# Continuous Speculative Decoding for Autoregressive Image Generation

Zili Wang[1,2]    Robert Zhang[3]    Kun Ding[1,2]    Qi Yang[1,2]    Fei Li[4]    Shiming Xiang[1,2]

[1]University of Chinese Academy of Sciences, China

[2]Institute of Automation, Chinese Academy of Sciences, China

[3]Independent Researcher    [4]China Tower Corporation Limited

Figure 1. Continuous speculative decoding accelerates the inference speed while maintaining the original generation quality.

## Abstract

*Continuous-valued Autoregressive (AR) image generation models have demonstrated notable superiority over their discrete-token counterparts, showcasing considerable reconstruction quality and higher generation fidelity. However, the computational demands of the autoregressive framework result in significant inference overhead. While speculative decoding has proven effective in accelerating Large Language Models (LLMs), their adaptation to continuous-valued visual autoregressive models remains unexplored. This work generalizes the speculative decoding algorithm from discrete tokens to continuous space. By analyzing the intrinsic properties of output distribution, we establish a tailored acceptance criterion for the diffusion distributions prevalent in such models. To overcome the inconsistency that occurred in speculative decoding output distributions, we introduce denoising trajectory alignment and token pre-filling methods. Additionally, we identify the hard-to-sample distribution in the rejection phase. To mitigate this issue, we propose a meticulous acceptance-rejection sampling method with a proper upper bound, thereby circumventing complex integration. Experimental results show that our continuous speculative decoding achieves a remarkable 2.33× speed-up on off-the-shelf models while maintaining the output distribution. Codes will be available at: https://github.com/MarkXCloud/CSpD*

## 1. Introduction

Autoregressive (AR) models have shown significant superiority in image generation tasks [4, 7, 10, 20, 30, 36, 38, 39, 45]. They predict next token sequentially based on previously generated tokens. Typically, input images are transformed from pixel space into a discrete token space through vector quantization (VQ), and then AR models predict the next token as a classification task. This approach has demonstrated considerable potential in image generation. However, the VQ operation induces instability during training and may be insufficient for capturing the nuanced image details [26, 45]. Recently, an alternative approach utilizes diffusion models [15, 28] to embed visual tokens into a continuous distribution [21, 37, 41]. The next token prediction is performed through a denoising process conditioned by the outputs of the autoregressive model. They not only mitigate the issues associated with discrete vector quantization but also facilitate the generation of higher-quality images [11, 13].

However, the inference of the autoregressive model is slow and expensive due to sequential decoding. Speculative decoding [5, 19] in LLMs aims to reduce the inference cost via draft & verification. It involves a smaller draft model generating a sequence of draft tokens $p(x)$. Then, a more accurate target model $q(x)$ would verify each token to decide whether to accept or reject it and then resam-
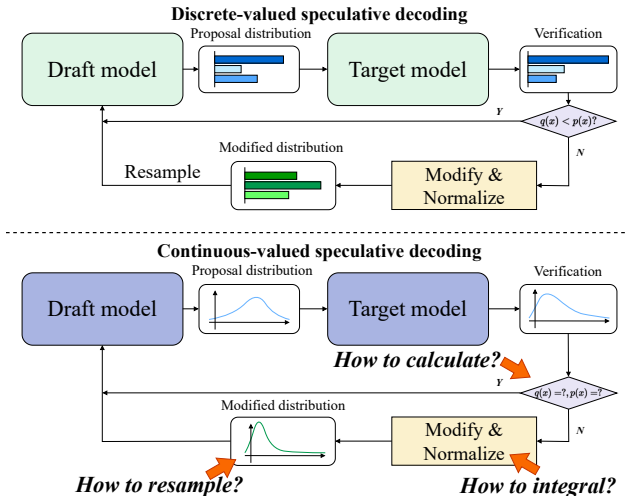
1

Figure 2. Comparison between discrete- and continuous-valued speculative decoding. Discrete models can conveniently compute output probabilities and be sampled from modified distributions. In contrast, continuous models require determining how to compute probabilities, and sampling from modified distributions via draft and target output distributions is often more challenging.

ple a new one. Nonetheless, speculative decoding has not been adapted to continuous-valued autoregressive models. Recent works [16, 35] are introduced to the discrete tokens, but the application to the continuous probability density function (PDF) is not feasible. As shown in Figure 2, it introduces the following challenges. a) In the continuous situation, we need to establish how to calculate the PDF of the draft and target distribution $q(x)$ and $p(x)$, and then calculate the acceptance criterion for verification, termed as $p(x)/q(x)$. Furthermore, distinct inconsistency exists between the draft and target outputs, leading to a low acceptance rate during speculative decoding. b) In the resampling phase after rejection, we must sample from the modified distribution to generate a new token. However, this distribution involves complex integration, which makes it challenging to derive an analytical form and to sample directly.

To this end, we introduce **Continuous Speculative Decoding**, the first to extend the speculative decoding algorithm to the continuous token space. To establish the acceptance criterion $p(x)/q(x)$, we derive an appropriate calculation method to obtain the two PDFs through a comprehensive analysis of the diffusion process at the output end. To enhance output consistency, we introduce a denoising trajectory alignment to align the output distributions of the draft and target models. Also, for the inconsistency in the autoregressive process, we introduce a token pre-filling strategy to mitigate the issue of low early acceptance rates and enhance the overall acceptance rate without compromising inference speed. For the rejection and resampling phase, we employ acceptance-rejection sampling [2] to the

modified PDF, which lacks an analytical form. By setting proper upper bounds, we successfully avoid the complex integration and simplify the sampling of the modified distribution. Similar to traditional speculative decoding, under these configurations, our method ensures that the output distribution of the speculative decoding process remains consistent with that of the target model alone.

Our continuous speculative decoding framework can be integrated seamlessly with existing continuous-valued visual autoregressive models without additional training procedures, alterations to the model architectures, or changes to the output distribution. Extensive experiments demonstrate the effectiveness of our algorithm via qualitative and quantitative evaluations. Furthermore, we measure and compare the wall-time improvements on open-sourced MAR [21] models with various configurations. We also report the performance metrics using Fréchet Inception Distance (FID) [14] and Inception Score (IS) [31] on ImageNet $256 \times 256$ generation. As shown in Figure 1, our continuous speculative decoding algorithm achieves an outstanding inference speedup by up to $2.33\times$ while maintaining the original generation quality to a large extent.

