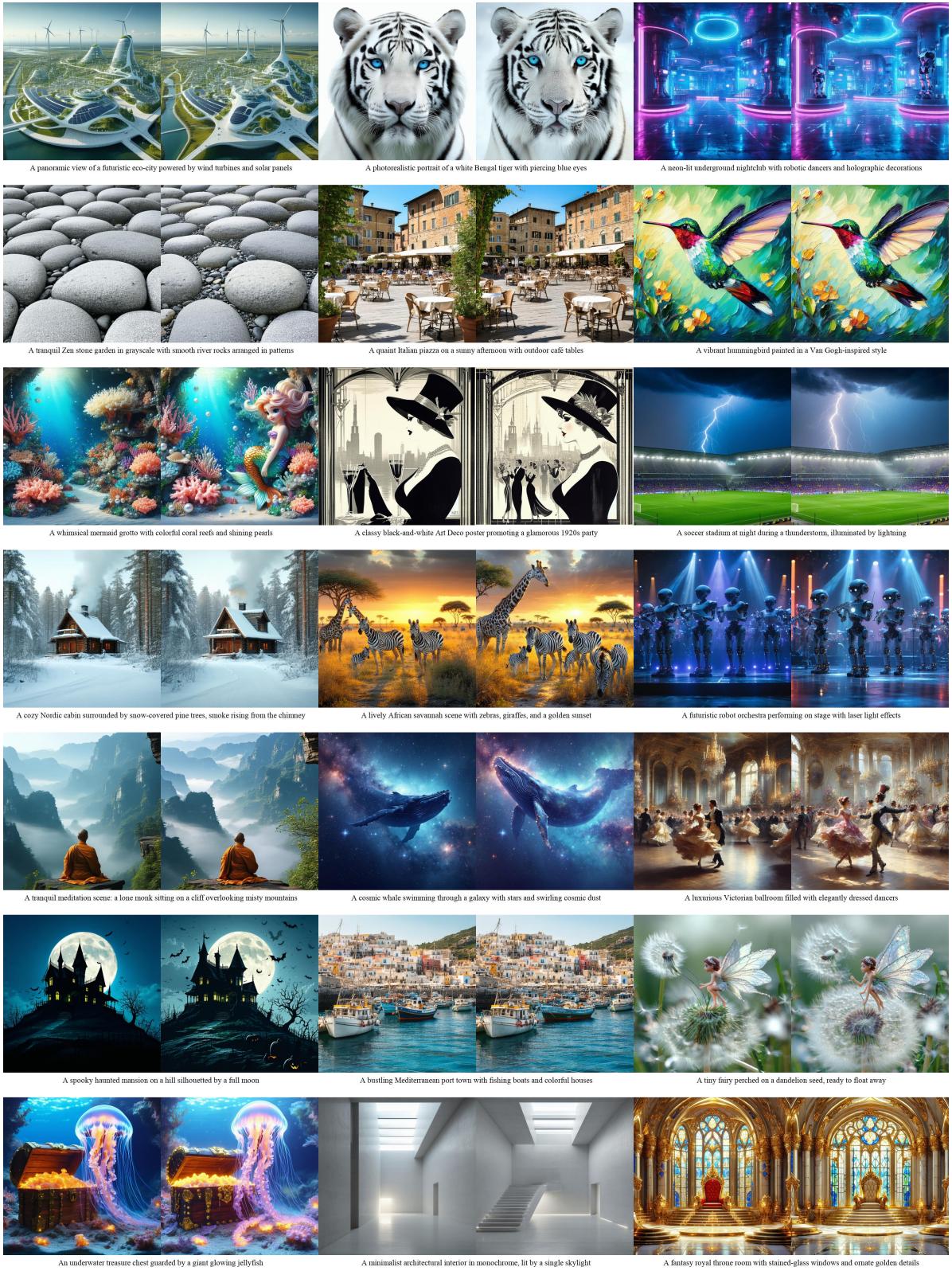


Figure A8. **Additional visual results.** Qualitative comparison between CFG (**left**) and CFG-Zero* (**right**). (Images generated by SD3.5.)