

Q-Sched: Pushing the Boundaries of Few-Step Diffusion Models with Quantization-Aware Scheduling

Natalia Frumkin
The University of Texas at Austin
nfrumkin@utexas.edu

Diana Marculescu
The University of Texas at Austin
dianam@utexas.edu

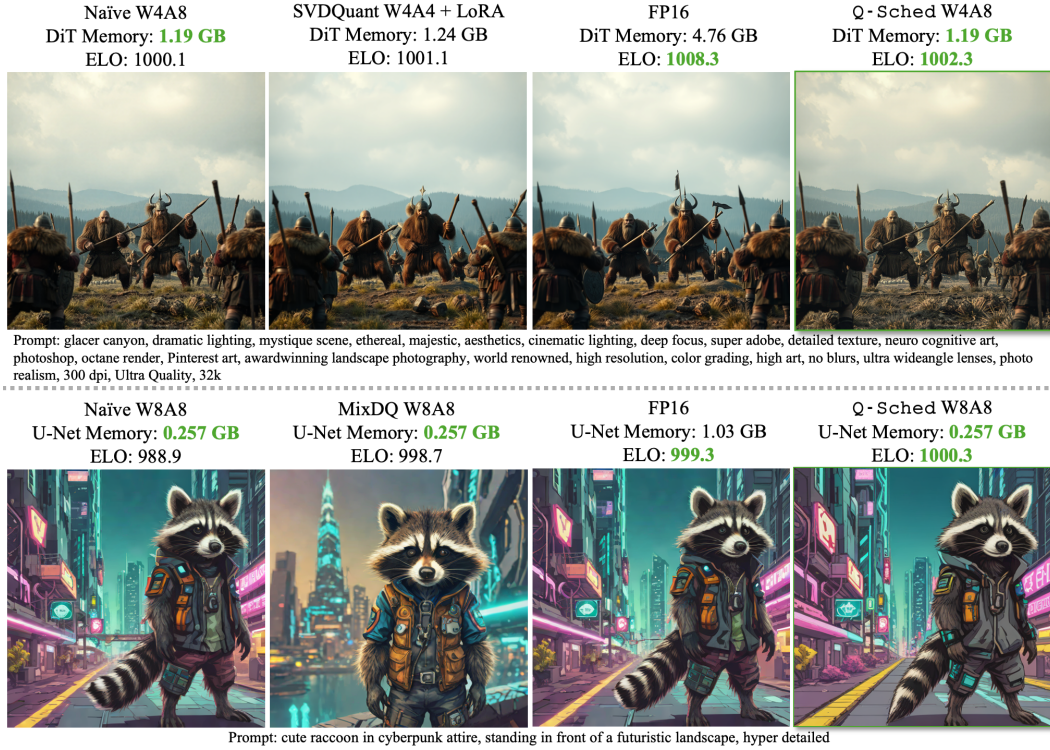
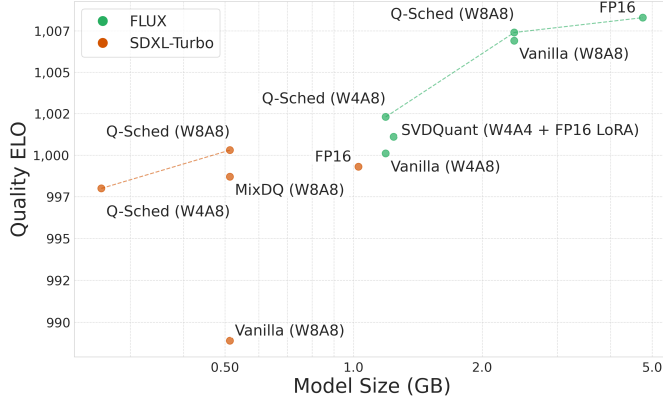


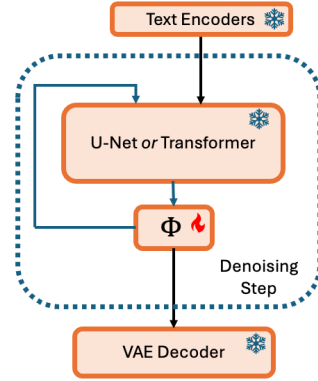
Figure 1: Q-Sched introduces a quantization-aware noise scheduler to few-step diffusion backbones and achieves excellent image fidelity. We find quantization and few-step diffusions to be complementary model compression strategies.

Abstract

Text-to-image diffusion models are computationally intensive, often requiring dozens of forward passes through large transformer backbones. For instance, Stable Diffusion XL generates high-quality images with 50 evaluations of a 2.6B-parameter model, an expensive process even for a single batch. Few-step diffusion models reduce this cost to 2-8 denoising steps but still depend on large, uncompressed U-Net or diffusion transformer backbones, which are often too costly for full-precision inference without datacenter GPUs. These requirements also limit existing post-training quantization methods that rely on full-precision calibration. We introduce Q-Sched, a new paradigm for post-training quantization that modifies the diffusion model scheduler rather than model weights. By adjusting the few-step sampling trajectory, Q-Sched achieves full-precision accuracy with a $4\times$ reduction in model size. To learn quantization-aware pre-conditioning coefficients,



(a) ELO Score vs. Model Size for various quantization methods on FLUX.1[schnell] [7] and SDXL-Turbo [38].



(b) Q-Sched optimizes the few-step diffusion scheduler, making it highly modular quantization method.

Figure 2: Q-Sched’s quantization-aware scheduling enables state-of-the-art image fidelity across multiple compressed few-step models. Q-Sched directly optimizes the few-step diffusion’s scheduler (see Figure 2b) whereas prior work directly optimizes the transformer or U-Net backbone.

we propose the JAQ loss, which combines text-image compatibility with an image quality metric for fine-grained optimization. JAQ is reference-free and requires only a handful of calibration prompts, avoiding full-precision inference during calibration.

Q-Sched delivers substantial gains: a 15.5% FID improvement over the FP16 4-step Latent Consistency Model and a 16.6% improvement over the FP16 8-step Phased Consistency Model, showing that quantization and few-step distillation are complementary for high-fidelity generation. A large-scale user study with more than 80,000 annotations further confirms Q-Sched’s effectiveness on both FLUX.1[schnell] and SDXL-Turbo. Code will be released upon publication. Our code is available at <https://github.com/enyac-group/q-sched>.

1 Introduction

Diffusion models are a powerful class of deep generative models with impressive results across computer vision [2, 6, 8, 17, 31, 50], natural language processing [4, 24], multi-modal modeling [5, 35], and interdisciplinary fields [3, 9]. State-of-the-art systems like Stable Diffusion XL [34, 31] and CogVideoX [51] demand server-grade GPUs and significant computational resources at inference time. This is due to the expensive denoising process, which often requires 40–1,000 steps per image or video, with each step invoking a deep, transformer-based noise estimation network.

To accelerate inference, two main strategies are typically considered: (1) reducing the number of function evaluations (*i.e.*, denoising steps), and (2) lowering the cost per evaluation. The latter can be addressed using standard model compression techniques such as quantization [13, 12], pruning [11], or knowledge distillation [18]. However, decreasing the number of denoising steps requires careful redesign of the sampling procedure, as naive reductions can severely degrade output quality.

Few-step diffusion models [43, 52, 30] reduce the number of steps taken along the sampling trajectory of a diffusion model’s ordinary or stochastic differential equation (ODE/SDE). These methods, along with a variety of other diffusion models, rely on an accurate probability flow ODE or variance-preserving SDE to express the relationship between the noise estimation network and the final generated image [42]. However, if the noise estimation network becomes corrupted, such as by compression, it can significantly change the ODE/SDE trajectory.

This work presents a novel quantization-aware noise schedule, dubbed Q-Sched, which modifies the few-step diffusion’s sampling trajectory to improve performance of few-step model.

