

# DIFFUSION MODELS AS INTRINSIC DISTRIBUTION ESTIMATORS FOR SELF-VERIFYING INFERENCE-TIME SCALING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

To enhance sample quality beyond their standard outputs, diffusion models typically rely on inference-time scaling, a process that necessitates external verifiers. We challenge this dependency by proposing a framework that reframes the generative model itself as an intrinsic distribution estimator. Our framework provides the theoretical base and empirical evidence for this, showing that the distance between independent noise and diffusion model output serves as a proxy for a sample’s distributional conformity. This insight enables our proposed method, Self-Verifying inference-time scaling method to directly assess at intermediate denoising step and to eliminate the need for external modules. Experiment results demonstrate that our scaling method achieves consistent improvements across diverse benchmarks in fidelity, preference, and compositionality. Our study establishes that the process of generating diffusion models is also an evaluative process, opening new avenues toward more resource-efficient and intrinsically aware generative models.

## 1 INTRODUCTION

Diffusion models have become a cornerstone of modern artificial intelligence, achieving state-of-the-art results across a multitude of domains. Their remarkable, modality agnostic expressive power has made them the mainstream choice for high-fidelity generation of images (Esser et al., 2024; Saharia et al., 2022; Betker et al., 2023), audio (Lee et al., 2025), and text (Li et al., 2022; Arriola et al., 2025). Beyond their generative prowess, the unique characteristics of diffusion models, such as their multi-step denoising process and their capacity for implicit likelihood estimation have enabled a range of intrinsic capabilities from zero-shot classification (Clark & Jaini, 2023; Li et al., 2023) and out-of-distribution (OOD) detection (Aithal et al., 2024; Zhang et al., 2025; Yao et al., 2024; Heng et al., 2024) to unsupervised image editing (Mokady et al., 2023; Huang et al., 2025).

A particularly promising application leveraging these properties is inference-time scaling, which aims to iteratively refine generation quality. However, prevailing methods typically rely on external pretrained reward model to evaluate a fully denoised sample Ma et al. (2025); Fernandes et al. (2025); Xie et al. (2025); Li et al. (2025). If the sample fails to meet certain criteria, the denoising process is reverted to an intermediate noisy state and corrected. This paradigm is inherently computationally expensive, as it requires at least one full denoising cycle for verification and introduces a dependency on auxiliary models.

This raises a fundamental question: can a diffusion model verify its own generation quality during the denoising process, without external supervision? In this work, we answer in the affirmative way by proposing a new perspective. We assume an objective as  $\|a - M(f(a, b))\|^2$ , where conditional generative model  $M$  minimizes the objective,  $a \sim \mathcal{N}$ ,  $b \sim q_{data}$ , and  $f$  is linear combination function. We prove well-trained model  $M^*$  can be viewed as an implicit distribution estimator and consequently distinguishing in-distribution (ID) and OOD samples, described in Eq. 5.

In addition, we demonstrate through reparameterization that this theoretical framework directly maps to the training objective of diffusion models. The model’s noise prediction error at any given timestep serves as a proxy for the sample’s conformity to the learned data distribution. We empirically validate this hypothesis using a pre-trained Stable Diffusion XL (SDXL) (Podell et al., 2023)

054 model, showing it can robustly differentiate ID ImageNet-1k (Deng et al., 2009) samples from OOD  
 055 dataset like ImageNet-A (Hendrycks et al., 2021) and ImageNet-C (Hendrycks & Dietterich, 2019).  
 056

057 Building on this principle, we introduce a self-verifying inference-time scaling method that can  
 058 estimate the reward at the middle of denoising steps with diffusion model itself without relying on  
 059 the external verifier, by estimating which candidates has smaller distance to independent noise than  
 060 the others. This approach leads to consistent performance improvements in unconditional generation  
 061 on ImageNet and LSUN (Yu et al., 2016), human preference alignment in text-to-image synthesis,  
 062 and compositional generation capabilities.

063 The main contributions of our work are as follows:

- 064 • We introduce a new framework that formalizes conditional noise generative models as implicit **distribution estimator**, capable of distinguishing ID from OOD components, by comparing between independent noise and model prediction.
- 065 • We demonstrate that this framework can be applied directly to diffusion models **theoretically** by reparameterization and **empirically** proven in ImageNet-1k, A, and C experiments.
- 066 • We propose a new inference-time scaling method that is **self-verifiable** and efficient, directly utilizing intermediate denoising sample to enhance generation quality **without requiring external verifiers or full generation steps**.

## 074 2 RELATED WORK

075 Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) are generative model that learn to  
 076 create new data by reversing a gradual noising process. This is accomplished by learning the reverse  
 077 of the predefined forward process, which is structured as a Markov chain. The forward process,  $q$ ,  
 078 progressively adds Gaussian noise to the original data  $x_0 \sim q(x_0)$  over  $T$  timesteps. The transition  
 079 at each timestep  $t$  is defined as:  $q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t\mathbf{I})$  where  $\beta_t \in (0, 1)$  are  
 080 small positive constants defined by a variance schedule. A key property of this process is that we  
 081 can sample  $x_t$  at an arbitrary timestep  $t$  directly from  $x_0$  in a closed form. Letting  $\alpha_t := 1 - \beta_t$  and  
 082  $\bar{\alpha}_t := \prod_{i=1}^t \alpha_i$ , expressed as:

$$083 q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \quad (1)$$

084 where  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ . The core of generation lies in learning the reverse process,  $p_\theta$ , which denoises  
 085 the data starting from pure noise  $x_T \sim \mathcal{N}(0, \mathbf{I})$  by iteratively sampling  $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$  until  $x_0$   
 086 is produced. This process is modeled as  $p_\theta(x_{t-1}|x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$  in the DDPM  
 087 (Denoising Diffusion Probabilistic Models) (Ho et al., 2020). Instead of optimizing the variational  
 088 lower bound (VLB) on log-likelihood, DDPM proposed a simplified objective that is proportional  
 089 to the VLB. The model,  $\epsilon_\theta(x_t, t)$ , is trained to predict the noise component  $\epsilon$  that was added to the  
 090 data at timestep  $t$ . The objective function is:

$$091 L_{\text{simple}}(\theta) := \mathbb{E}_{t, x_0, \epsilon} [||\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)||^2] \quad (2)$$

092 This noise prediction network  $\epsilon_\theta$  has predominantly been implemented using a U-Net (Ronneberger  
 093 et al., 2015; Ho et al., 2020; Dhariwal & Nichol, 2021; Rombach et al., 2022). More recently,  
 094 architectures like the DiT (Diffusion transformer) (Peebles & Xie, 2023; Esser et al., 2024) have  
 095 replaced the U-Net backbone with a Transformer, demonstrating scalability.