Our contributions can be summarized as follows:
- We are the first to extend speculative decoding to continuous space and adapt it for continuous-valued autoregressive image generation.
- We establish an appropriate acceptance criterion for continuous PDF and an integral-free sampling method through acceptance-rejection sampling of the modified distribution, which lacks an analytical form.
- We derive denoising trajectory alignment to align the distributions and a token pre-filling method to address the issue of low initial step acceptance rates.
- Without additional training, modifications on model architecture, or changes of output distribution, our method achieves $2.33\times$ speedup and comparable image quality.

## 2. Related Work

### 2.1. Autoregressive Image Generation

Autoregressive (AR) model is employed for next-token prediction and has important applications in image generation. Early works [6, 7, 9, 12, 29, 38, 39] perform autoregressive image prediction on pixel sequence level with CNN [7, 38], RNN [39] and Transformer [6, 9, 29]. Inspired by LLMs, subsequent studies [10, 30, 40, 43] regard image generation as discrete token classification. Based on this, mask generative models [3, 4, 20] train generation models via random masking. VAR [36] modifies the autoregressive paradigm into next-scale prediction to gradually increase the scale of predictions. Similar to the autoregressive language model, AR image generation through discrete token prediction is scalable to text-conditioned im-
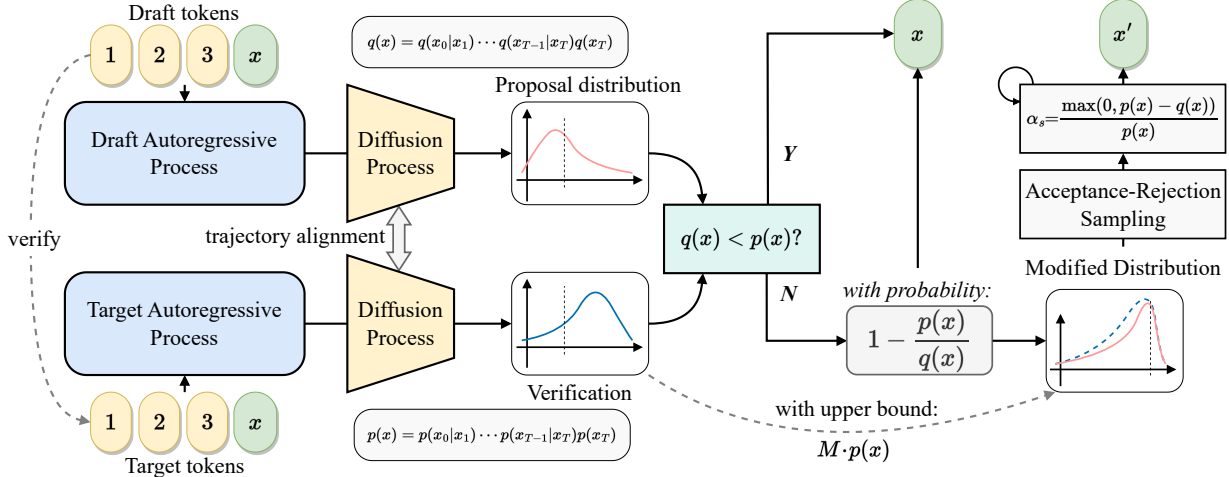
Figure 3. The overview of our proposed continuous speculative decoding. Continuous speculative decoding leverages the diffusion model component of continuous AR models. Tokens $1 \sim 3$ are prefix tokens, and token $x$ is to be verified. Upon obtaining and comparing the probability density values from the draft and target model, if $q(x) < p(x)$, $x$ is accepted. Otherwise, $x$ is rejected with probability $1 - \frac{p(x)}{q(x)}$, followed by sampling from the modified distribution via acceptance-rejection sampling to obtain $x'$.

age generation [24, 33, 44, 46]. However, training discrete image tokenizer is difficult, and its ability to convey detailed visuals may be limited [26, 45]. GIVT [37] represents continuous tokens via Gaussian mixture models. MAR [21] and DisCo-Diff [41] generate tokens via diffusion process [15, 28] conditioned by the autoregressive model. Visual AR models with diffusion process also show scalability to text-conditioned image generation [11]. DART [13] performs autoregressive steps within the whole image denoising process. HART [34] employs discrete and continuous tokenizer to generate images, with classification for discrete tokens and denoising for the residual between primitive visual tokens and discrete tokens. RAR [47] employs a randomized permutation objective to improve bidirectional context learning. However, autoregressive models suffer from heavy inference overhead. The inference speed is slowed down by step-by-step generation.

## 2.2. Speculative Decoding

Speculative decoding [5, 19] achieves lossless acceleration by verifying the draft model with the target model. Following this, previous works mainly focus on reducing draft model overhead and strengthening the consistency between the draft and target models. SpecInfer [27] employs multiple small draft models and aggregates their predictions into a tree structure to be verified through tree-based parallel decoding. Medusa [1] predicts multiple drafts by using the original LLM's features and training another set of classification heads. Then, these drafts are verified by tree attention. Eagle [22, 23] improves the draft accuracy through the prediction at the feature level instead of the token level

to tackle the feature uncertainty problem. Jacobi iteration is also employed to reduce inference overhead in the decoding process [18, 32, 48, 49]. Online Speculative Decoding [25] and DistillSpec [50] align the output from the draft model with the target model with more training. Speculative decoding can achieve lossless acceleration theoretically, but the generation quality may be affected under a larger speed-up ratio. BiLD [17] proposes a more relaxed acceptance condition. Besides, finetuning the target model can also improve generation quality [1, 18, 42].

Speculative decoding can be applied directly to visual AR models that use discrete tokens. SJD [35] improves the Jacobi iteration process by adding speculative decoding while keeping the variety of image generation. LANTERN [16] looks at distribution ambiguity and uses relaxation to add more flexible candidate tokens, maintaining high image quality. These studies have greatly helped speed up visual autoregressive models. However, they mainly focus on discrete tokens and have not explored continuous space. In contrast, our work demonstrates speculative decoding in AR models in continuous output space.