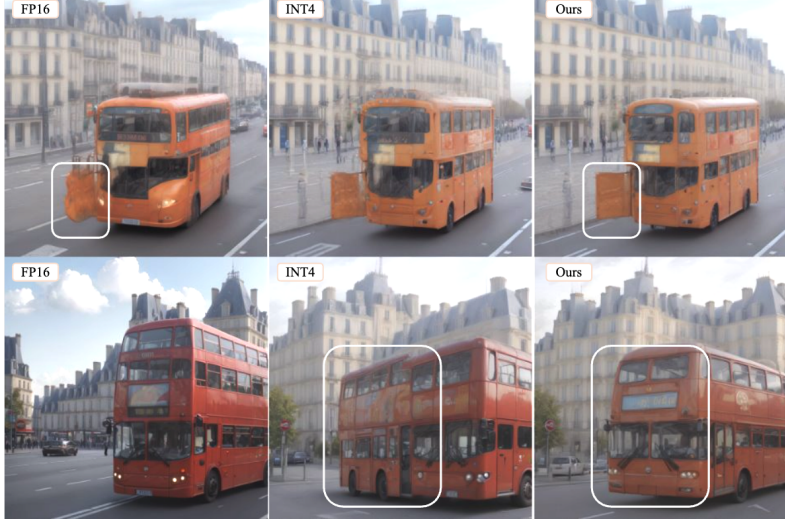


Figure 3: 4-Step (top row) and 8-Step (bottom row) LCMs. Prompt: "a car and a bus on a french highway". Q-Sched is capable of avoiding artifacts present in the FP16 or INT4 generative images. Q-Sched is close to the original schedule since it generates similar images yet our optimized schedule allows for Q-Sched to avoid some artifacts generated from the original schedule.

Our contributions are summarized as follows:

1. We present a novel quantization-aware scheduler that can be applied to existing few-step diffusion models. As shown in Figure 8, Q-Sched is capable of outperforming the original few-step diffusion, achieving an impressive **15.5% FID improvement** over an existing 4-step LCM *while* reducing model size.
2. Q-Sched’s novel preconditioning coefficients enable the quantized model to diverge from the potentially overfit few-step diffusion model and bypass artifacts created from the distillation process and quantization.
3. We propose the Joint Alignment-Quality (JAQ) loss, a reference-free metric that balances perceptual fidelity and text-image alignment. JAQ enables fine-grained control over visual attributes (*e.g.*, texture, detail, saturation) without requiring access to the full-precision model.
4. We perform more than 80,000 human preference annotations and find that Q-Sched outperforms MixDQ [53] on SDXL-Turbo and SVDQuant [23] on FLUX.1 in terms of perceived image quality.

In Figure 1, we illustrate that Q-Sched can achieve the best ELO rating when stacked up against other images in a pair-wise image quality comparison. Furthermore, we show in Figure 2a that Q-Sched is Pareto optimal with respect to both ELO and model size. These results highlight Q-Sched’s ability to balance perceptual quality and efficiency better than competing methods.

2 Background and Related Work

Diffusion models estimate an unknown data distribution $p_{data}(\mathbf{x})$ through a perturbation scheme known as denoising [16]. In this setup, we consider a set of sequential Gaussian perturbations which transform an image x_0 to random noise x_T . The forward process involves perturbing an image to noise, and the reverse process (denoising) involves removing noise from x_T until we obtain \hat{x}_0 . As shown by [42], the *continuous-time* denoising process can be modeled as the solution to the following stochastic differential equation:

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)dt + g(t)d\mathbf{w} \quad (1)$$

where $\mathbf{f}(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$ and $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ correspond to the drift and diffusion coefficients. \mathbf{w} is the Wiener process and $t \in [0, T]$ where T is a positive constant representing the number of

diffusion steps. In order to generate \hat{x}_0 from the denoising process, we sample from the probability flow ODE using standard numerical methods (Euler [42], Heun [20], *etc.*). The probability flow ODE is defined as:

$$d\mathbf{x}_t = \left[\mathbf{f}(x_t, t) - \frac{1}{2}g(t)^2 \nabla_x \log p_t(\mathbf{x}_t, \sigma_t) \right] dt \quad (2)$$

where $\nabla_x \log p_t(x_t, \sigma_t)$ is the Stein’s score. In the case of EDM [20] and Consistency Models [43], the drift and diffusion coefficients are set to $\mathbf{f}(x_t, t) = 0$ and $g(t) = \sigma_t = \sqrt{2t}$. It follows from these design choices that $p_t(x_t, \sqrt{2t}) = p_{data}(x) \otimes \mathcal{N}(0, t^2 I)$ where \otimes is the convolution operator.

The Stein’s score is learned through denoised score matching using a denoiser D_θ :

$$\nabla_x \log p_t(x_t, \sigma_t) = \frac{D_\theta(x_t, \sigma_t) - x_t}{\sigma_t^2} \quad (3)$$

which is parametrized by a neural network \mathcal{E}_θ . A common preconditioning scheme developed by EDM [20] is:

$$D_\theta(\mathbf{x}_t, \sigma_t) = c_{skip}(\sigma_t)\mathbf{x}_t + c_{out}(\sigma_t)\mathcal{E}_\theta(x_t, t) \left(c_{in}(\sigma_t)\mathbf{x}_t, c_{noise}(\sigma_t) \right) \quad (4)$$

where $c_{skip}(\cdot)$, $c_{out}(\cdot)$, $c_{in}(\cdot)$, and $c_{noise}(\cdot)$ are hyperparameters chosen prior to training. The DDIM [41] and Stable Diffusion [36] preconditioning is slightly different and will be described in more detail when discussing Consistency Models. In the following section, we will present a common method for compressing diffusion models.

Few-Step Diffusion Models Distillation is a leading approach for compressing diffusion models by reducing denoising steps. These *few-step* models offer substantial inference speedups with acceptable image quality, making them ideal for large-scale deployment. For instance, Instaflow [28] distills StableDiffusion v1-4 from $T = 50$ to $T = 1$, cutting inference time from 2.9s to 0.09s—enabling servers to handle orders of magnitude more users.

Consistency Models [43] use a consistency function to guide distillation—covered in the next section. Another popular approach, Rectified Flow, enables 1-step distillation by straightening the ODE trajectory, with Instaflow achieving this via coupling between noise and image spaces. Adversarial Diffusion Distillation (ADD) [38] further improves sample quality by introducing a discriminator that encourages the student model to generate outputs indistinguishable from the teacher’s.

Consistency Models In consistency models, the consistency function maps any point along the sampling trajectory to the final image x_0 :

$$x_0 = \mathcal{F}(x_t, t) \quad \forall t \in [0, T] \quad (5)$$

which is equivalent to the self-consistency property:

$$\mathcal{F}(x_t, t) = \mathcal{F}(x_s, s) \quad \forall t, s \in [0, T]. \quad (6)$$

Consistency Model variants arise from how one parametrizes the consistency function $\mathcal{F}(x_t, t)$. Borrowing from EDM, Latent Consistency Models (LCMs) [29] parameterize the consistency function using the same preconditioning as described in Equation (4). LCMs have grown to be a mainstream generative model due to their impressive performance with modest computational costs.

Following LCMs success, Trajectory Consistency Distillation (TCD) [54] developed a new trajectory consistency function to further generalize the self-consistency function. They parametrize the consistency function without using c_{skip} , c_{out} and instead opt for the DDIM-style parametrization:

$$\mathcal{F}_\theta(x_t, t) = \frac{\alpha_0}{\alpha_t} \mathbf{x}_t - \alpha_0 \left(\frac{\sigma_t}{\alpha_t} - \frac{\sigma_0}{\alpha_0} \right) \mathcal{E}_\theta(x_t, t). \quad (7)$$

The TCD schedule is used in Phased Consistency Models (PCMs) [45], which improve upon LCMs by adding an improved classifier-free guidance (CFG) solver, higher CFG controllability, along with better distillation stability. Our proposed Q-Sched scheduler can be applied to the preconditioning of both LCMs and PCMs to achieve highly distilled, quantized generative models.

Quantization for Diffusion Models There are a variety of methods applying model compression on standard text-to-image diffusion models. To compress the noise estimation network, \mathcal{E}_θ , PTQ4DM [39], ADP-DM [44] and Q-Diffusion [25] use timestep aware calibration methods to generate a calibration set for post-training quantization (PTQ). These methods consider how to learn the optimal quantization scheme for a given network across the many function evaluations required during the forward process. Prior work [46, 44, 10] also addresses activation quantization since model activations change with each timestep. Quantization schemes can be applied dynamically, such as using the TDQ module [40], to learn a set of quantization parameters for each timestep. These "timestep aware" methods do not consider how quantization incurs a distribution shift at each timestep.