096 Following the initial success of DDPMs, subsequent work has focused on improving sample quality,  
 097 computational efficiency, and sampling speed. Sample fidelity was enhanced through techniques like  
 098 guidance (Dhariwal & Nichol, 2021; Ho & Salimans, 2022; Karras et al., 2024), while Latent Diffu-  
 099 sion Models (LDMs) (Rombach et al., 2022) drastically reduced computational costs by operating  
 100 in a compressed latent space. To address the slow sampling speed, Denoising Diffusion Implicit  
 101 Models (DDIMs) (Song et al., 2020) introduced a deterministic sampling process, and reformu-  
 102 lating diffusion as an SDE/ODE enabled the use of fast numerical solvers like DPM-Solver (Lu  
 103 et al., 2022). More recently, Consistency Models (Song et al., 2023; Luo et al., 2023) and Recti-  
 104 fied Flow (Liu et al., 2022; 2023) have achieved generation in a single or very few steps without  
 105 significant degradation by learning a direct noise-to-data mapping.  
 106  
 107

The ability of diffusion models to precisely model complex data distributions has also made them a powerful tool for applications beyond generation, such as anomaly detection (AD) or Out-of-Distribution (OOD) prediction. Current diffusion-based AD methods typically fall into two categories from Liu et al. (2025b) of reconstruction, density-based, and the other is the the flow-based approach. Reconstruction-based methods (Bercea et al., 2023; Zhang et al., 2025; Yao et al., 2024) leverage the principle that a model trained on normal data will yield higher reconstruction errors for anomalous inputs. In contrast, density-based approaches (Livernoche et al., 2023; Luo, 2023) use the learned distribution directly, often employing the magnitude of the score function or an estimated diffusion time as an anomaly score. Flow-based research (Heng et al., 2024; Aithal et al., 2024) utilizes the denoising trajectories to distinguish In-Distribution (ID) samples and OOD samples.

Finally, significant research has focused on aligning model outputs with human preferences and further enhancing image quality. One prominent approach involves fine-tuning the model directly on preference data or utilizing a reward model, techniques successful in language models. For instance, DPO-diffusion (Wallace et al., 2024) adapts Direct Preference Optimization (DPO) (Rafailov et al., 2023) and Flow-GRPO (Liu et al., 2025a) leverages Group Relative Policy Optimization (GRPO) (Shao et al., 2024) with Rectified Flow model to few-step denoising in order to achieve online-RL with an external verifier. Complementary to fine-tuning, another line of work explores inference-time scaling techniques. These methods typically generate multiple candidates and use an external verifier to select the optimal output Xie et al. (2025); Fernandes et al. (2025). For example, research has explored techniques such as Best-of-N sampling, Zero-Order Search, and Search over Paths, which compare fully denoised samples to identify the one that best aligns with desired criteria or quality metrics (Ma et al., 2025). This strategy of leveraging an external verifier to guide generation at inference time has become a common approach for scaling the performance of diffusion models.

### 3 IDENTIFYING OUT-OF-DISTRIBUTION VIA DIFFUSION MODEL AS A DISTRIBUTION ESTIMATOR

#### 3.1 CONDITIONAL GENERATIVE MODEL AS A DISTRIBUTION ESTIMATOR

In this subsection, we investigate the robustness of the distribution estimation model to out-of-distribution (OOD) perturbations in a controlled generative setting. Our goal is to formalize and analyze how model performance degrades when a component of the input data is drawn from outside its training distribution.

**Generative Process and Learning Objective** Let  $a \in \mathbb{R}^d$  be a latent variable representing a signal, drawn from a standard multivariate Gaussian distribution  $a \sim \mathcal{A} = \mathcal{N}(0, \mathbf{I}_d)$ . Let  $b \in \mathbb{R}^d$  be a structured data component, drawn from an arbitrary high-dimensional data distribution  $B$  with mean  $\mu_b$  and covariance  $\Sigma_b$ . We assume  $a$  and  $b$  are statistically independent.

The model input  $x \in \mathbb{R}^d$  is generated by a linear combination of these components,  $x = f(a, b) = c_1 a + c_2 b$ , where  $0 \leq c_1, c_2 \leq 1$  are scalar coefficients controlling the relative influence of  $a$  and  $b$ . We consider a model  $M : \mathbb{R}^d \rightarrow \mathbb{R}^d$  that learns to produce samples following the conditional distribution of  $a$  given  $x$ .

Instead of comparing the model output with the same ground-truth sample  $a$ , we introduce an independent draw  $a' \sim \mathcal{A}$  to treat  $M$  as a distributional estimator of  $a$ : a well-trained model should produce outputs whose distribution matches that of  $a$ . The optimal model  $M^*$  minimizes the Mean Squared Error (MSE), which is equivalent to the statistical risk, with respect to this independent reference  $a'$ .

$$M^* = \underset{M}{\operatorname{argmin}} \mathcal{L}(M), \quad \mathcal{L}(M) = \mathbb{E}_{a,b,a'} [\|a' - M(f(a,b))\|^2] \quad (3)$$

This risk decomposes as,

$$\mathbb{E}_{a,b,a'} [\|a' - M(f(a,b))\|^2] = \underbrace{\mathbb{E} \|a'\|^2}_{\text{prior variance (constant)}} + \mathbb{E}_{a,b} \|M(f(a,b))\|^2 \quad (4)$$

Since the first term is constant, minimizing  $\mathcal{L}(M)$  is equivalent to matching the second moment of  $M(f(a,b))$  to that of the true prior  $\mathcal{A}$ .

**In-Distribution (ID) vs Out-Of-Distribution (OOD) Risk** We now study model performance when faced with an outlier sample  $b^*$  not representative of  $\mathcal{B}$ .

Such a  $b^*$  may lie in a low probability region under  $\mathcal{B}$  or originate from a different distribution. We define two key quantities. First, ID Risk, the expected loss achieved by the optimal model  $M^*$  on data sampled from the training distribution:  $R_{ID} = \mathbb{E}_{a,b,a'}[\|a' - M^*(c_1a + c_2b)\|^2]$  and second OOD Risk, the expected loss when  $b$  is replaced by fixed outlier  $b^*$ :  $R_{OOD}(b^*) = \mathbb{E}_{a,a'}[\|a' - M^*(c_1a + c_2b^*)\|^2]$ .

Because the constant term  $\mathbb{E}\|a'\|^2$  cancels in the comparison, the difference between  $R_{ID}$  and  $R_{OOD}(b^*)$  is fully determined by the change in the second moment of the model outputs.