## 3. Methodology

### 3.1. Preliminaries

In LLM practice, speculative decoding utilizes a draft model $M_q$ with output distribution $q(x)$ and a target model $M_p$ with output distribution $p(x)$. Firstly, the draft model predicts draft tokens by sampling $x \sim q(x)$. Then, the target model verifies these tokens to output $p(x)$ in parallel. If $\frac{p(x)}{q(x)} > 1$, the draft token is accepted. Otherwise, we

reject the token with probability $1 - \frac{p(x)}{q(x)}$ and resample $x$ from a modified distribution $p'(x) = norm(max(0, p(x) - q(x))) = \frac{max(0, p(x) - q(x))}{\sum_{x'} max(0, p(x') - q(x'))}$. Through this procedure, speculative decoding preserves the original output distribution $x \sim p(x)$ as the target model.

In previous works of speculative decoding, the outputs of the models follow discrete distributions. However, the continuous-valued visual AR model, like MAR [21], generates tokens by sampling from a continuous distribution. This difference prevents speculative decoding from being directly applied to it. In this section, we will elaborate on speculative decoding within the context of the continuous-valued visual AR model, as shown in Figure 3. More detailed proofs and correctness of continuous speculative decoding can be found in supplemental material.

### 3.2. How to determine the acceptance criterion by draft and target distributions?

Unlike discrete situations, $p(x)$ in continuous space represents the probability density function (PDF) of $x$. Our goal is to maintain the same distribution as PDF $p(x)$ through our algorithm. Therefore, the ratio of the PDFs from draft and target output distributions can be directly obtained as a criterion of acceptance. That is, if $\frac{p(x)}{q(x)} > 1$, the draft token is accepted. Otherwise, we reject the token with probability $1 - \frac{p(x)}{q(x)}$. The next step is to calculate $\frac{p(x)}{q(x)}$.

Since the output distribution of the continuous autoregressive model is usually obtained via diffusion process determined by DDPM [15], to sample a token $x = x_0$, we sample $p(x_0|x_T)$ instead via reverse diffusion (denoising) process, where $x_T$ is a Gaussian noise. A nice property of the denoising process is that the sampling process is a Markov chain. We can calculate $\frac{p(x)}{q(x)}$ through:

$$\frac{p(x_0|x_T)}{q(x_0|x_T)} = \frac{p(x_T) \prod_{t=1}^{T} p(x_{t-1}|x_t)}{q(x_T) \prod_{t=1}^{T} q(x_{t-1}|x_t)}, \quad (1)$$

with the conditioned probability distributions $p(x_{t-1}|x_t)$ approximated as Gaussian by a neural network $\theta$ [28]:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)), \quad (2)$$

where $\mu_\theta$ and $\Sigma_\theta$ are mean and variance predicted by $\theta$. However, significant inconsistency exists in denoising processes due to the differences between draft and target models. As shown in Figure 4, the denoising trajectories of the draft model and the target model diverge to different outputs, leading to distinctly different final generated distributions—the lack of consistency results in an exceedingly low acceptance probability.

We propose a denoising trajectory alignment strategy to enhance the consistency of the distributions generated by the denoising trajectories. The alignment ensures that
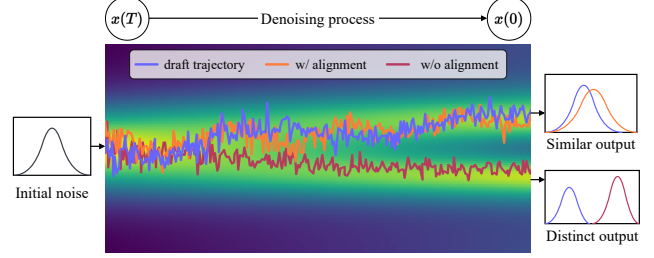


Figure 4. Illustration of denoising trajectory alignment. The denoising process maps the noise distribution to data distribution through gradual denoising. These denoising steps generate a trajectory. Aligned trajectories lead to similar output distribution, while unaligned one produces a more distinct distribution.

both models utilize the same sampling point from the standard Gaussian distribution $\varepsilon_t \sim \mathcal{N}(0, I)$ at each denoising step and reparameterizes each $x_{t-1} = \sqrt{\Sigma_\theta(x_t, t)} \cdot \varepsilon_t + \mu_\theta(x_t, t)$, thereby aligning the final output distributions towards a consistent trend. Through the property of Gaussian distribution reparameterization, we have:

$$p(x_{t-1}|x_t) = \frac{1}{\sqrt{|\Sigma_\theta(x_t, t)|}} p(\varepsilon_t). \quad (3)$$

Based on this, the alignment can simplify subsequent computations. Let $\Sigma_t^q$ and $\Sigma_t^p$ be the draft and target model variance at timestep $t$. Note that $p_\theta(x_{t-1}^p|x_t^p)/q_\theta(x_{t-1}^q|x_t^q) = \sqrt{|\Sigma_t^q|}/\sqrt{|\Sigma_t^p|}$ and $p(x_T) = q(x_T)$, thus we define:

$$\Sigma = \prod_{t=2}^{T} \sqrt{|\Sigma_t^q|} / \prod_{t=2}^{T} \sqrt{|\Sigma_t^p|}. \quad (4)$$

Substituting into Equation (1) yields:

$$\frac{p(x_0|x_T)}{q(x_0|x_T)} = \Sigma \cdot \frac{p_\theta(x_0|x_1^p)}{q_\theta(x_0|x_1^q)}. \quad (5)$$

Inconsistency also exists in autoregressive steps. The tokens generated by the draft model during the initial autoregressive steps have a lower acceptance rate, which will be discussed in the ablation study. This phenomenon can be attributed to the different prefix embedding. MAR models of different sizes have learnable prefix embeddings of their own [21]. Naturally, these models' autoregressive predictions diverge because they have different prefix conditions.

To address this issue, we propose pre-filling a portion (say 5%) of the tokens from the target model during autoregressive generation to obtain a relatively consistent prefix. The pre-filling does not compromise inference speed, as speculative decoding at a low acceptance rate is functionally equivalent to the step-by-step decoding performed by the target model alone [19]. Furthermore, pre-filling establishes a better starting point for autoregressive generation, thereby improving the overall acceptance rate.

Finally, we can calculate $p(x)/q(x)$ using the Gaussian distribution of the last denoising step and the cumulative product of variances along previous steps.