To address distribution shifts, Q-DM [26] considers how the noise function accumulates errors across time and proposes a noise-estimation scheme to reduce the quantization error. Similarly, PTQD [13] addresses the error accumulation across timesteps by applying bias correction directly to the sampled image x_t . PTQD models the distribution shift from full precision to quantized model using a linear approximation:

$$\mathcal{E}_\theta^Q(x_t, t) = (1 + \gamma) \cdot \mathcal{E}_\theta + \delta \quad (8)$$

and aims to learn γ by standard deviation correction and considers δ as uncorrelated quantization noise to be modelled as Gaussian. PTQD provides the most practical scheduler modification for quantized diffusions which can be adapted for few-step diffusions. We explain how to adapt PTQD to the TCD scheduler in Appendix G. Next, we demonstrate how to apply Q-Sched to TCD and show how the same approach generalizes to other few-step samplers.

3 Quantization-Aware Scheduling

To prepare the TCD scheduler for optimization with Q-Sched, let us consider sampling with a quantized network. TCD's Strategic Stochastic Sampling (SSS) [54] using a quantized network $\mathcal{E}_\theta^Q(x_t, t)$ is given by:

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t \mathcal{E}_\theta^Q(x_t, t)}{\alpha_t} + \sigma_{s'} \mathcal{E}_\theta^Q(x_t, t) \right) + \eta \mathbf{z} \quad (9)$$

where a denoised image, x_s , is generated from a previous denoised image, x_t , at timestep t . The noise schedule is given by σ, α and the sampler injects stochastic noise sampled from a distribution $\mathbf{z} \sim N(0, I)$. The stochastic control parameter, η , is defined as:

$$\eta = \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \quad (10)$$

and can be manually overridden during sampling to generate images with different levels of stochasticity. The TCD sampler in Equation (9) is a state-of-the-art few-step diffusion sampler which is used in Phased Consistency Models. Note that the schedule relies on two inputs from the previous time-step: x_t and $\mathcal{E}_\theta^Q(x_t, t)$. These will be particularly relevant when applying Q-Sched.

Q-Sched: A Learnable Schedule Pre-Conditioner We introduce Q-Sched, a lightweight, post-training method that modifies the noise schedule of few-step diffusion models using two learnable scalar preconditioning coefficients: c_x and c_ϵ , applied to x_t and the quantized noise prediction \mathcal{E}_θ^Q , respectively. As illustrated in Figure 2b, Q-Sched operates independently of the model backbone (U-Net or transformer), making it broadly compatible with any few-step diffusion model whose scheduler resembles the TCD scheduler.

Quantized diffusion models often suffer from artifacts such as distortions and hallucinations (Figure 7). Traditional methods mitigate this by modeling distributional shifts from full-precision outputs using calibration data. In contrast, Q-Sched does not require access to the original model or any distribution statistics. Instead, it directly learns noise schedule corrections by optimizing a reference-free image quality metric, the JAQ loss.

The modified sampling procedure is:

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{c_x \mathbf{x}_t - \sigma_t c_\epsilon \mathcal{E}_\theta^Q(x_t, t)}{\alpha_t} + \sigma_{s'} c_\epsilon \mathcal{E}_\theta^Q(x_t, t) \right) + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z} \quad (11)$$

During distillation, few-step diffusion models compress the standard diffusion process into a couple of steps, removing the Gaussian denoising assumption and producing models that are sensitive to further model compression. We find that two coefficients are sufficient to learn a corrected sampling schedule which rivals its full precision counterparts for some few-step models. Our quantization-aware schedule is a new paradigm because it does not rely on any distribution or bias correction, but rather optimizes the preconditioning coefficients with respect to the JAQ loss, our proposed image quality metric. In using an image quality metric rather than a mean or standard deviation adjustment, our noise schedule has the ability to learn the best sampling trajectory rather than adhering to the full precision models trajectory, which for distilled few-step diffusions, is likely to be overfit.

To learn our preconditioning coefficients, c_x, c_ϵ , we apply grid search and evaluate each coefficient combination with respect to the JAQ loss. Next, we will discuss our new reference-free loss function, JAQ, and its benefits over existing image assessment tools.

JAQ: A Joint Alignment Quality Loss Function Reference-free metrics such as CLIPScore [14] have become essential for quick evaluation of text-to-image generation models. Unlike FID [15], SSIM [48], and other comparative metrics, reference-free metrics do not rely on a ground truth reference image and therefore are very useful in generative tasks when a ground truth is not available. When quantizing these generative models, an image generated from a quantized model, x_q , has an altered sampling trajectory from the original full precision model. This is evident in Figure 3, where x_q produces a different, sometimes cleaner image than full precision. In short, the quantized model takes an altered sampling trajectory which coarsely follows the full precision model yet generates sufficient differences that reference-based metrics do not capture the image’s detail.

Our Joint Alignment Quality loss combines a text-to-image compatibility score with a pure image quality score to achieve better results than simply optimizing with respect to metrics such as CLIPScore or CLIP-IQA [47] independently. We design the JAQ loss so that it can better differentiate between images that are highly similar to one another, whereas standard image quality metrics are designed to rank images that come from a much larger distribution. Given a text-to-image compatibility metric, $\text{TC}(x)$, and a pure image quality metric, $\text{IQ}(x)$, JAQ combines them as follows:

$$\text{JAQ}(x) = \text{TC}(x) + k \cdot \text{IQ}(x) \quad (12)$$

Only optimizing with respect to a text-to-image compatibility metric such as CLIPScore causes loss of image details as we only optimize for how closely the image follows the prompt (see Figure 6). Furthermore, we find that CLIPScore lacks adequate sensitivity to image artifacts due to quantization (examples of artifact types are shown in Figure 7).

Conversely, if we only use an image quality metric, our Q-Sched scheduler will optimize for generating details that may not otherwise be in the prompt. JAQ allows us to balance these two conflicting objectives using a linear combination with the parameter k controlling the tradeoff between text fidelity and image detail.

4 Experiments

Experimental Setup We apply Q-Sched across diverse few-step diffusion models, including both U-Net [37] and DiT [33] backbones, as well as consistency (LCM [29], PCM [45]) and flow-matching strategies. We further compare against recent state-of-the-art models such as SDXL-Turbo [38] and

Table 1: Comparison of different schedulers on Phased Consistency Models and Latent Consistency Models using a Stable Diffusion v1-5 backbone. The original schedule is TCD [54] for Phased Consistency Models and the Multi-step Consistency Sampling [29] for Latent Consistency Models. The FID and CLIPScore are calculated with respect to the COCO-30k dataset. NFEs stands for number of function evaluations referring to the number of passes through the network $\mathcal{E}_\theta^Q(x_t, t)$.

NFEs	Precision	Schedule	Calibration Size	PCMs		LCMs	
				FID	CLIPScore	FID	CLIPScore
2	FP16	Original	-	24.17	25.489	38.74	25.155
	W4A8	Original	-	28.70	25.343	40.93	24.886
	W4A8	PTQD	1024	23.33	25.265	37.59	24.919
	W4A8	Q-Sched	5	22.24	25.543	32.50	25.152
4	FP16	Original	-	23.29	25.482	31.94	25.969
	W4A8	Original	-	23.08	25.557	38.41	25.456
	W4A8	PTQD	1024	19.42	25.639	39.72	24.678
	W4A8	Q-Sched	5	17.39	25.715	26.98	25.336
8	FP16	Original	-	20.15	25.714	27.34	26.052
	W4A8	Original	-	18.48	25.664	27.55	25.397
	W4A8	PTQD	1024	15.85	25.770	28.06	25.241
	W4A8	Q-Sched	5	16.83	25.698	25.82	25.214

FLUX.1 [7]. We evaluate two regimes: 4W8A and 8W8A. Most models remain robust, with 4W8A outputs often matching full-precision quality. Due to metric limitations on high-fidelity images [19], we primarily report FID for 4W8A. Only the U-Net or DiT backbone is quantized, as it dominates model size (Table 4).

LCM and PCM are evaluated at 2, 4, and 8 denoising steps on COCO-30k [27], using FID (vs. real images), CLIPScore (for prompt-image alignment), and FID-SD (vs. Stable Diffusion outputs). FLUX.1 and SDXL-Turbo are assessed on the SVDQuant [23, 22] subset of MJHQ-30k, with FID and a human preference study for evaluation. MJHQ-30k is a collection of 5,000 high-quality Midjourney prompts from 10 common categories. We report FID for SDXL-Turbo and conduct a human preference study to fully evaluate these models.