We hypothesize that a model optimized for the statistics of  $\mathcal{B}$  will suffer a degradation in performance when evaluated on  $b^*$ , leading to:

$$\begin{aligned} R_{ID} &\leq R_{OOD}(b^*) \\ \Rightarrow \mathbb{E}_{a,b,a'}[\|a' - M^*(c_1a + c_2b)\|^2] &\leq \mathbb{E}_{a,a'}[\|a' - M^*(c_1a + c_2b^*)\|^2] \end{aligned} \quad (5)$$

**Conditions for the Inequality** The inequality  $R_{ID} \leq R_{OOD}(b^*)$  is guaranteed under the following sufficient conditions:

- Bayes-Optimal Model:  $M^*$  exactly reproduces the true conditional distribution  $M^*(x) \stackrel{d}{=} a|x$ . In this case, the expected loss reduces to the sum of the prior variance and the conditional second moment of  $a$ . Replacing  $b$  with  $b^*$  induces a conditional distribution different from the training distribution, which can only increase (or leave unchanged) the expected second moment, yielding  $R_{ID} \leq R_{OOD}(b^*)$ .
- Well-Specified Generative Model: The assumed generative process  $x = f(a, b)$  and prior  $p(a) = \mathcal{N}(0, I_d)$  match the true data-generating process. If the model is mis-specified or underfitted, there may exist  $b^*$  such that  $R_{OOD}(b^*) < R_{ID}$ .
- Atypicality of  $b^*$ : The sample  $b^*$  lies in a statistically atypical region of  $\mathcal{B}$ , e.g.,  $p_{\mathcal{B}}(b^*) \ll p_{\mathcal{B}}(\mu_b)$ . If  $b^*$  is typical, then  $R_{OOD}(b^*) \approx R_{ID}$ .

Under these conditions, the OOD risk must be at least as large as the ID risk, with strict inequality whenever  $p(a|x)$  under  $b^*$  sufficiently deviates from the training-time conditional distribution.

### 3.2 DIFFUSION MODEL AS A DISTRIBUTION ESTIMATOR

Let the objective from our assumed model,  $M$ , be defined as:

$$L = \mathbb{E}_{t \sim p(t)} [\mathbb{E}_{a \sim \mathcal{N}(0, I), b \sim q} [\|a - M(c_{1,t}a + c_{2,t}b, c_{1,t}, c_{2,t})\|^2]] \quad (6)$$

where  $(c_{1,t}, c_{2,t})$  are coefficients indexed by a timestep  $t$  drawn from a distribution  $p(t)$ .

We begin by establishing a set of one-to-one correspondences to align our proposed circumstances with the DDPM framework. First, we map the variable  $a \leftrightarrow \epsilon$  and  $b \leftrightarrow x_0$ , where both  $a$  and  $\epsilon$  are sampled from  $\mathcal{N}(0, I)$ , and both  $b$  and  $x_0$  are from the data distribution  $q$ . Second, the coefficients  $(c_{1,t}, c_{2,t})$  are defined to match the DDPM schedule:  $c_{1,t} := \sqrt{1 - \bar{\alpha}_t}$  and  $c_{2,t} := \sqrt{\bar{\alpha}_t}$ . Given that the timestep  $t$  uniquely determines these coefficients, our model  $M$  (conditioned on coefficients) becomes functionally identical to the DDPM model  $\epsilon_\theta$  (conditioned on  $t$ ). This equivalence is formally expressed as  $M(\cdot, c_{1,t}, c_{2,t}) \equiv \epsilon_\theta(\cdot, t)$ .

By substituting these mappings directly into our redefined objective (Eq. 6), we obtain:

$$L = \mathbb{E}_{t \sim p(t)} [\mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), x_0 \sim q} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2]] \quad (7)$$

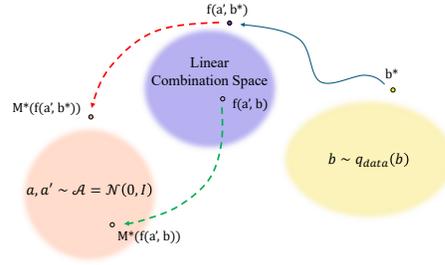


Figure 1: Visualization of the model’s distribution estimation for ID and OOD data. The **green dashed arrow** shows that the optimal model  $M^*$  successfully estimates the true distribution when conditioned on an ID sample  $f(a', b)$ . Conversely, the **red dashed arrow** illustrates the model’s failure to estimate the distribution when conditioned on an OOD sample  $f(a', b^*)$ .

This resulting expression is identical in form and substance to the DDPM objective (Eq. 2). Therefore, we have formally shown that  $L = L_{DDPM}$ . The proposed framework is an exact reparameterization of the DDPM noise-prediction training paradigm.

### 3.3 IS YOUR DIFFUSION MODEL REALLY DISTRIBUTION ESTIMATOR?

To empirically validate our theoretical framework (Eq. 5), we conduct experiments to determine if a pre-trained diffusion model can effectively distinguish between in-distribution (ID) and out-of-distribution (OOD) data by measuring noise prediction discrepancies.

**Distribution Estimation Setup** We employ Stable Diffusion XL (SDXL) (Podell et al., 2023), a publicly available latent diffusion model, as our off-the-shelf distribution oracle. For our in-distribution data, we use the validation set of ImageNet-1k (Deng et al., 2009), which represents the model’s learned data manifold. To challenge the model with out-of-distribution samples, we leverage two datasets: ImageNet-A (Hendrycks et al., 2021), a curated dataset of ImageNet-1k that are naturally adversarial to classifiers, and ImageNet-C (Hendrycks & Dietterich, 2019), a dataset for evaluating robustness against common visual corruptions.

Based on our assumption, the diffusion model will exhibit a higher noise prediction error for OOD samples than for ID samples. Given a pair of images, we first provide the model with a prompt "a photo of {class\_name}", corresponding to the image’s ground-truth label. Next, we compute the L2 distance between the independent noise  $\epsilon'$  and the model’s predicted noise  $\epsilon_\theta(x_t(x, t, \epsilon))$  at several distinct timesteps, where  $\epsilon, \epsilon'$  are independently sampled from the Gaussian distribution. This "noise distance" serves as a proxy for how well the image conforms to the model’s learned distribution. Finally, we classify the image with the lower average noise distance across these timesteps as the in-distribution sample. A prediction is deemed correct if the model successfully identifies the ImageNet-1k sample (ID).

For the ImageNet-A comparison, we randomly sample 10 pairs of images for each of ImageNet-A’s 200 classes, comparing each ImageNet-A sample against a randomly selected ImageNet-1k sample from the same class. For the ImageNet-C evaluation, we compare each corrupted image against its original, clean counterpart from ImageNet-1k, using 5 such pairs for every class. Due to computational constraints, we limit our ImageNet-C analysis to four corruption types (defocus blur, contrast, elastic transform, saturate) at severity levels 1 and 3.