### 3.3. How to resample a new token from the modified distribution after $x$ is rejected?

When the sample from the draft model is rejected, we re-sample a new token from the modified distribution $p'(x) = norm(max(0, p(x) - q(x)))$ in LLMs. This distribution is easy to sample because it is discrete. However, the modified distribution is hard to obtain in continuous space, and the sampling operation is difficult. Note that $p'(x) = \frac{max(0, p(x) - q(x))}{\sum_{x'} max(0, p(x') - q(x'))}$, we convert this formula in continuous-valued PDF to obtain:

$$p'(x) = \frac{max(0, p(x) - q(x))}{\int_{x'} max(0, p(x') - q(x'))dx'}. \quad (6)$$

Equation (6) represents the normalized distribution of the subtraction between the target and draft model, with only the positive part remaining. The functional curve of this distribution is illustrated in Figure 5. From this equation, we cannot compute the exact expression of this distribution because of the integration of the normalizing factor, and therefore, it cannot be directly sampled.

To derive a method for sampling from Equation (6), we employ acceptance-rejection sampling [2] on $p'(x)$. Specifically, we sample from the target model distribution $p(x)$, which is easy to sample as a proposal distribution. Then, we calculate the threshold $\alpha_s$:

$$\alpha_s = \frac{p'(x)}{M \cdot p(x)}, \quad (7)$$

where $M$ is the upper bound factor that holds $M \cdot p(x) \geq p'(x)$ for any $x$. $M$ ensures the correctness of acceptance-rejection sampling [2].

Afterward, we sample $\epsilon \sim U(0, 1)$. We accept the sample from proposal distribution if $\epsilon < \alpha_s$. Otherwise, we reject it and resample another $x$. This sampling approach is equivalent to sampling from the modified distribution [2].

Note that the upper bound $M$ affects the sampling quality and should be determined properly. Given $Z = \int_x max(0, p(x) - q(x))dx$, consider that:

$$p'(x) = \frac{max(0, p(x) - q(x))}{Z}. \quad (8)$$

Given $p(x) - q(x) \geq 0$, we have:

$$\frac{p(x) - q(x)}{Z} \leq \frac{p(x)}{Z} \mapsto M \cdot p(x). \quad (9)$$

Therefore, we can set $M := 1/Z$ as the upper bound. However, the calculation of $Z$ still requires integration of

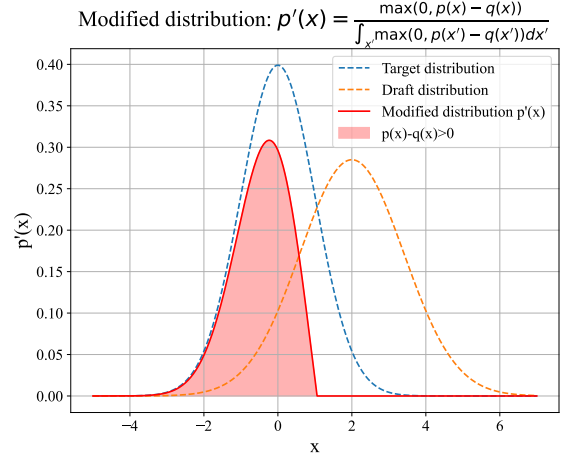Modified distribution: $p'(x) = \frac{max(0, p(x) - q(x))}{\int_{x'} max(0, p(x') - q(x'))dx'}$



Figure 5. Illustration of the modified distribution (unnormalized), where the dashed lines represent the output distributions of the draft and target models, and the red area denotes the modified distribution. The integral of this area is hard to compute, and there is no analytical expression available, which complicates sampling.

$max(0, p(x) - q(x))$. Performing integration over the entire space introduces additional complex computations and potential errors. To this end, we substitute $M = 1/Z$ into Equation (7) to obtain:

$$\alpha_s = \frac{max(0, p(x) - q(x))/Z}{p(x)/Z} = \frac{max(0, p(x) - q(x))}{p(x)}. \quad (10)$$

Equation (10) reveals that with $M = 1/Z$, we can calculate $\alpha_s$ directly by the two PDFs. Calculating the integral $Z$ can be avoided, which reduces computational complexity while maintaining the correctness of the results.

To calculate $\alpha_s$, we substitute Equation (5) to obtain:

$$\alpha_s = \frac{max(0, \Sigma \cdot p_\theta(x_0|x_1^p) - q_\theta(x_0|x_1^q))}{\Sigma \cdot p_\theta(x_0|x_1^p)}. \quad (11)$$

Therefore, we can resample from the target distribution, calculate $\alpha_s$, and reject with $\epsilon \sim U(0, 1)$ to simulate sampling from the original modified distribution.

## 4. Experiment

### 4.1. Implementation Details

We systematically conduct experiments with open-sourced continuous-valued visual autoregressive model MAR [21] on ImageNet [8] $256 \times 256$ generation. Due to the limited open-sourced models available, we evaluate MAR-B (208M) and MAR-L (479M) as draft models, as well as MAR-L and MAR-H (943M) as target models. We use official pretrained checkpoints. We report FID [14] and IS [31] as well as wall-time speedup ratio on a single NVIDIA A100 GPU. The runtime settings strictly conform to the prescribed configurations [21]. More ablation studies and quantitative results can be found in supplemental material.
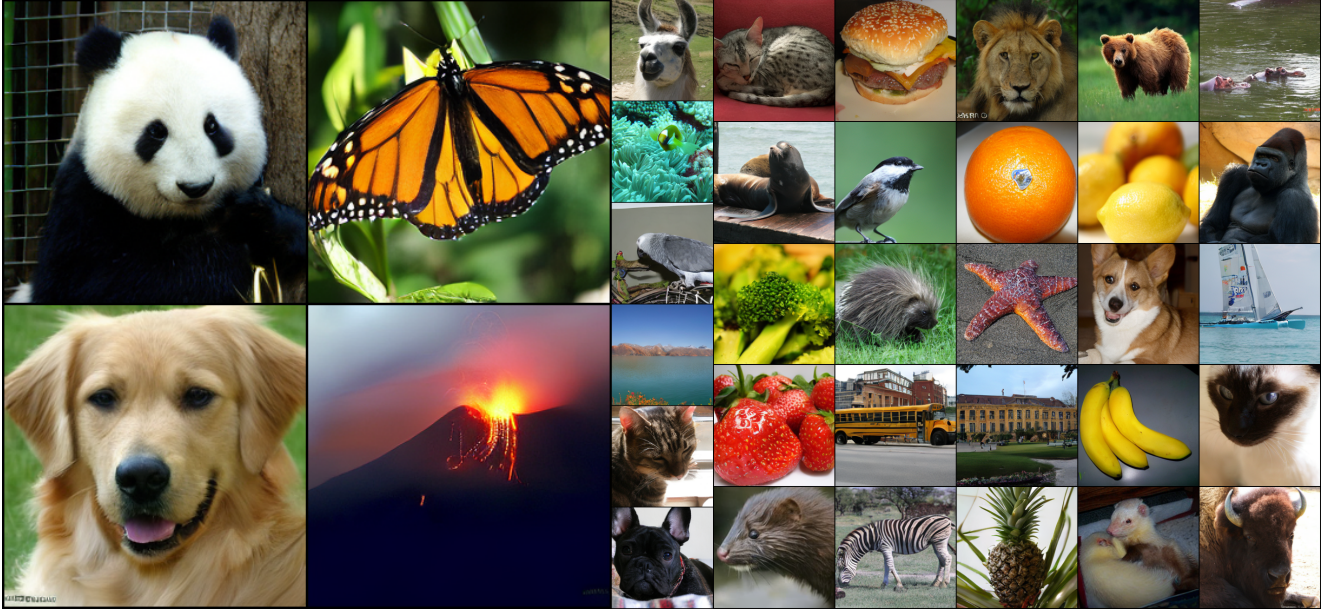
Figure 6. Qualitative Results. We show the images generated under continuous speculative decoding with MAR.