We use two variants of the Joint Alignment Quality (JAQ) loss: one based on CLIP-derived metrics, the other on human preference scores. The CLIP-based version uses $TC(x) = CLIPScore(x)$ and $IQ(x) = CLIP - IQA(x)$. For SDXL-Turbo and FLUX.1, we apply a preference-based variant using $TC(x) = AQ-MAP(x)$ and $IQ(x) = HPSV2(x)$, where AQ-MAP [21] produces a spatial alignment score and HPSV2 [49] is fine-tuned on real human preferences. We set $k = 2$ in both cases.

Calibration Set We compare our results to PTQD [13], the only other quantization-aware noise scheduler designed for few-step diffusion. While PTQD relies on a 1,024-image calibration set generated from a full-precision model, our approach requires only a small set of representative prompts. Specifically, we hand-curate calibration prompts for the latent and phased consistency models, and sample from the sDCI prompt set [23] for SDXL-Turbo and FLUX.1 [schnell]. Our compact 20-image calibration set is reused across multiple evaluations, as larger sets would be prohibitively costly in post-training quantization.

Results: Latent and Phased Consistency Models In Figure 3, we consider three different quantized noise schedules and their performance on two consistency model families. We show that Q-Sched is able to overcome some artifacts present in the full precision and quantized models, illustrating that it learns a new few-step sampling trajectory which rivals the original model. Our method can produce greater detail than both FP16 and 4W8A (see Figure 9) for a variety of images. Furthermore, Q-Sched, provides excellent FID results in comparison to full precision. This illustrates that modifying the few-step diffusion’s noise schedule can allow the diffusion model to generate *better* images than the full precision model. Q-Sched outperforms PTQD noise schedules in 4/6 consistency model variants on StableDiffusion v1-5 with a fraction of the calibration set (see Table 1).

PTQD performs correction to the quantized model’s distribution *with respect to its full precision counterpart*. In contrast, Q-Sched does not rely on images from a full precision model and can be adjusted to outperform a flagship few-step diffusion model with only five prompts. As consistency models are distilled from a large diffusion model, we’d expect them to be very sensitive to quantization as often distilled models are much more sensitive to compression. Surprisingly, the opposite is true. Q-Sched can outperform a full precision few-step diffusion model by **16.1%, 15.5%, 5.6%** for a **2-step, 4-step, and 8-step model** respectively, illustrating that quantization and few-step distillation are complementary model compression techniques.

Scheduler	Precision	FID	FID-SD	CLIPScore
TCD	FP16	18.65	10.45	26.531
TCD	W4A8	22.70	12.51	26.241
PTQD	W4A8	161.96	176.29	25.910
Q-Sched	W4A8	<u>18.89</u>	<u>12.17</u>	<u>26.513</u>

(a) Comparison on a 2-step Phased Consistency model using the Stable Diffusion XL backbone. FID-SD is computed relative to images generated by Stable Diffusion XL using corresponding COCO-30k prompts.

Method	Precision	FID
-	FP16	25.48
Naive	W4A8	25.75
MixDQ	W4A8	25.36
Q-Sched	W4A8	<u>21.41</u>
Naive	W8A8	25.49
MixDQ	W8A8	25.16
Q-Sched	W8A8	26.34

(b) Quantized model comparison on SDXL-Turbo under varying bitwidths. FID is computed on the MJHQ dataset.

Table 2: Quantitative evaluation of large-Scale few-step diffusion models with a Stable Diffusion XL backbone. W4A8 and W8A8 are a $4\times$ and $8\times$ model size reduction in comparison to FP16, yet our method improves over baseline. As FID, FID-SD, and CLIPScore may exhibit reduced reliability at large model scales, we complement these metrics with user preference studies in Figure 2a.

In Table 2a, we present results on a large-scale 2-step Phased Consistency Model built on the Stable Diffusion XL backbone. Q-Sched achieves only a 1.2% FID degradation in the 4W8A setting, demonstrating that quantization-aware preconditioning can maintain high performance even under aggressive compression. In contrast, PTQD fails in this regime, yielding significantly higher FID and visibly poor image quality. We attribute this to its reliance on Gaussian noise assumptions, which break down in few-step diffusion—especially for large models—where each step approximates a segment of the ODE trajectory rather than a Gaussian denoising step. This is expected, given that the true data distribution p_{data} is highly complex and often non-Gaussian.

Results: SDXL-Turbo and FLUX.1[schnell] In Table 2b, we compare quantization strategies on SDXL-Turbo (4-step inference) using the FID metric on the MJHQ dataset, evaluating two bitwidth settings: W4A8 and W8A8. Under W4A8, Q-Sched achieves a FID of 21.41, significantly outperforming MixDQ [53] (25.36) and Naive (25.75), demonstrating strong robustness to aggressive quantization. However, at W8A8, Q-Sched shows a higher FID (26.34) than both MixDQ (25.16) and Naive (25.49), suggesting that its advantages are most pronounced in lower-bit regimes, where other methods degrade more severely.

Despite this, a small-scale user study with 52 anonymized participants showed that Q-Sched was preferred by 56% of users over MixDQ in the W8A8 setting, indicating that perceptual quality may remain competitive even when FID is less favorable. Additionally, in Figure 2a, we present user preference results for Q-Sched applied to both SDXL-Turbo and FLUX.1 [schnell], showing that it outperforms MixDQ [53] and SVDQuant [23], respectively, at similar model sizes (see Appendix A for details). We compute an ELO rating, a relative quality ranking inspired by chess scoring, by aggregating all pairwise 1v1 image comparisons across models, where a higher score reflects consistent user preference.

In Figure 4, we compare Q-Sched across bit widths using a user study. W4A4 proved too aggressive, but W4A5 and W4A6

Q-Sched w4a5	40.5%	59.5%	FP16
Q-Sched w4a6	43.4%	56.6%	FP16
Q-Sched w4a8	42.7%	57.3%	FP16
Q-Sched w8a8	54.3%	45.7%	FP16

User Preference (%)

Figure 4: Comparing Q-Sched across various bit widths.

produced images comparable to full precision. 1v1 comparisons with full-precision FLUX.1 [7] follow the protocol in Appendix A.

Ablation on Pre-Conditioning Coefficients and Loss Function Choice We ablate the choice of pre-conditioning coefficients in the Phased Consistency Model by comparing performance when optimizing only the model-side coefficient c_ϵ , the sample-side coefficient c_x , or both jointly. As shown in Figure 5b, jointly optimizing both c_ϵ and c_x consistently yields the best results across all three metrics: PickScore, HPSv2, and JAQ Loss. These findings highlight the importance of treating both denoising and reconstruction terms as tunable components rather than fixing one a priori. All metrics are averaged over 1024 images generated with the SDXL backbone.

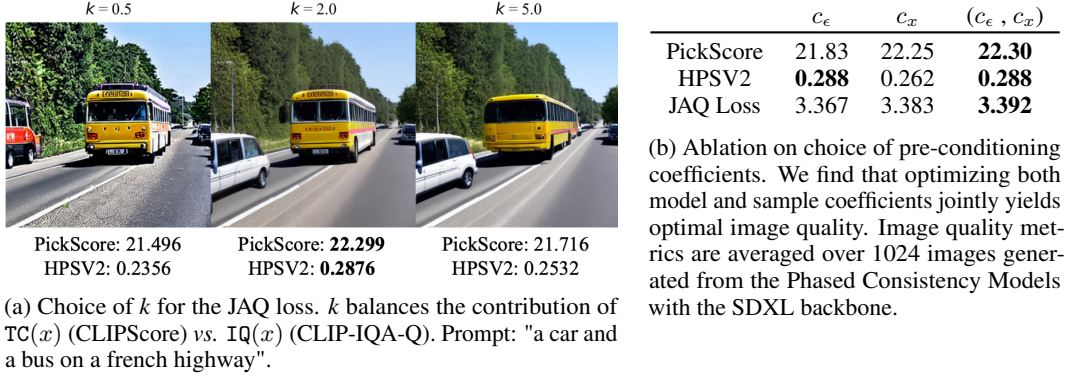


Figure 5: Ablation studies on various design choices for Q-Sched.