**Distribution Estimation Results** Our experiments confirm that the noise prediction error of SDXL is a reliable indicator for distinguishing ID from OOD data.

| ImageNet-A | ImageNet-C |       |          |       |                   |       |          |       |
|------------|------------|-------|----------|-------|-------------------|-------|----------|-------|
|            | defocus    |       | contrast |       | elastic transform |       | saturate |       |
|            | 1          | 3     | 1        | 3     | 1                 | 3     | 1        | 3     |
| 56.5%      | 98.6%      | 99.3% | 98.4%    | 99.6% | 92.8%             | 89.7% | 89.3%    | 80.7% |

Table 1: Accuracy of our OOD detection method on subsets of the ImageNet-A and ImageNet-C dataset. The model classifies which image in a pair is ID based on noise prediction error.

When tasked with discriminating between ImageNet-1k and the challenging, naturally adversarial samples of ImageNet-A, our method achieved an accuracy of **56.5%**. While modest, the result is significantly above 50% of random choice, demonstrating that the model can perceive near-OOD images, even though these OOD images came from the ID set.

The model’s discriminative power was substantially more pronounced on the ImageNet-C dataset, as detailed in Table 1. As the results show, the model achieves near-perfect accuracy on simpler corruptions like defocus and contrast. Interestingly, for structurally complex corruptions such as elastic transform and saturate, the accuracy remains high but shows a slight decrease as severity intensifies. A plausible explanation is that the severe artifacts introduce high-frequency distortions that statistically mimic the Gaussian noise the model is trained to predict, making the OOD samples harder to distinguish from their ID counterparts.

**Algorithm 1** Self-Verifying Inference-Time Scaling

---

**Require:** Diffusion model  $\epsilon_\theta$ , denoising timestep set  $\mathbf{t} = \{T..1\}$ , normal dist.  $\mathcal{N}$ , non-deterministic scheduler  $S$ , inference-time scaling factor at its corresponding timestep  $n_t$ .

- 1: Sample  $\epsilon = \{\epsilon_1, \dots, \epsilon_{n_T}\}$  from  $\mathcal{N}$  ▷ Sample  $n$  noise vectors
- 2: **for**  $i = 1 \rightarrow n_T$  **do**
- 3:      $x_i \leftarrow \epsilon_i$  ▷ Initialization
- 4: **end for**
- 5: **for**  $t$  in  $\mathbf{t}$  **do**
- 6:     **for**  $i = 1 \rightarrow n_t$  **do**
- 7:         noise\_dist $_i \leftarrow \|\epsilon_\theta(x_{i,t}(\epsilon, t), t) - \epsilon'\|^2$  ▷ where  $\epsilon, \epsilon' \sim \mathcal{N}$
- 8:     **end for**
- 9:      $i_{\text{best}} \leftarrow \text{argmin}_i(\text{noise\_dist}_i)$  ▷ Choose best sample
- 10:      $x_t \leftarrow x_{i_{\text{best}}, t}$
- 11:     **for**  $i = 1 \rightarrow n_t$  **do**
- 12:          $x_{i,t-1} \leftarrow S.\text{step}(\epsilon_\theta(x_t, t), x_t, t)$  ▷ Denoise  $n$  times in parallel
- 13:     **end for**
- 14: **end for**
- 15: **return**  $x_0$

---

#### 4 INFERENCE TIME SCALING THROUGH DISTRIBUTION ESTIMATION OF THE DIFFUSION MODEL

The remarkable performance of large language models (LLMs) has been significantly enhanced by inference-time strategies (Muennighoff et al., 2025) that expand the search space for solutions, such as chain-of-thought prompting (Kojima et al., 2022) and self-consistency (Bartsch et al., 2023), followed by a verification or selection step (Yao et al., 2023; Liu et al., 2025c). These methods leverage the model’s extensive knowledge to generate and evaluate multiple candidate outputs, ultimately improving final performance without costly retaining (Snell et al., 2024).

Current inference-time scaling approaches for diffusion models often employ an external verifier to score and select the best candidates. This approach typically requires intermediate samples to be fully denoised into a clean image at each evaluation step, consequently costs intensive.

We propose a reasonable approach that circumvents these complexities. We have shown that the diffusion model can be a distribution estimator (Eq. 5), hypothetically and empirically. Building on this following, we further posit that this inherent capability can be repurposed as an internal quality verifier, allowing the model to assess the plausibility at its intermediate generative trajectories, by comparing the noise distance of the denoising candidates.

**Self-Verifying Inference-Time Scaling for Diffusion Model** Our algorithm, Self-Verifying Inference-Time Scaling, operationalizes this principle through a step-wise choice and branching procedure. The process begins with setting a number of candidates at timestep  $n_t$  and samples  $n_T$  noises from a Gaussian distribution. At each timestep  $t$  from  $T$  down to 1, we evaluate all current candidates using verification score, compute independent noise distance with model predict noise (line 7 from Algorithm.1). The single candidate with the lowest noise distance is selected, and all other are pruned. A stochastic denoising scheduler is then used to branch from this single best candidate, generating a diverse set of  $n_t$  candidates for the subsequent step. This iterative cycle of verifying, selecting, and branching efficiently guides the generation process along the most self-preferred trajectory to produce the final image.

## 5 EXPERIMENTS

We conduct a comprehensive set of experiments to validate the effectiveness of our proposed method, which is called Self-Verify. We evaluate its performance against a vanilla baseline on both unconditional and text-to-image generation tasks across a variety of benchmarks. Furthermore, we perform a detailed ablation study to analyze the impact of different hyperparameter choices and to identify an efficient configuration for our method.

## 5.1 EXPERIMENTAL SETUP

Our method, Self-Verify, is compared against a standard vanilla sampling baseline. Based on our ablation studies (detailed in Section 5.2), we adopt an efficient default configuration for Self-Verify: we maintain 4 candidates for the first four timesteps and then prune to a single trajectory for the remaining steps. This results in a modest increase in the number of function evaluations (NFE) from 50 for the vanilla sampler to 62 for Self-Verify.

**FID evaluation** For unconditional image generation, we evaluate Fréchet Inception Distance (FID) (Heusel et al., 2017) and Inception score (IS) (Salimans et al., 2016). We use the pre-trained ADM models (Dhariwal & Nichol, 2021) and their official checkpoints on the ImageNet 256x256 (Deng et al., 2009), LSUN-Cat 256x256, and LSUN-Bedroom 256x256 (Yu et al., 2016) datasets. Following standard evaluation protocol, we use the DDIM (Song et al., 2020) sampler with 50 steps and  $\eta = 1$ . We generate 10,000 samples for each experiment and compute FID and IS scores against the pre-computed reference statistics.