| $M_p$ | $M_q$ | $\gamma$ | $\alpha$ | Speedup ratio | | | |
|---|---|---|---|---|---|---|---|
| | | | | bs=1 | bs=8 | bs=128 | bs=256 |
| MAR-L | MAR-B | 32 | 0.26 | **1.18×** | **1.21×** | **1.44×** | **1.49×** |
| MAR-L | MAR-B | 16 | 0.31 | 1.10× | 1.17× | 1.39× | 1.42× |
| MAR-L | MAR-B | 8 | 0.36 | 1.05× | 1.12× | 1.29× | 1.32× |
| MAR-L | MAR-B | 4 | 0.39 | 1.01× | 1.00× | 1.13× | 1.15× |
| MAR-H | MAR-B | 32 | 0.19 | **1.44×** | **1.61×** | **2.17×** | **2.33×** |
| MAR-H | MAR-L | 32 | 0.18 | 1.26× | 1.34× | 1.47× | 1.53× |
| MAR-H | MAR-B | 16 | 0.26 | 1.37× | 1.51× | 2.07× | 2.20× |
| MAR-H | MAR-L | 16 | 0.24 | 1.24× | 1.29× | 1.41× | 1.46× |
| MAR-H | MAR-B | 8 | 0.27 | 1.26× | 1.44× | 1.88× | 1.96× |
| MAR-H | MAR-L | 8 | 0.28 | 1.11× | 1.21× | 1.32× | 1.33× |
| MAR-H | MAR-B | 4 | 0.30 | 1.11× | 1.20× | 1.56× | 1.62× |
| MAR-H | MAR-L | 4 | 0.30 | 1.00× | 1.03× | 1.15× | 1.18× |

Table 1. Results of speedup ratio on MAR [21] under different model size, draft number and batch size. The bs refers to batch size. The acceptance rate $\alpha$ of each setting is also represented.

## 4.2. Main Results

**Speedup results.** Table 1 shows the speedup ratio under different batch sizes, with draft numbers ranging from 8 to 32, along with the overall acceptance rate. As the batch size grows, the efficacy of speculative decoding becomes evident. Compared with the original model, our approach achieves impressive up to 2.33× acceleration at most. Due to the limited open-sourced models, the current disparity

| $M_p$ | $M_q$ | w/o CFG | | w/ CFG | |
|---|---|---|---|---|---|
| | | FID↓ | IS↑ | FID↓ | IS↑ |
| MAR-L | | 2.60 | 221.4 | 1.78 | 296.0 |
| MAR-L | MAR-B | 2.59 ±0.04 | 218.4±3.4 | 1.81±0.05 | 303.7±4.3 |
| MAR-H | | 2.35 | 227.8 | 1.55 | 303.7 |
| MAR-H | MAR-B | 2.36±0.05 | 228.5±2.2 | 1.60±0.05 | 301.6±2.6 |
| MAR-H | MAR-L | 2.34±0.04 | 228.9±2.8 | 1.57±0.04 | 301.4±2.5 |

Table 2. Evaluation of FID and IS comparison on ImageNet $256 \times 256$ unconditional and conditional generation. Continuous speculative decoding achieves acceleration while maintaining performance within a reasonable interval.

between draft and target models in scale is not extensive (208M vs. 943M). As the scale disparity increases, the acceleration effect is expected to become more pronounced, suggesting more potential for further enhancement.

**Quantitative results.** The class-conditioned and unconditioned FID and IS metrics for our continuous speculative decoding and the original MAR models are shown in Table 2. Through multiple repeated experiments, the evaluation performance consistently falls within an acceptable range. These results confirm the validity of our theoretical analysis and demonstrate that our algorithm effectively maintains the target output distribution. Consequently, the quality of the generated outputs is preserved, which will be further discussed in subsequent sections. Furthermore, these characteristics establish our algorithm as a robust solution for efficient and reliable model inference.
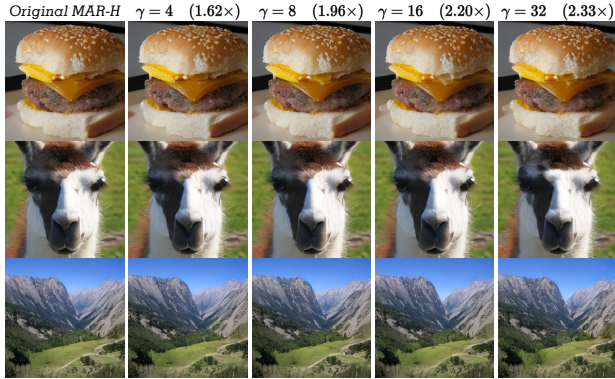
Figure 7. Qualitative Comparison Results. We show the generated images using the algorithm at various draft length $\gamma$.