How Do We Choose k For The JAQ Loss? We optimize the Q-Sched preconditioners using the JAQ loss, which balances image quality and text-image consistency via a tradeoff hyperparameter, k . As shown in Figure 5a, small k values can lead to color distortion, while larger values (*e.g.*, $k = 5$) cause outputs to drift from the true data distribution. In such cases, the JAQ loss behaves similarly to CLIP-IQA-Q, which lacks sensitivity to concept alignment. We find that a hand-tuned value of k is sufficient for producing a high-quality noise schedule, and the final results are not highly sensitive to its exact choice. Throughout our experiments, we use $k = 2$.

For a detailed comparison of loss functions, see Appendix C.1.

5 Broader Impacts

Model compression enables wider and more ubiquitous AI usage as we compress large foundation models to run on resource-constrained GPUs. Our work’s potential societal consequences are similar to those of prior work as both quantization and few-step diffusions are model compression methods for text-to-image generative models. Generated images have the potential to mislead, mis-represent and cause social harm. We conduct a user preference on a crowd-sourcing platform where generated content is shown to users worldwide and has the potential to cause harm.

6 Conclusion

Few-step diffusion models dramatically reduce inference cost by distilling large generative models—such as Stable Diffusion XL—into versions requiring only 2–8 denoising steps, achieving a 5–25× speedup. However, these models typically reduce runtime without addressing model size. Our method, Q-Sched, pushes this efficiency frontier further by introducing quantization into the few-step regime. Through noise-aware preconditioning coefficients, Q-Sched enables effective quantization with minimal performance loss. We report **8.0% and 16.1% FID improvements** over full-precision baselines for PCMs and LCMs, respectively. A user preference study also shows that Q-Sched **outperforms existing quantization methods on FLUX.1[schnell] and SDXL-Turbo in perceived image quality**. These results demonstrate that quantization and few-step distillation are complementary, enabling substantial efficiency gains without compromising generation quality.

7 Acknowledgments

This work was supported in part by NSF CCF Grant No. 2107085, iMAGiNE - the Intelligent Machine Engineering Consortium at UT Austin, and UT Cockrell School of Engineering Doctoral Fellowships. Human Evaluation studies were conducted using Rapidata.

References

- [1] Rapidata: An api that provides fast access to large-scale human evaluations, 2025. URL <https://www.rapidata.ai/>. Accessed: 2025-05-16.
- [2] Tomer Amit, Eliya Nachmani, Tal Shaharabany, and Lior Wolf. Segdiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*, 2021.
- [3] Namrata Anand and Tudor Achim. Protein structure and sequence generation with equivariant denoising diffusion probabilistic models. *arXiv preprint arXiv:2205.15019*, 2022.
- [4] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- [5] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022.
- [6] Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khrulkov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. *arXiv preprint arXiv:2112.03126*, 2021.
- [7] Black Forest Labs. Flux.1-schnell. <https://huggingface.co/black-forest-labs/FLUX.1-schnell>, 2024. Accessed: 2025-05-14.
- [8] Emmanuel Asiedu Brempong, Simon Kornblith, Ting Chen, Niki Parmar, Matthias Minderer, and Mohammad Norouzi. Denoising pretraining for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4175–4186, 2022.
- [9] Chentao Cao, Zhuo-Xu Cui, Shaonan Liu, Dong Liang, and Yanjie Zhu. High-frequency space diffusion models for accelerated mri. *arXiv preprint arXiv:2208.05481*, 2022.
- [10] Lei Chen, Yuan Meng, Chen Tang, Xinzhu Ma, Jingyan Jiang, Xin Wang, Zhi Wang, and Wenwu Zhu. Q-dit: Accurate post-training quantization for diffusion transformers. *arXiv preprint arXiv:2406.17343*, 2024.
- [11] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. *Advances in neural information processing systems*, 36, 2024.
- [12] Cong Guo, Yuxian Qiu, Jingwen Leng, Xiaotian Gao, Chen Zhang, Yunxin Liu, Fan Yang, Yuhao Zhu, and Minyi Guo. Squant: On-the-fly data-free quantization via diagonal hessian approximation. *arXiv preprint arXiv:2202.07471*, 2022.
- [13] Yefei He, Luping Liu, Jing Liu, Weijia Wu, Hong Zhou, and Bohan Zhuang. Ptdq: Accurate post-training quantization for diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [14] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- [15] 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. *Advances in neural information processing systems*, 30, 2017.
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [17] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23(47):1–33, 2022.

- [18] Tao Huang, Yuan Zhang, Mingkai Zheng, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge diffusion for distillation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [19] Sadeep Jayasumana, Srikumar Ramalingam, Andreas Veit, Daniel Glasner, Ayan Chakrabarti, and Sanjiv Kumar. Rethinking fid: Towards a better evaluation metric for image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9307–9315, 2024.
- [20] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *arXiv preprint arXiv:2206.00364*, 2022.
- [21] Chunyi Li, Haoning Wu, Zicheng Zhang, Hongkun Hao, Kaiwei Zhang, Lei Bai, Xiaohong Liu, Xiongkuo Min, Weisi Lin, and Guangtao Zhai. Q-refine: A perceptual quality refiner for ai-generated image. In *2024 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2024.
- [22] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation, 2024.
- [23] Muiyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Junxian Guo, Xiuyu Li, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdquant: Absorbing outliers by low-rank component for 4-bit diffusion models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [24] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems*, 35:4328–4343, 2022.
- [25] Xiuyu Li, Long Lian, Yijiang Liu, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. *arXiv preprint arXiv:2302.04304*, 2023.
- [26] Yanjing Li, Sheng Xu, Xianbin Cao, Xiao Sun, and Baochang Zhang. Q-dm: An efficient low-bit quantized diffusion model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [27] 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 *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- [28] Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. Instaflo: One step is enough for high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2023.
- [29] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- [30] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- [31] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [32] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing*, 21(12):4695–4708, 2012.
- [33] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.

- [34] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- [35] 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.
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- [38] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, pages 87–103. Springer, 2024.
- [39] Yuzhang Shang, Zhihang Yuan, Bin Xie, Bingzhe Wu, and Yan Yan. Post-training quantization on diffusion models. *arXiv preprint arXiv:2211.15736*, 2022.
- [40] Junhyuk So, Jungwon Lee, Daehyun Ahn, Hyungjun Kim, and Eunhyeok Park. Temporal dynamic quantization for diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [41] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [42] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations, 2021.
- [43] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- [44] Changyuan Wang, Ziwei Wang, Xiuwei Xu, Yansong Tang, Jie Zhou, and Jiwen Lu. Towards accurate data-free quantization for diffusion models. *arXiv preprint arXiv:2305.18723*, 2(5), 2023.
- [45] Fu-Yun Wang, Zhaoyang Huang, Alexander William Bergman, Dazhong Shen, Peng Gao, Michael Lingelbach, Keqiang Sun, Weikang Bian, Guanglu Song, Yu Liu, et al. Phased consistency model. *arXiv preprint arXiv:2405.18407*, 2024.
- [46] Haoxuan Wang, Yuzhang Shang, Zhihang Yuan, Junyi Wu, and Yan Yan. Quest: Low-bit diffusion model quantization via efficient selective finetuning. *arXiv preprint arXiv:2402.03666*, 2024.
- [47] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 2555–2563, 2023.
- [48] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4): 600–612, 2004.
- [49] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *CoRR*, 2023.
- [50] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for video generation. *arXiv preprint arXiv:2203.09481*, 2022.
- [51] 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. *arXiv preprint arXiv:2408.06072*, 2024.

- [52] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. *arXiv preprint arXiv:2311.18828*, 2023.
- [53] Tianchen Zhao, Xuefei Ning, Tongcheng Fang, Enshu Liu, Guyue Huang, Zinan Lin, Shengen Yan, Guohao Dai, and Yu Wang. Mixdq: Memory-efficient few-step text-to-image diffusion models with metric-decoupled mixed precision quantization. In *European Conference on Computer Vision*, pages 285–302. Springer, 2024.
- [54] Jianbin Zheng, Minghui Hu, Zhongyi Fan, Chaoyue Wang, Changxing Ding, Dacheng Tao, and Tat-Jen Cham. Trajectory consistency distillation. *arXiv preprint arXiv:2402.19159*, 2024.

A Details on User Preference Assessment

We design our evaluation setup following the user preference study methodology from SDXL-Turbo [38], with several improvements. For each model pair in this study, we perform 1-vs-1 comparisons based on shared prompts. Human responses, collected via Rapidata [1], come from evaluators who are presented with two images, each generated by a different model for the same prompt, and are asked: “Which image is of higher quality and more aesthetically pleasing?”