**Text-to-Image Human Preference Evaluation** To assess performance on text-to-image generation, we use the Stable Diffusion XL (Podell et al., 2023) model with the DPMSolverSDE (Lu et al., 2022) sampler. We use a set of 500 prompts from the test split of the Pick-a-Pic dataset (Kirstain et al., 2023), which is designed for modeling human preferences. We evaluate the generated images using a suite of automated scoring models: CLIPscore (Hessel et al., 2022) for image-text alignment, and PickScore (Kirstain et al., 2023), ImageReward (Xu et al., 2023), HPSv2 (Wu et al., 2023), and Aesthetic score, which serve as strong proxies for human aesthetic and compositional preferences.

**Object-Oriented Evaluation** To measure the model’s ability to follow compositional instructions, we utilize the GenEval (Ghosh et al., 2023). This benchmark evaluates performance across several fine-grained categories, including the rendering of colors, object count, spatial position, and color attribute. Experiments are conducted with the SDXL model and DPMSolverSDE sampler.

**Ablation** We conduct an ablation study to determine the most effective and efficient candidate scheduling function,  $n_t$ . This study uses the same environment as the text-to-image human preference evaluation. We investigate several strategies: (1) Constant Candidates: maintaining a constant number of candidates throughout, in this experiment 4, (2) maintaining  $n$  candidates for a fixed number of initial timestep  $t$  before reverting to vanilla sampling, and (3) dynamically decreasing the number of candidates from 16 to 1 using linear and exponential pruning schedules. We report the HPS score as the primary metric, with PickScore and ImageReward available in the Appendix A.2

## 5.2 EXPERIMENTAL RESULTS

**Main Results** Our experiments imply that Self-Verify shows consistent improvement in the quality of generated samples across a range of benchmarks and tasks. While the improvements sometimes modest, they demonstrate a systematic positive effect, validating the efficacy of our approach.

| Method      | ImageNet     |               | LSUN-Cat     |              | LSUN-Bedroom |              |
|-------------|--------------|---------------|--------------|--------------|--------------|--------------|
|             | FID ↓        | IS ↑          | FID ↓        | IS ↑         | FID ↓        | IS ↑         |
| Vanilla     | 20.73        | 86.346        | 19.89        | <b>4.974</b> | 8.54         | 2.368        |
| Self-Verify | <b>20.59</b> | <b>87.510</b> | <b>19.66</b> | 4.957        | <b>8.45</b>  | <b>2.373</b> |

Table 2: Unconditional generation results. Self-Verify improves FID scores across all datasets with marginal compute overhead.

For **unconditional generation**, we observe a consistent trend of improved FID scores across the ImageNet, LSUN-Cat, and LSUN-Bedroom datasets (Table 2). This suggests that by pruning less plausible paths, our method guides the sampling process along a trajectory that better aligns with the learned data distribution, leading to enhanced overall sample fidelity.

In the **text-to-image domain**, Self-Verify systemically leads to systemically higher scores on metrics that serve as strong proxies for human preference, including PickScore, ImageReward, and HPS (Table 3). It is particularly compelling that this enhancement is achieved without guidance from any external, human-aligned reward models during the inference process. We hypothesize this stems from the model’s learned distribution implicitly encoding human aesthetic biases, akin to a mere exposure effect (Zajonc, 2001). High-quality and aesthetically pleasing im-

| Method            | CLIP         | PickScore    | ImageReward   | HPS           | Aesthetic     |
|-------------------|--------------|--------------|---------------|---------------|---------------|
| Vanilla           | 8.71         | 22.02        | 0.8361        | 0.2906        | 5.8928        |
| Self-Verify (4/4) | <b>28.78</b> | <b>22.12</b> | <b>0.8972</b> | 0.2935        | 5.9329        |
| Self-Verify (8/1) | 28.70        | 22.10        | 0.8929        | <b>0.2936</b> | <b>5.9732</b> |

Table 3: Text-to-image human preference evaluation on prompts from the test split of the Pick-a-pic dataset. Self-verify (n/t) implies maintaining  $n$  candidates for a fixed number of initial timesteps  $t$ . Our method achieves the best performance across human preference metrics.



Figure 2: **Qualitative Comparison on Text-to-image domain.** Our approach yields outputs with stronger alignment to the textual prompt and higher aesthetic quality, whereas vanilla generations miss key prompt details or appear less coherent.

ages are more prevalent in the training data than flawed ones. Consequently, the model naturally learns a probabilistic distribution where desired samples are assigned higher density than undesired ones ( $p(x_{undesired}) < p(x_{desired})$ ). By seeking the most self-consistent generative trajectory, our method is inherently guided toward these high-density, high-preference regions. The qualitative comparison in Figure 2 highlights the effectiveness of our method across three representative prompts. These examples show that Self-Verify yields more visually consistent generations.

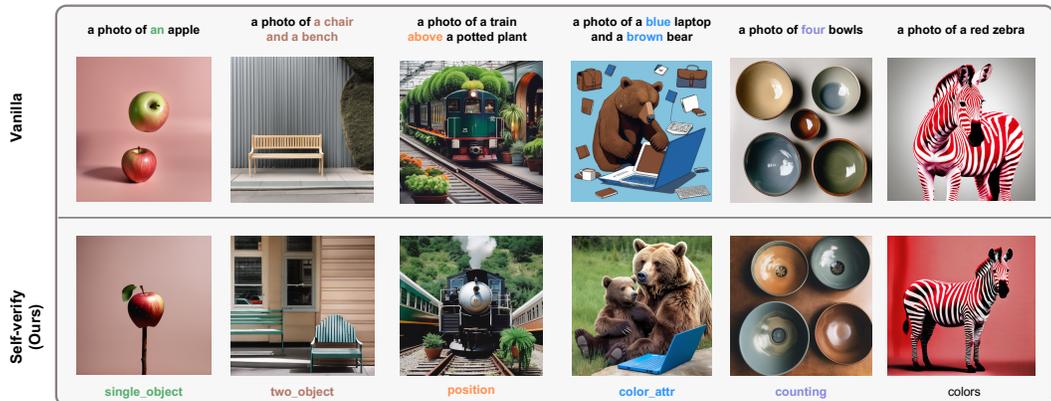


Figure 3: **Qualitative comparisons on the GenEval benchmark.** Compared to vanilla generations, Self-Verify (Ours) more faithfully follows prompt specifications across categories, leading to semantically accurate and coherent outputs.

| Category      | Vanilla | Self-Verify   |
|---------------|---------|---------------|
| single_object | 97.81%  | <b>98.44%</b> |
| two_object    | 75.76%  | <b>78.28%</b> |
| position      | 11.00%  | <b>11.25%</b> |
| color_attr    | 18.25%  | <b>20.50%</b> |
| counting      | 39.38%  | <b>41.88%</b> |
| colors        | 84.31%  | 84.31%        |
| Overall score | 0.5442  | <b>0.5578</b> |

Table 4: Object-oriented evaluation on GenEval benchmark. Self-Verify consistently improves compositional correctness.