**Qualitative results.** To demonstrate the quality of image generation by our method, we present more visualizations, as shown in Figure 6 and 7. Figure 6 demonstrates the primary generation results. Figure 7 showcases the results of the original MAR-H with autoregressive step 256 compared with various draft lengths $\gamma$. In addition to a significant acceleration up to $2.33\times$, our method has primarily maintained the quality of the generated images, which is consistent with our theoretical proof.
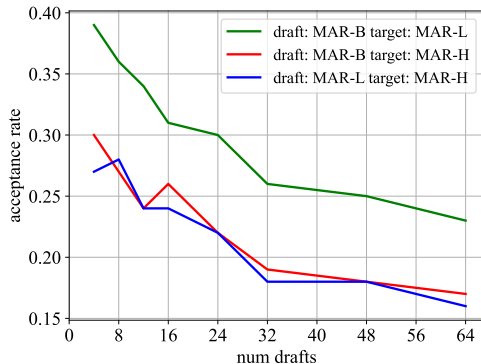


Figure 8. Acceptance ratio under different number of drafts. Larger number of drafts leads to the decay of acceptance ratio.

### 4.3. Ablation Study

**Effectiveness of Draft & Verification.** We present the comparative results of the generation process utilizing a pure draft model and the draft & verification paradigm, as illustrated in Figure 9. In the verification phase, regions within the draft that exhibit suboptimal token generation quality are systematically identified and substituted with tokens of higher quality. This methodology not only upholds the overall compositional integrity but also significantly enhances the richness and detail of the generated output.
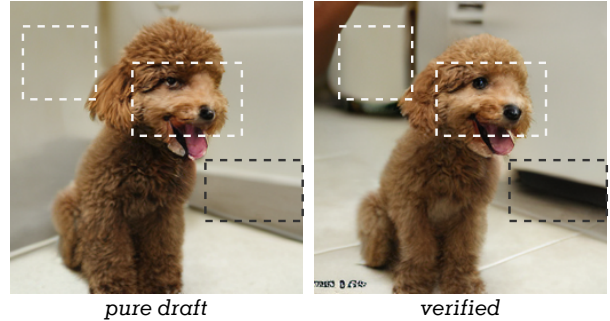


pure draft                                   verified

Figure 9. Comparison on pure draft (left) and verified (right) generation results. Regions of rejected tokens are roughly marked out.

| $M_p$ | $M_q$ | $\gamma$ | $\alpha$ | |
| --- | --- | --- | --- | --- |
| | | | w/o align | w/ align |
| MAR-L | MAR-B | 32 | 0.10 | **0.34** |
| MAR-L | MAR-B | 16 | 0.12 | **0.37** |
| MAR-L | MAR-B | 8 | 0.12 | **0.39** |
| MAR-L | MAR-B | 4 | 0.13 | **0.37** |
| MAR-H | MAR-B | 32 | 0.07 | **0.30** |
| MAR-H | MAR-L | 32 | 0.06 | **0.33** |
| MAR-H | MAR-B | 16 | 0.07 | **0.33** |
| MAR-H | MAR-L | 16 | 0.08 | **0.35** |
| MAR-H | MAR-B | 8 | 0.13 | **0.31** |
| MAR-H | MAR-L | 8 | 0.12 | **0.34** |
| MAR-H | MAR-B | 4 | 0.14 | **0.32** |
| MAR-H | MAR-L | 4 | 0.12 | **0.34** |

Table 3. Ablation study on the impact of acceptance rate with and without denoising trajectory alignment.

**The $\alpha$ vs. $\gamma$.** The relationship between acceptance rates and draft length is depicted in Figure 8. As the length of the draft increases, the acceptance rate tends to decline. This observation suggests that while longer drafts can substantially mitigate inference overhead, they are intrinsically constrained by the capabilities of the draft model itself. Consequently, an increase in the number of draft tokens is associated with greater deviations from the target model's distribution, ultimately leading to reduced acceptance rates.

**Effectiveness of trajectory alignment.** Table 3 shows the acceptance rate of the randomly sampled trajectory with our aligned trajectory. Denoising trajectory alignment enhances the consistency between the draft and target output distributions, thereby increasing the acceptance rate and simplifying the calculation of $p(x)/q(x)$. Figure 10 illustrates the image generated with and without alignment. The influence of alignment on the overall quality appears to be negligible. In contrast, utterly random sampling leads to a lower acceptance rate because of the distinct output distribution.

*w/o denoising trajectory alignment*
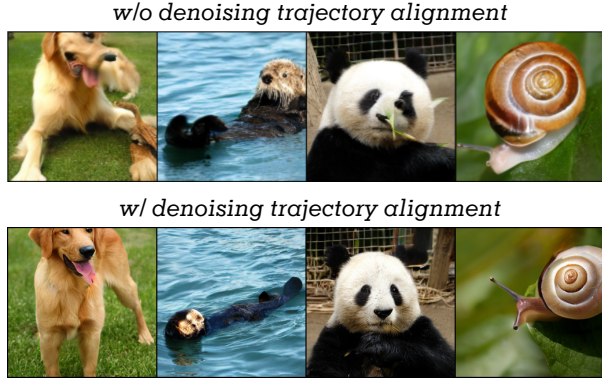
*w/ denoising trajectory alignment*

Figure 10. The examples without (upper) and with (lower) denoising trajectory alignment. After alignment, the generated images exhibit a reduction in deformations and artifacts, thereby achieving higher quality.

| $M_p$ | $M_q$ | $\gamma$ | $\alpha$/Speed | | |
|---|---|---|---|---|---|
| | | | 0% | 5% | 15% |
| MAR-L | MAR-B | 32 | 0.27/1.24× | 0.34/1.22× | **0.37**/1.21× |
| MAR-L | MAR-B | 16 | 0.35/1.19× | 0.37/1.19× | **0.38**/1.17× |
| MAR-L | MAR-B | 8 | 0.35/1.14× | **0.39**/1.13× | **0.39**/1.12× |
| MAR-L | MAR-B | 4 | 0.32/1.04× | 0.37/1.02× | **0.39**/1.00× |
| MAR-H | MAR-B | 32 | 0.25/1.63× | 0.30/1.63× | **0.33**/1.61× |
| MAR-H | MAR-L | 32 | 0.24/1.36× | **0.33**/1.35× | 0.32/1.34× |
| MAR-H | MAR-B | 16 | 0.32/1.53× | 0.33/1.52× | **0.34**/1.51× |
| MAR-H | MAR-L | 16 | 0.34/1.32× | **0.35**/1.29× | **0.35**/1.29× |
| MAR-H | MAR-B | 8 | 0.33/1.47× | 0.31/1.47× | **0.34**/1.44× |
| MAR-H | MAR-L | 8 | 0.34/1.21× | 0.34/1.21× | **0.35**/1.21× |
| MAR-H | MAR-B | 4 | 0.31/1.21× | 0.32/1.21× | **0.34**/1.20× |
| MAR-H | MAR-L | 4 | 0.31/1.05× | **0.34**/1.03× | **0.34**/1.03× |

Table 4. Ablation study on pre-filling ratio. The experimental configuration remains the same as Table 1. Underline indicates the highest speedup. **Bold** means the highest $\alpha$.