Evaluators are globally sourced and must pass a set of validation questions designed to assess annotation quality. Only those who successfully complete this qualification step are allowed to rate the models.

ELO scores are computed using the same approach as SDXL-Turbo [38], with $K = 32$. We find that this value of K enables more noticeable ranking adjustments, especially when models have similar performance levels.

All models in our study are evaluated using 1,000 prompts sampled from the MJHQ-30k dataset. We release this subset, which we call the Q-Sched split, to enable consistent benchmarking of future quantization methods. Each prompt is evaluated by four unique annotators. Therefore, each 1-vs-1 comparison results in 4,000 total human annotations.

B Compute Resources & Statistical Significance

We conduct all our experiments on a high-end AI server with eight Nvidia A6000s. Each model can be run independently on one A6000 and Q-Sched takes approximately twenty minutes to run the full grid search.

Our main experiments are averaged over two-three runs but we do not report error bars at this time.

C Ablation Studies

C.1 Comparing Loss Functions for Q-Sched

To evaluate the overall image quality for text-image generative modeling, CLIPScore [14] is specifically designed to capture text-image compatibility and does not consider overall image quality. In Figure 6, we illustrate that Q-Sched optimized with CLIPScore produces an updated noise schedule that is over saturated and lacks image depth. In contrast, Brisque [32] is often used as a standard reference-free image quality metric, but when used in Q-Sched it creates images with smoother and less detailed features. We consider three variants of CLIP-IQA [47] and find that CLIP-IQA using the predefined quality prompt (we denote this version by CLIP-IQA-Q) achieves a noise schedule with high-fidelity images. However, CLIP-IQA-Q has a significant weakness: it cannot properly score images with hallucinations because it does not have an understanding of the underlying image prompt or concept. Therefore, we combine the benefits of CLIPScore and CLIP-IQA-Q into the JAQ loss and find that the resulting schedule fares extremely well with respect to raw image quality as well as to concept adherence.

C.2 Adding Stochasticity

Phased Consistency Model’s implementation of the original sampler, TCD, is deterministic, meaning that there is no additive noise during sampling. The controllable noise parameter, η , allows a practitioner to adjust the additive noise during the sampling process and is defined in Equation (9). In order to compare PTQD’s correction to our method, we ablate across different levels of stochasticity and report performance for six stochasticity levels in Table 3. $\eta = 0$ refers to deterministic sampling and PTQD’s uncorrelated noise correction is not used since it adds stochastic noise by construction. Please see the appendix for more details on PTQD’s implementation in both deterministic and stochastic sampling regimes.

We find that Q-Sched outperforms PTQD for all stochasticity regimes on the 2-step phased consistency model. With a simple grid search using our JAQ loss, we can outperform PTQD and the original TCD scheduler in different sampling regimes.

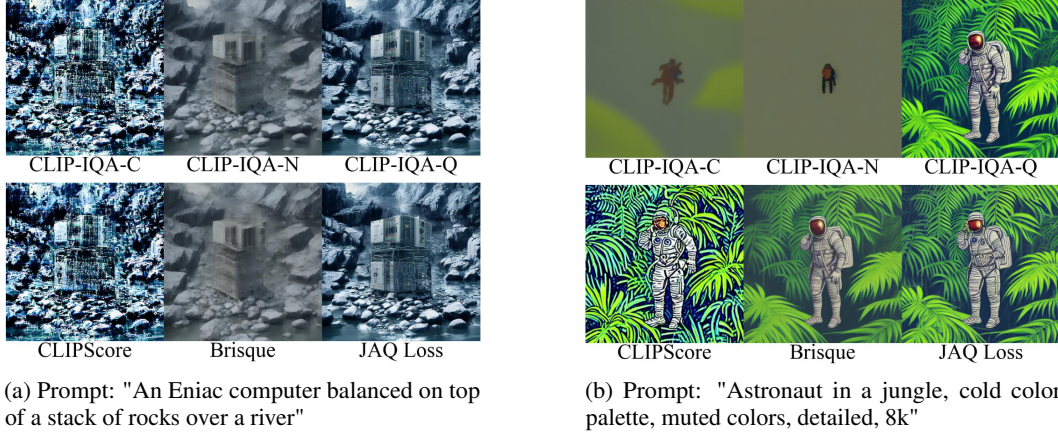


Figure 6: Optimizing Q-Sched with various reference-less image quality metrics. Our loss function, JAQ, is a linear combination of CLIPScore and CLIP-IQA-Q. We compare against three CLIP-IQA prompts: Complexity, Noisiness, and Quality denoted as -C, -N, -Q respectively.

Table 3: Adding Stochasticity and its effect on W4A8 Quantization for PCM using a Stable Diffusion v1-5 backbone. We report FID on COCO-30k. The stochasticity term, η , controls the amount of added Gaussian noise. $\eta = 0$ is deterministic sampling.

Method	$\eta =$					
	0	0.1	0.3	0.5	0.7	0.9
TCD	28.70	24.06	23.44	22.97	26.74	22.40
PTQD	23.33	25.59	24.95	25.69	24.53	26.71
Q-Sched	22.24	19.29	23.44	19.67	19.46	17.87

D Quantization-Induced Artifacts

In our preliminary analysis using a two-step Consistency Model, we observed several characteristic ways in which quantization degrades image quality. As shown in Figure 7, quantized models tend to exhibit three prominent types of artifacts: color distortion, image degradation, and hallucinated structures. These issues are especially pronounced in low-bit settings and appear consistently across a variety of models and prompts.

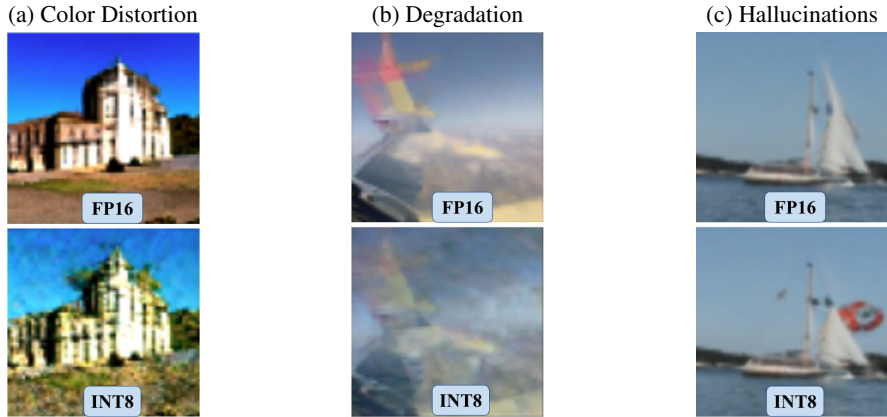


Figure 7: Three types of image artifacts that occur when quantizing image generation models. Images are unconditionally generated from a Two-Step Consistency Model [43].

E Model Size Analysis

In Table 4, we show the full model size breakdown for the diffusion model backbone, text encoders, and the VAE decoder. During inference, either one or both text encoders are used, and we do not need the VAE encoder, since this is for training exclusively.

Table 4: FP16 Diffusion Model Size Breakdown (in GB)

	LCM	PCM	SDXL-Turbo	FLUX.1[schnell]
UNet/DiT	1.72	4.84	1.03	4.76
Text Encoder(s)	0.25	0.29	0.33	1.95
VAE Decoder	0.07	0.13	0.02	0.02
Total	2.04 GB	5.26 GB	1.37 GB	6.73 GB

For our ELO *vs.* Model Size Pareto front in Figure 2a, we consider the DiT memory and compute model size by taking the parameter count and multiplying it by the number of bytes required per parameter. For W4A4 + LoRA 64, the setup used for SVDQuant [23], we compute the number of LoRA parameters using the back-of-the-envelope calculation provided in SVDQuant and add it to this calculation. We provide raw data for clarity in Table 5.

Table 5: DiT Memory (in GB) for various bitwidths.

Precision	SDXL-Turbo	FLUX.1[schnell]
FP16	1.03	4.76
W8A8	0.51	2.38
W4A4 + LoRA 64	0.28	1.24
W4A8	0.26	1.19

F Additional Analysis on COCO-30k

This result reinforces the core finding of our paper: quantization, when paired with a scheduler designed to account for noise sensitivity (as in Q-Sched), can be synergistic with few-step diffusion rather than detrimental. Notably, our quantized model achieves a lower FID than the original full-precision model, suggesting that Q-Sched helps overcome limitations introduced by both step reduction and bit-level compression.