Furthermore, when assessing compositional understanding with the **GenEval benchmark**, our method achieves a higher overall score (Table 4). The gains are consistent for all categories. This indicates that the internal verification process helps the model better adhere to intricate prompt details, resulting in more semantically accurate and coherent images. Self-Verify consistently improves compositional understanding by adhering to object count, position, color attributes, and counting, while vanilla generations fail on these fine-grained details (Figure 3).

| HPS                 |    | Candidates          |        |        |        |        |
|---------------------|----|---------------------|--------|--------|--------|--------|
|                     |    | 1                   | 2      | 4      | 8      | 16     |
| Timesteps           | 1  | 0.2906<br>(vanilla) | 0.2904 | 0.2915 | 0.2936 | 0.2929 |
|                     | 2  | -                   | 0.2915 | 0.2916 | 0.2918 | 0.2925 |
|                     | 4  | -                   | 0.2904 | 0.2935 | 0.2920 | 0.2917 |
|                     | 8  | -                   | 0.2893 | 0.2929 | 0.2930 | 0.2933 |
|                     | 16 | -                   | 0.2913 | 0.2911 | 0.2928 | 0.2940 |
| Constant Candidates |    | 0.2921              |        |        |        |        |
| Linear Pruning      |    | 0.2923              |        |        |        |        |
| Exponential Pruning |    | 0.2927              |        |        |        |        |

Table 5: Ablation study on candidate scheduling, reporting HPS score. The setting with 4 candidates for the first 4 steps offers a strong performance-cost trade-off. Constant candidates, Linear Pruning, and Exponential Pruning are described at Sec. 5.1

**Ablation Study** Our ablation study on the candidate scheduling function provides insight into the method’s cost-performance trade-off. The results confirm that applying verification for even a few initial steps yields sufficient performance gain over the vanilla baseline (Table 5). We found that the relationship between quality and computational cost is not linear; simply increasing the number of candidates or verification steps does not guarantee proportionally better results. The configurations chosen for our main experiments, which were 4 candidates for the first 4 steps and 8 candidates for the first step, were identified as an effective sweet spot. Additionally, the adaptive pruning results (below 3 rows) tell us these strategies cannot be strong substitutes.

## 6 LIMITATIONS AND FUTURE WORK

We acknowledge two primary limitations. First, our empirical validations of the suggested assumption and proposed method are confined to the visual domain. Validating the broader applicability of our framework is thus a critical next step. Investigating whether distribution estimation and its subsequent inference-time scaling holds for diffusion models in other modalities, such as audio synthesis and language modeling, will be essential to ascertain the domain generalizability of our approach.

Furthermore, the performance gains afforded by our method, while consistent, are incremental rather than transformative compare to approaches using an external verifier. A promising direction is therefore to bridge this performance gap, potentially by developing more nuanced aggregation methods for the intrinsic noise signal or by utilizing external verifiers to combine the schemes to amplify the model’s capacity.

## 7 CONCLUSION

In this work, we present a new framework for interpreting conditional generative models as implicit distribution estimators (Sec. 3.1). Through the reparameterization of conditional generative model, we established that the diffusion model can serve this role, using their noise prediction to estimate the likelihood that a sample belongs to the learned data distribution (Sec. 3.2). Our initial experiments empirically confirmed this hypothesis, showing that a standard pre-trained diffusion model can effectively distinguish between in-distribution and out-of-distribution samples (Sec. 3.3).

Building on this principle, we proposed an inference-time scaling methodology that leverages the diffusion model itself as an intrinsic quality verifier (Sec. 4). This self-verifying approach adjusts the generation process without relying on any external reward models. Our comprehensive experiments demonstrated that this method yields consistent improvements across a range of tasks, enhancing unconditional generation fidelity on ImageNet and LSUN, improving human preference scores in text-to-image synthesis, and boosting compositional capabilities (Sec. 5.2).

## REFERENCES

- 486  
487  
488 Sumukh K Aithal, Pratyush Maini, Zachary Chase Lipton, and J Zico Kolter. Understanding hallucinations in diffusion models through mode interpolation. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=aNTnHBkw4T>.  
489  
490  
491
- 492 Marianne Arriola, Subham Sekhar Sahoo, Aaron Gokaslan, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Justin T Chiu, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=tyEyYT267x>.  
493  
494  
495
- 496 Henning Bartsch, Ole Jorgensen, Domenic Rosati, Jason Hoelscher-Obermaier, and Jacob Pfau. Self-consistency of large language models under ambiguity. In Yonatan Belinkov, Sophie Hao, Jaap Jumelet, Najoung Kim, Arya McCarthy, and Hosein Mohebbi (eds.), *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pp. 89–105, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.blackboxnlp-1.7. URL <https://aclanthology.org/2023.blackboxnlp-1.7/>.  
497  
498  
499  
500  
501
- 502 Cosmin I Bercea, Michael Neumayr, Daniel Rueckert, and Julia A Schnabel. Mask, stitch, and re-sample: Enhancing robustness and generalizability in anomaly detection through automatic diffusion models. *arXiv preprint arXiv:2305.19643*, 2023.  
503  
504  
505
- 506 James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafull Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. Improving image generation with better captions. OpenAI Technical Report, 2023. URL <https://cdn.openai.com/papers/dall-e-3.pdf>.  
507  
508  
509  
510
- 511 Kevin Clark and Priyank Jaini. Text-to-image diffusion models are zero shot classifiers. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 58921–58937. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/b87bdcf963cad3d0b265fcb78ae7d11e-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/b87bdcf963cad3d0b265fcb78ae7d11e-Paper-Conference.pdf).  
512  
513  
514  
515
- 516 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.  
517  
518
- 519 Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=AAWuCVzaVt>.  
520  
521  
522  
523
- 524 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 *Forty-first international conference on machine learning*, 2024.  
525  
526  
527
- 528 Guilherme Fernandes, Vasco Ramos, Regev Cohen, Idan Szpektor, and João Magalhães. Latent beam diffusion models for generating visual sequences, 2025. URL <https://arxiv.org/abs/2503.20429>.  
529  
530  
531
- 532 Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36: 52132–52152, 2023.  
533  
534
- 535 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations, 2019. URL <https://arxiv.org/abs/1903.12261>.  
536  
537
- 538 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15262–15271, June 2021.  
539