**Influence of pre-filled tokens.** The ablation study of pre-filling ratios at 0%, 5%, and 15% is illustrated in Figure 11. Pre-filling can compensate for the low acceptance rates observed during the initial stages of autoregressive sampling and enhance the overall acceptance rate, as shown in Table 4. Moreover, Figure 12 shows the visualizations under different pre-filling ratios. Notably, the discrepancies between the draft and the target model result in certain artifacts and reduced image quality at 0% pre-filling. However, introducing a modest proportion of pre-filled tokens from the target model has effectively mitigated these artifacts. As the pre-filling ratio increases, the advantages conferred by this approach exhibit diminishing returns.
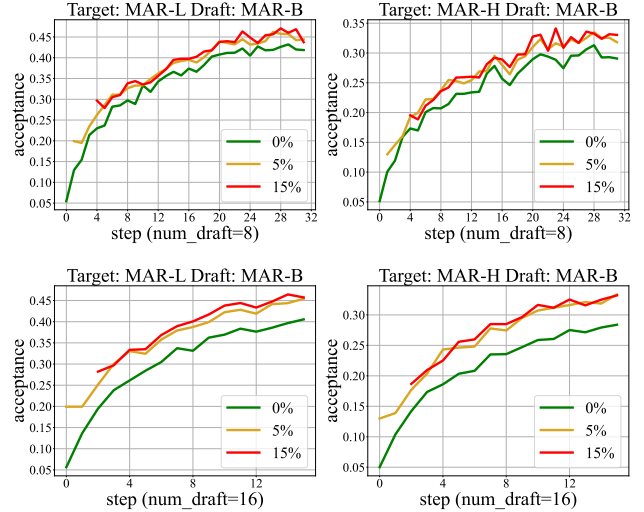


Figure 11. Per-step acceptance $\alpha$ under different pre-filling ratios. Acceptance rate per step is averaged on 1000 samples.



Figure 12. Comparing image generation quality under different token pre-filling portions.

## 5. Conclusion

We explore a novel approach to adapt speculative decoding to continuous-valued visual AR models. We analyze the critical challenges that hinder the algorithm's application in continuous space. To this end, the acceptance criterion is established. Moreover, the denoising trajectory alignment and token pre-filling enhance the acceptance rate during the diffusion and autoregressive process. The problem of sampling from the modified distribution that doesn't have an analytical form is tackled through acceptance-rejection sampling with an appropriately defined upper bound. Our continuous speculative decoding achieves up to $2.33\times$ speedup while maintaining output distribution and high generation fidelity. We hope our work will provide more thoughts and insights into the inference acceleration with continuous-valued autoregressive models in visual and other domains.

# References

[1] Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D Lee, Deming Chen, and Tri Dao. Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv preprint arXiv:2401.10774*, 2024. 3

[2] George Casella, Christian P Robert, and Martin T Wells. Generalized accept-reject sampling schemes. *Lecture notes-monograph series*, pages 342–347, 2004. 2, 5

[3] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 2

[4] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 1, 2

[5] Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*, 2023. 1, 3

[6] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *International Conference on Machine Learning*, pages 1691–1703. PMLR, 2020. 2

[7] Xi Chen, Nikhil Mishra, Mostafa Rohaninejad, and Pieter Abbeel. Pixelsnail: An improved autoregressive generative model. In *International Conference on Machine Learning*, pages 864–872. PMLR, 2018. 1, 2

[8] 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*, pages 248–255. Ieee, 2009. 5, 2

[9] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021. 2

[10] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 12873–12883, 2021. 1, 2

[11] Lijie Fan, Tianhong Li, Siyang Qin, Yuanzhen Li, Chen Sun, Michael Rubinstein, Deqing Sun, Kaiming He, and Yonglong Tian. Fluid: Scaling autoregressive text-to-image generative models with continuous tokens. *arXiv preprint arXiv:2410.13863*, 2024. 1, 3

[12] Karol Gregor, Ivo Danihelka, Andriy Mnih, Charles Blundell, and Daan Wierstra. Deep autoregressive networks. In *International Conference on Machine Learning*, pages 1242–1250. PMLR, 2014. 2

[13] Jiatao Gu, Yuyang Wang, Yizhe Zhang, Qihang Zhang, Dinghuai Zhang, Navdeep Jaitly, Josh Susskind, and Shuangfei Zhai. Dart: Denoising autoregressive transformer for scalable text-to-image generation. *arXiv preprint arXiv:2410.08159*, 2024. 1, 3

[14] 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. 2, 5

[15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 1, 3, 4

[16] Doohyuk Jang, Sihwan Park, June Yong Yang, Yeonsung Jung, Jihun Yun, Souvik Kundu, Sung-Yub Kim, and Eunho Yang. Lantern: Accelerating visual autoregressive models with relaxed speculative decoding. *arXiv preprint arXiv:2410.03355*, 2024. 2, 3

[17] Sehoon Kim, Karttikeya Mangalam, Suhong Moon, Jitendra Malik, Michael W Mahoney, Amir Gholami, and Kurt Keutzer. Speculative decoding with big little decoder. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[18] Siqi Kou, Lanxiang Hu, Zhezhi He, Zhijie Deng, and Hao Zhang. Cllms: Consistency large language models. *arXiv preprint arXiv:2403.00835*, 2024. 3

[19] Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR, 2023. 1, 3, 4, 2

[20] Tianhong Li, Huiwen Chang, Shlok Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan. Mage: Masked generative encoder to unify representation learning and image synthesis. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 2142–2152, 2023. 1, 2

[21] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024. 1, 2, 3, 4, 5, 6

[22] Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle: Speculative sampling requires rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*, 2024. 3

[23] Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle-2: Faster inference of language models with dynamic draft trees. *arXiv preprint arXiv:2406.16858*, 2024. 3

[24] Dongyang Liu, Shitian Zhao, Le Zhuo, Weifeng Lin, Yu Qiao, Hongsheng Li, and Peng Gao. Lumina-mgpt: Illuminate flexible photorealistic text-to-image generation with multimodal generative pretraining. *arXiv preprint arXiv:2408.02657*, 2024. 3