These findings complement the results on SDXL-Turbo and FLUX.1[schnell] discussed in the main paper, and further establish Q-Sched as a general-purpose solution for high-fidelity, compressed diffusion generation.

G Applying PTQD to the TCD Scheduler

Using PTQD’s linear parameterization for the quantization error, we substitute $\mathcal{E}_\theta^Q(x_t, t) = (1 + \gamma) \cdot \mathcal{E}_\theta + \delta$ into Equation (9):

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t \mathcal{E}_\theta(x_t, t)}{\alpha_t} + \sigma_{s'} \mathcal{E}_\theta(x_t, t) \right) + \frac{\alpha_s}{\alpha_{s'}(1 + \gamma)} \left(\sigma_{s'} - \frac{\alpha_{s'} \sigma_t}{\alpha_t} \right) \delta + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z}. \quad (13)$$

PTQD assumes the uncorrelated noise is sampled from a normal distribution $\delta \sim N(\mu_\delta, \sigma_\delta)$. This method applies bias correction to handle the mean deviation, μ_δ , and analytically compute standard deviation, σ_δ . We adapt PTQD’s approach to the TCD schedule and use the new standard deviation, σ_δ for sampling δ :

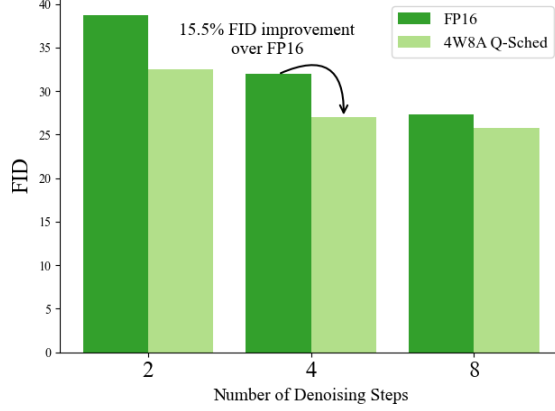


Figure 8: FID on COCO-30k. A 4W8A compressed model with our Q-Sched scheduler outperforms its FP16 counterpart with a $4\times$ reduction in model size.

$$\sigma_\delta^2 = 1 - \frac{\alpha_s^2}{\alpha_{s'}^2} \left(1 - \left(\frac{\delta(\sigma_{s'} - \frac{\alpha_{s'}\sigma_t}{\alpha_t})}{(1+\gamma)} \right)^2 \right). \quad (14)$$

For the edge case, where $\sigma_\delta < 0$, the deviation is set to zero ($\sigma_\delta = 0$). The proof for extending PTQD to the TCD scheduler is in the appendix.

PTQD attempts to model the distribution shift from a full precision to quantized model using two assumptions:

1. The quantized model’s distribution shift can be modeled through a linear correction term.
2. The uncorrelated quantization noise is normally distributed.

While these assumptions are similar to prior work on diffusion models, they are likely to break down on the few-step diffusions where the denoising process is distilled from many steps and is not expected to be linear nor follow a Gaussian distribution.

Quantization Noise Correction using PTQD Based on the PTQD quantization noise assumption, the quantization error is linearly parametrized as $\Delta\mathcal{E}_\theta = \gamma \cdot \mathcal{E}_\theta + \delta$ where γ, δ are learnable parameters corresponding to the correlated noise w.r.t. full precision and the uncorrelated noise respectively. PTQD models the uncorrelated noise as Gaussian (*i.e.*, $\delta \sim \mathcal{N}(\mu_q, \sigma_q)$).

Variance Schedule Calibration for Trajectory Consistency Distillation (TCD) TCD’s Strategic Stochastic Sampling (SSS) using a quantized network $\mathcal{E}_\theta^Q(x_t, t)$ is given by:

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t \mathcal{E}_\theta^Q(x_t, t)}{\alpha_t} + \sigma_{s'} \mathcal{E}_\theta^Q(x_t, t) \right) + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z} \quad (15)$$

Using PTQD’s linear parametrization for the quantization error, we substitute $\mathcal{E}_\theta^Q(x_t, t) = (1 + \gamma) \cdot \mathcal{E}_\theta + \delta$:

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t((1+\gamma) \cdot \mathcal{E}_\theta(x_t, t) + \delta)}{\alpha_t} + \sigma_{s'}((1+\gamma) \cdot \mathcal{E}_\theta(x_t, t) + \delta) \right) + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z} \quad (16)$$

$$= \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t(1+\gamma)\mathcal{E}_\theta(x_t, t)}{\alpha_t} + \sigma_{s'}(1+\gamma)\mathcal{E}_\theta(x_t, t) - \frac{\alpha_{s'}\sigma_t\delta}{\alpha_t} + \sigma_{s'}\delta \right) + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z} \quad (17)$$

$$= \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t(1+\gamma)\mathcal{E}_\theta(x_t, t)}{\alpha_t} + \sigma_{s'}(1+\gamma)\mathcal{E}_\theta(x_t, t) \right) + \frac{\alpha_s}{\alpha_{s'}} \left(\sigma_{s'} - \frac{\alpha_{s'}\sigma_t}{\alpha_t} \right) \delta + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z} \quad (18)$$

$$(19)$$

The correlated noise can be corrected by applying:

$$\frac{\mathcal{E}_\theta^Q(x_t, t)}{1+\gamma} = \frac{(1+\gamma)\mathcal{E}_\theta(x_t, t) + \delta}{1+\gamma} \quad (20)$$

$$= \mathcal{E}_\theta(x_t, t) + \frac{\delta}{1+\gamma} \quad (21)$$

The resultant SSS sampling step becomes:

$$\mathbf{x}_s = \frac{\alpha_s}{\alpha_{s'}} \left(\alpha_{s'} \frac{\mathbf{x}_t - \sigma_t\mathcal{E}_\theta(x_t, t)}{\alpha_t} + \sigma_{s'}\mathcal{E}_\theta(x_t, t) \right) + \frac{\alpha_s}{\alpha_{s'}(1+\gamma)} \left(\sigma_{s'} - \frac{\alpha_{s'}\sigma_t}{\alpha_t} \right) \delta + \sqrt{1 - \frac{\alpha_s^2}{\alpha_{s'}^2}} \mathbf{z} \quad (22)$$

$$(23)$$

The variance schedule becomes:

$$\sigma_\delta^2 = 1 - \frac{\alpha_s^2}{\alpha_{s'}^2} - \left(\frac{\alpha_s}{\alpha_{s'}(1+\gamma)} \left(\sigma_{s'} - \frac{\alpha_{s'}\sigma_t}{\alpha_t} \right) \right)^2 \delta^2 \quad (24)$$

$$= 1 - \frac{\alpha_s^2}{\alpha_{s'}^2} \left(1 - \frac{(\sigma_{s'} - \frac{\alpha_{s'}\sigma_t}{\alpha_t})^2}{(1+\gamma)^2} \delta^2 \right) \quad (25)$$

$$= 1 - \frac{\alpha_s^2}{\alpha_{s'}^2} \left(1 - \left(\frac{\delta(\sigma_{s'} - \frac{\alpha_{s'}\sigma_t}{\alpha_t})}{(1+\gamma)} \right)^2 \right) \quad (26)$$

Since $\mathbf{z} \sim N(\mu_\delta, \sigma_\delta)$, we must handle the edge case when $\sigma_\delta < 0$. If the variance is negative, we simply set $\sigma_\delta = 0$.

Upon comparing Q-Sched to PTQD you may ask "Why is Q-Sched able to learn a better noise schedule when it is also a linear correction?" Q-Sched learns scalar coefficients on x_t and \mathcal{E}_θ that are optimized with respect to the reference-free JAQ loss. This allows us to learn a new schedule with linear corrections to improve our overall noise schedule, rather than matching the existing full precision schedule. This is an important distinction from PTQD, which tries to learn a linear correction with respect to full precision, which may not be possible since quantization produces a nonlinear distortion on the diffusion model. In short, PTQD attempts to match the full precision sampling trajectory, whereas Q-Sched aims to learn a new sampling trajectory given a compressed \mathcal{E}_θ .