- 540 Alvin Heng, Alexandre H. Thiery, and Harold Soh. Out-of-distribution detection with  
541 a single unconditional diffusion model. In A. Globerson, L. Mackey, D. Bel-  
542 grave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural In-*  
543 *formation Processing Systems*, volume 37, pp. 43952–43974. Curran Associates, Inc.,  
544 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/4dc37a7bc61057252ce043fa3b83aac2-Paper-Conference.pdf)  
545 [file/4dc37a7bc61057252ce043fa3b83aac2-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/4dc37a7bc61057252ce043fa3b83aac2-Paper-Conference.pdf).
- 546 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A  
547 reference-free evaluation metric for image captioning, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2104.08718)  
548 [abs/2104.08718](https://arxiv.org/abs/2104.08718).
- 549 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochre-  
550 iter. Gans trained by a two time-scale update rule converge to a local nash equilibrium.  
551 In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and  
552 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran  
553 Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/](https://proceedings.neurips.cc/paper_files/paper/2017/file/8ald694707eb0fefe65871369074926d-Paper.pdf)  
554 [paper/2017/file/8ald694707eb0fefe65871369074926d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/8ald694707eb0fefe65871369074926d-Paper.pdf).
- 555 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*  
556 *arXiv:2207.12598*, 2022.
- 557 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
558 *neural information processing systems*, 33:6840–6851, 2020.
- 559 Yi Huang, Jiancheng Huang, Yifan Liu, Mingfu Yan, Jiayi Lv, Jianzhuang Liu, Wei Xiong,  
560 He Zhang, Liangliang Cao, and Shifeng Chen. Diffusion model-based image editing: A sur-  
561 vey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(6):4409–4437, 2025.  
562 doi: 10.1109/TPAMI.2025.3541625.
- 563 Tero Karras, Miika Aittala, Tuomas Kynkäänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine.  
564 Guiding a diffusion model with a bad version of itself. *Advances in Neural Information Processing*  
565 *Systems*, 37:52996–53021, 2024.
- 566 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
567 a-pic: An open dataset of user preferences for text-to-image generation. In *Thirty-seventh Confer-*  
568 *ence on Neural Information Processing Systems*, 2023. URL [https://openreview.net/](https://openreview.net/forum?id=G5RwHpBUv0)  
569 [forum?id=G5RwHpBUv0](https://openreview.net/forum?id=G5RwHpBUv0).
- 570 Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
571 language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*,  
572 volume 35, pp. 22199–22213, 2022.
- 573 Keon Lee, Dong Won Kim, Jaehyeon Kim, Seungjun Chung, and Jaewoong Cho. DiTTo-TTS: Dif-  
574 fusion transformers for scalable text-to-speech without domain-specific factors. In *The Thirteenth*  
575 *International Conference on Learning Representations*, 2025. URL [https://openreview.](https://openreview.net/forum?id=hQvX9MBowC)  
576 [net/forum?id=hQvX9MBowC](https://openreview.net/forum?id=hQvX9MBowC).
- 577 Alexander C. Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your dif-  
578 fusion model is secretly a zero-shot classifier. In *Proceedings of the IEEE/CVF International*  
579 *Conference on Computer Vision (ICCV)*, pp. 2206–2217, October 2023.
- 580 Shufan Li, Konstantinos Kallidromitis, Akash Gokul, Arsh Koneru, Yusuke Kato, Kazuki Kozuka,  
581 and Aditya Grover. Reflect-dit: Inference-time scaling for text-to-image diffusion transformers  
582 via in-context reflection, 2025. URL <https://arxiv.org/abs/2503.12271>.
- 583 Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. Diffusion-  
584 LM improves controllable text generation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,  
585 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL  
586 <https://openreview.net/forum?id=3s9IrEsjLyk>.
- 587 Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,  
588 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl, 2025a.  
589 URL <https://arxiv.org/abs/2505.05470>.
- 590  
591  
592  
593

- 594 Jing Liu, Zhenchao Ma, Zepu Wang, Chenxuanyin Zou, Jiayang Ren, Zehua Wang, Liang Song,  
595 Bo Hu, Yang Liu, and Victor Leung. A survey on diffusion models for anomaly detection. *arXiv*  
596 *preprint arXiv:2501.11430*, 2025b.  
597
- 598 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and  
599 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.  
600
- 601 Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. InstafLOW: One step is enough for  
602 high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference*  
603 *on Learning Representations*, 2023.
- 604 Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma, Chong Ruan, Peng Li, Yang Liu, and Yu Wu.  
605 Inference-time scaling for generalist reward modeling, 2025c. URL [https://arxiv.org/](https://arxiv.org/abs/2504.02495)  
606 [abs/2504.02495](https://arxiv.org/abs/2504.02495).  
607
- 608 Victor Livernoche, Vineet Jain, Yashar Hezaveh, and Siamak Ravanbakhsh. On diffusion modeling  
609 for anomaly detection. *arXiv preprint arXiv:2305.18593*, 2023.  
610
- 611 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast  
612 ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in neural*  
613 *information processing systems*, 35:5775–5787, 2022.
- 614 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthe-  
615 sizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.  
616
- 617 Weijian Luo. A comprehensive survey on knowledge distillation of diffusion models. *arXiv preprint*  
618 *arXiv:2304.04262*, 2023.  
619
- 620 Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang,  
621 Yandong Li, Tommi Jaakkola, Xuhui Jia, and Saining Xie. Scaling inference time compute for  
622 diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
623 *Recognition (CVPR)*, pp. 2523–2534, June 2025.  
624
- 625 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for  
626 editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference*  
627 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 6038–6047, June 2023.
- 628 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke  
629 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time  
630 scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.  
631
- 632 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
633 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.  
634
- 635 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
636 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
637 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 638 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
639 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*  
640 *in neural information processing systems*, 36:53728–53741, 2023.  
641
- 642 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
643 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
644 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.  
645
- 646 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-  
647 ical image segmentation. In *International Conference on Medical image computing and computer-*  
*assisted intervention*, pp. 234–241. Springer, 2015.

- 648 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-  
649 yar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan  
650 Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion mod-  
651 els with deep language understanding. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,  
652 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL  
653 <https://openreview.net/forum?id=08Yk-n5l2Al>.
- 654 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and  
655 Xi Chen. Improved techniques for training gans. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon,  
656 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Cur-  
657 ran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper\\_files/  
658 paper/2016/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf).
- 659 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,  
660 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathe-  
661 matical reasoning in open language models, 2024. URL [https://arxiv.org/abs/2402.  
662 03300](https://arxiv.org/abs/2402.03300).
- 663 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally  
664 can be more effective than scaling model parameters, 2024. URL [https://arxiv.org/  
665 abs/2408.03314](https://arxiv.org/abs/2408.03314).
- 666 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
667 learning using nonequilibrium thermodynamics. In *International conference on machine learn-  
668 ing*, pp. 2256–2265. pmlr, 2015.
- 669 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv  
670 preprint arXiv:2010.02502*, 2020.
- 671 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. 2023.
- 672 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,  
673 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using  
674 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
675 and Pattern Recognition*, pp. 8228–8238, 2024.
- 676 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.  
677 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-  
678 image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- 679 Enze Xie, Junsong Chen, Yuyang Zhao, Jincheng YU, Ligeng Zhu, Yujun Lin, Zhekai Zhang,  
680 Muyang Li, Junyu Chen, Han Cai, Bingchen Liu, Daquan Zhou, and Song Han. SANA  
681 1.5: Efficient scaling of training-time and inference-time compute in linear diffusion trans-  
682 former. In *Forty-second International Conference on Machine Learning*, 2025. URL [https:  
683 //openreview.net/forum?id=27hOkXzy9e](https://openreview.net/forum?id=27hOkXzy9e).
- 684 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao  
685 Dong. Imagereward: learning and evaluating human preferences for text-to-image generation. In  
686 *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pp.  
687 15903–15935, 2023.
- 688 Hang Yao, Ming Liu, Zhicun Yin, Zifei Yan, Xiaopeng Hong, and Wangmeng Zuo. Glad: Towards  
689 better reconstruction with global and local adaptive diffusion models for unsupervised anomaly  
690 detection. In *European Conference on Computer Vision*, pp. 1–17. Springer, 2024.
- 691 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik  
692 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In  
693 A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in  
694 Neural Information Processing Systems*, volume 36, pp. 11809–11822. Curran Associates, Inc.,  
695 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/  
696 file/271db9922b8d1f4dd7aaef84ed5ac703-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/271db9922b8d1f4dd7aaef84ed5ac703-Paper-Conference.pdf).

Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop, 2016. URL <https://arxiv.org/abs/1506.03365>.

Robert B Zajonc. Mere exposure: A gateway to the subliminal. *Current directions in psychological science*, 10(6):224–228, 2001.

Hui Zhang, Zheng Wang, Dan Zeng, Zuxuan Wu, and Yu-Gang Jiang. Diffusionad: Norm-guided one-step denoising diffusion for anomaly detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.

## A APPENDIX

### A.1 LLM USAGE

We utilize LLMs, Gemini-2.5-pro and ChatGPT, to polish our writing and to get help experimental code implementation.

### A.2 MORE RESULTS FROM ABLATION STUDY

The results from the ablation study in Sec.5.2 on PickScore (Kirstain et al., 2023) from table 6 and ImageReward (Xu et al., 2023) from table. 7. The findings show a similar trend to the results from HPS (Table.5). Consequently, we selected the "4 candidates from first 4 timesteps" approach. For the text-to-image human preference alignment experiments, we also included the results for "8 candidates for only the first timestep".

| PickScore           |    | Candidates         |       |       |       |       |
|---------------------|----|--------------------|-------|-------|-------|-------|
|                     |    | 1                  | 2     | 4     | 8     | 16    |
| Timesteps           | 1  | 22.02<br>(vanilla) | 22.04 | 22.07 | 22.08 | 22.03 |
|                     | 2  | -                  | 22.04 | 22.06 | 22.08 | 22.07 |
|                     | 4  | -                  | 22.02 | 22.12 | 22.06 | 22.06 |
|                     | 8  | -                  | 22.04 | 22.04 | 22.10 | 22.09 |
|                     | 16 | -                  | 22.04 | 22.07 | 22.06 | 22.13 |
| Constant Candidates |    |                    | 22.05 |       |       |       |
| Linear Pruning      |    |                    | 22.05 |       |       |       |
| Exponential Pruning |    |                    | 22.07 |       |       |       |

Table 6: Ablation study on candidate scheduling, report PickScore score. Constant candidates, Linear Pruning, and Exponential Pruning are described at Sec. 5.1

| ImageReward         |    | Candidates          |        |        |        |        |
|---------------------|----|---------------------|--------|--------|--------|--------|
|                     |    | 1                   | 2      | 4      | 8      | 16     |
| Timesteps           | 1  | 0.8361<br>(vanilla) | 0.8543 | 0.8707 | 0.8929 | 0.8985 |
|                     | 2  | -                   | 0.8490 | 0.8389 | 0.8706 | 0.8600 |
|                     | 4  | -                   | 0.8356 | 0.8972 | 0.8681 | 0.8541 |
|                     | 8  | -                   | 0.8253 | 0.9201 | 0.8660 | 0.9024 |
|                     | 16 | -                   | 0.8619 | 0.8619 | 0.8516 | 0.8795 |
| Constant Candidates |    |                     | 0.8724 |        |        |        |
| Linear Pruning      |    |                     | 0.8746 |        |        |        |
| Exponential Pruning |    |                     | 0.8908 |        |        |        |

Table 7: Ablation study on candidate scheduling, report ImageReward score. Constant candidates, Linear Pruning, and Exponential Pruning are described at Sec. 5.1

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

### A.3 SAMPLES FROM UNCONDITIONAL GENERATION

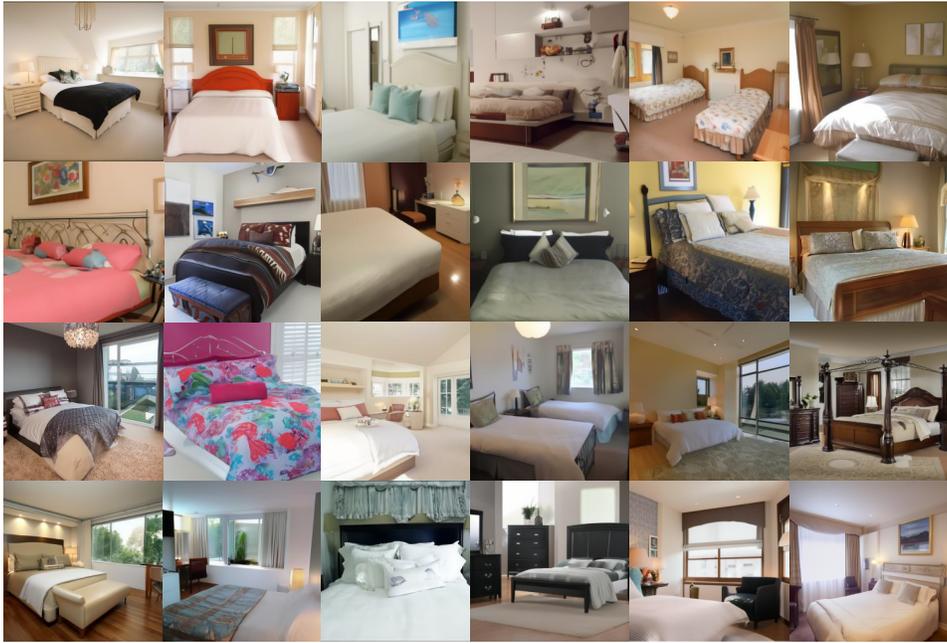


Figure 4: Unconditionally generated samples from model trained on LSUN-Bedroom 256x256 with our sampling methods. FID = 8.45