H Strict Guarantees for Quantization-Aware Scheduling

Let us consider the few-step sampling trajectories for the pre-trained and quantized models, parametrized by $\mathcal{E}_\theta(t)$ and $\mathcal{E}_\theta^Q(t)$ respectively. These two few-step diffusion models sample at

the same time-steps, $0 = t_0 < t_1, t_2 \dots t_N = T$, where N represents the number of steps in the few-step model. For ease of notation, we will use the time-step 0 to refer to t_0 and 1 to refer to t_1 , etc. A denoising step going from time $t + 1 \rightarrow t$, produces a partially denoised image, x_t , and its quantized counterpart, x_t^Q . Following directly from Equation 9, the denoising error, $\Delta x_t = x_t - x_t^Q$, can be explicitly computed as:

$$\Delta x_t = \frac{\alpha_t}{\alpha_{t'}} \left(\alpha_{t'} \frac{\Delta x_{t+1} - \sigma_{t+1}(\mathcal{E}_\theta(t+1) - \mathcal{E}_\theta^Q(t))}{\alpha_{t+1}} + \sigma_{t'}(\mathcal{E}_\theta(t+1) - \mathcal{E}_\theta^Q(t+1)) \right) \quad (27)$$

$$= \frac{\alpha_t}{\alpha_{t+1}} \Delta_{t+1} + \frac{\alpha_t}{\alpha_{t'}} \left(\sigma_{t'} - \frac{\sigma_{t+1}}{\alpha_{t+1}} \right) (\mathcal{E}_\theta(t+1) - \mathcal{E}_\theta^Q(t+1)) \quad (28)$$

$$= k_t \Delta x_{t+1} + m_t \Delta \mathcal{E}_\theta(t+1) \quad (29)$$

where we define the sampler coefficients as $k_t = \frac{\alpha_t}{\alpha_{t+1}}$, $m_t = \frac{\alpha_t}{\alpha_{t'}} \left(\sigma_{t'} - \frac{\sigma_{t+1}}{\alpha_{t+1}} \right)$ and denote the change in the network as $\Delta \mathcal{E}_\theta(t) = \mathcal{E}_\theta(t) - \mathcal{E}_\theta^Q(t)$. Assuming the initial denoised image is the same ($x_N = x_N^Q$), the error in the final denoised image, Δx_0 , is given by:

$$\Delta x_0 = k_0 \Delta x_1 + m_0 \Delta \mathcal{E}_\theta(1) \quad (30)$$

$$= k_0 k_1 k_2 \dots (k_N \Delta x_N + m_{N-1} \Delta \mathcal{E}_\theta(N)) + \dots + k_0 k_1 m_2 \Delta \mathcal{E}_\theta(3) + k_0 m_1 \Delta \mathcal{E}_\theta(2) + m_0 \Delta \mathcal{E}_\theta(1) \quad (31)$$

$$= \sum_{s=1}^S \left(\prod_{v=0}^{s-2} k_v \right) m_{s-1} \Delta \mathcal{E}_\theta(s) \quad (32)$$

H.1 Expected Quantization Error

The average expected error over all images in a given dataset, $x_0 \in \mathcal{D}$, is given by:

$$E[||\Delta x_0||] = \sum_{s=1}^S \left(\prod_{v=0}^{s-2} k_v \right) m_{s-1} E[||\Delta \mathcal{E}_\theta(s)||] \quad (33)$$

since $E[||\Delta x_0||]$ is a homogeneous function.

In Q-Sched, we apply our quantization-aware pre-conditioning on every noise coefficient: $\tilde{m}_t = c_t^x \cdot m_t$ and $\tilde{k}_t = c_t^x \cdot k_t$. Let us denote the expected error induced by Q-Sched with respect to the pre-trained model's x_0 as $E[||\Delta \tilde{x}_0||]$.

We empirically show in Tables 1 and 2 that $E[||\Delta x_0||] \neq 0$ since the images produced by the naive quantization method produce a different FID from the original pre-trained model's image distribution. Since Equation 33 is a linear function of $k_t, m_t, \forall t \in 1 \dots N$, and there is a global minimum at $E[||x_0 - x_0||] = 0$, it must be that $\exists \tilde{m}_t^*, \tilde{k}_t^* \forall t$ such that $E[||\Delta \tilde{x}_0||] < E[||\Delta x_0||]$. In short, we guarantee that there exists quantization-aware coefficients that strictly improve our expected quantization error over naive quantization.

H.2 Aside: Positive Sampler Coefficients

The TCD Scheduler has $\beta_0 = 0.0085, \beta_N = 0.012, \alpha_t = 1 - \beta_t, \sigma_t = \prod_{i=0}^t \alpha_i$ with a scaled linear schedule:

$$\beta_t = \left(\sqrt{\beta_0} + t \cdot (\sqrt{\beta_N} - \sqrt{\beta_0}) \right)^2 \quad (34)$$

Therefore: $1 > \alpha_0 > \alpha_1 > \dots > \alpha_N > 0$ and $1 > \sigma_0 > \sigma_1 > \dots > \sigma_N > 0$. We note the $t' = (1 - \gamma)t$ where $\gamma \in [0, 1]$, so $t' \leq t$. This implies that $\sigma_{t'} > \sigma_{t+1}$ so:

$$k_t > 0 \quad , \quad m_t = \frac{\alpha_t}{\alpha_{t'}} \left(\sigma_{t'} - \frac{\sigma_{t+1}}{\alpha_{t+1}} \right) > 0 \quad (35)$$

This illustrates that $k_t, m_t \in \mathbb{R}^+$.

I Qualitative Analysis

We present qualitative results from a few hand-picked examples of good and bad results for our method in Figure 9 and Figure 10 respectively. Good examples were easier to find than bad examples, as our optimized Q-Sched scheduler tends to provide more detail and an enhanced version of the 4W8A model. In particular, we see additional generative details in the first and third prompts in Figure 9, where our method generates fine-grained details that are not present in either full precision or 4W8A. We also note that our method provides additional texture such as on the fried chicken sandwich on seventh prompt and more realistic hair textures in the fourth prompt.

Our method still struggles with multi-person images such as the first and second prompts in Figure 10. We note that plain text is not improved with our method, as shown in the third prompt. We also find that some artifacts are not rectified with our model, such as in the sixth prompt where the plane silhouette is not properly generated. We also notice a slight divergence in some images from the original prompt such as in the seventh prompt, where a vase is not present in our generated image. Overall, these bad examples are fewer than the good examples and our significant FID improvement shows the Q-Sched is capable of providing very solid results.

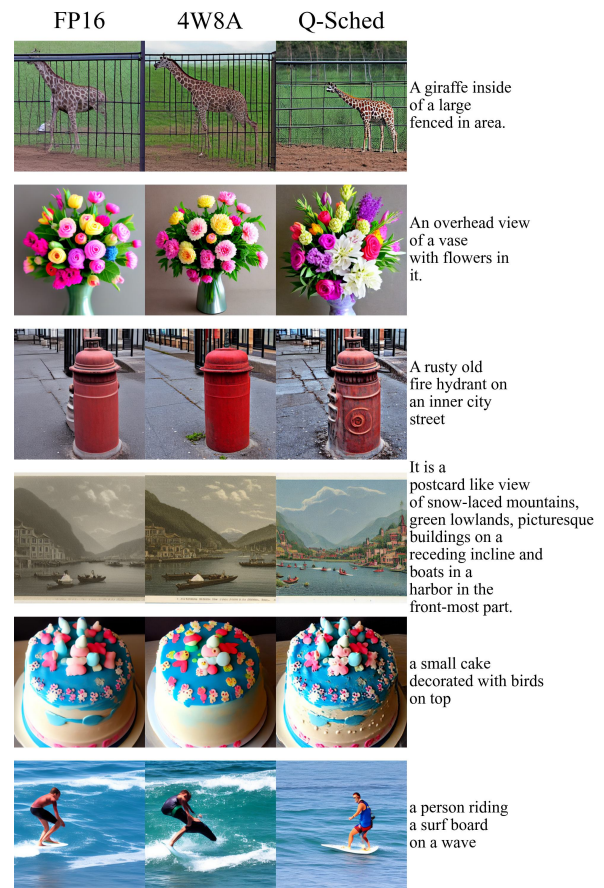


Figure 9: Selected prompts where Q-Sched outperforms full precision.



Figure 10: Selected prompts where Q-Sched fails to improve over 4W8A.