[25] Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Alvin Cheung, Zhijie Deng, Ion Stoica, and Hao Zhang. Online speculative decoding. *arXiv preprint arXiv:2310.07177*, 2023. 3

[26] Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization: Vq-vae made simple. *arXiv preprint arXiv:2309.15505*, 2023. 1, 3

[27] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, et al. Specinfer: Accelerating generative large language model serving with tree-based speculative inference and verification. *arXiv preprint arXiv:2305.09781*, 2023. 3

[28] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 1, 3, 4

[29] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International Conference on Machine Learning*, pages 4055–4064. PMLR, 2018. 2

[30] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 1, 2

[31] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in Neural Information Processing Systems*, 29, 2016. 2, 5

[32] Andrea Santilli, Silvio Severino, Emilian Postolache, Valentino Maiorca, Michele Mancusi, Riccardo Marin, and Emanuele Rodolà. Accelerating transformer inference for translation via parallel decoding. *arXiv preprint arXiv:2305.10427*, 2023. 3

[33] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 3

[34] Haotian Tang, Yecheng Wu, Shang Yang, Enze Xie, Junsong Chen, Junyu Chen, Zhuoyang Zhang, Han Cai, Yao Lu, and Song Han. Hart: Efficient visual generation with hybrid autoregressive transformer. *arXiv preprint arXiv:2410.10812*, 2024. 3

[35] Yao Teng, Han Shi, Xian Liu, Xuefei Ning, Guohao Dai, Yu Wang, Zhenguo Li, and Xihui Liu. Accelerating autoregressive text-to-image generation with training-free speculative jacobi decoding. *arXiv preprint arXiv:2410.01699*, 2024. 2, 3

[36] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. 2024. 1, 2

[37] Michael Tschannen, Cian Eastwood, and Fabian Mentzer. Givt: Generative infinite-vocabulary transformers. In *European Conference on Computer Vision*, pages 292–309. Springer, 2025. 1, 3

[38] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in Neural Information Processing Systems*, 29, 2016. 1, 2

[39] Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International Conference on Machine Learning*, pages 1747–1756. PMLR, 2016. 1, 2

[40] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in Neural Information Processing Systems*, 30, 2017. 2

[41] Yilun Xu, Gabriele Corso, Tommi Jaakkola, Arash Vahdat, and Karsten Kreis. Disco-diff: Enhancing continuous diffusion models with discrete latents. In *International Conference on Machine Learning*, 2024. 1, 3

[42] Hanling Yi, Feng Lin, Hongbin Li, Peiyang Ning, Xiaotian Yu, and Rong Xiao. Generation meets verification: Accelerating large language model inference with smart parallel auto-correct decoding. *arXiv preprint arXiv:2402.11809*, 2024. 3

[43] Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627*, 2021. 2

[44] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022. 3

[45] Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video transformer. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 10459–10469, 2023. 1, 3

[46] Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion–tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023. 3

[47] Qihang Yu, Ju He, Xueqing Den, Xiaohui Shen, and Liang-Chieh Chen. Randomized autoregressive visual generation. *arXiv preprint arXiv:2411.00776*, 2023. 3

[48] Weilin Zhao, Yuxiang Huang, Xu Han, Chaojun Xiao, Zhiyuan Liu, and Maosong Sun. Ouroboros: Speculative decoding with large model enhanced drafting. *arXiv preprint arXiv:2402.13720*, 2024. 3

[49] Yao Zhao, Zhitian Xie, Chen Liang, Chenyi Zhuang, and Jinjie Gu. Lookahead: An inference acceleration framework for large language model with lossless generation accuracy. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6344–6355, 2024. 3

[50] Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat, Aditya Krishna Menon, Afshin Rostamizadeh, Sanjiv Kumar, Jean-François Kagy, and Rishabh Agarwal. Distillspec: Improving speculative decoding via knowledge distillation. *arXiv preprint arXiv:2310.08461*, 2023. 3

# Continuous Speculative Decoding for Autoregressive Image Generation

## Supplementary Material

## A. Detailed Proof

We will provide a more detailed process and proof of continuous speculative decoding.

### A.1. Denoising Trajectory Alignment

For $\frac{p(x)}{q(x)}$, we obtain $x = x_0$ through the denoising process:

$$p(x_0|x_T) = p(x_T) \prod_{t=1}^{T} p(x_{t-1}|x_t), \tag{12}$$

with the conditioned probability distributions as Gaussian approximated by a neural network:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)). \tag{13}$$

Therefore, $p_\theta(x_{t-1}|x_t)$ can be calculated by the PDF of Gaussian distribution. The calculation of $q_\theta(x_{t-1}|x_t)$ is the same.

Empirically, $x_{t-1}$ is obtained by sampling from the Gaussian distribution on the right-hand side via **reparameterization**. That is, we first sample $\varepsilon_t \sim \mathcal{N}(0, \mathrm{I})$, and then we obtain the result by scale and shift $x_{t-1} = \sqrt{\Sigma_\theta(x_t, t)} \cdot \varepsilon_t + \mu_\theta(x_t, t)$. To this end, we can calculate $p(x)$ and $q(x)$ to obtain the ratio.

However, as described in Sec. 3, directly calculating the $p(x)$ and $q(x)$ is algebraically correct but may lead to a low acceptance rate due to a distinct denoising trajectory. Thus, we employ the same $\epsilon_t$ on both $p(x)$ and $q(x)$ to align their trajectory as closely as possible without affecting the denoising procedure and results.

Additionally, alignment also brings the simplification of calculating $\frac{p(x)}{q(x)}$. Note that in Gaussian distribution, we have:

$$p(x) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\Sigma|}} \exp\left\{\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\} \tag{14}$$

$$= \frac{1}{(\sqrt{2\pi})^n \sqrt{|\Sigma|}} \exp\left\{\varepsilon_t^T \varepsilon_t\right\}. \tag{15}$$

Since we have the same $\epsilon_t$ of both $p(x)$ and $q(x)$, the exponential term can be eliminated to obtain:

$$\frac{p(x)}{q(x)} = \frac{\frac{1}{(\sqrt{2\pi})^n \sqrt{|\Sigma_p|}} \exp\left\{\frac{1}{2}(x - \mu_p)^T \Sigma_p^{-1}(x - \mu_p)\right\}}{\frac{1}{(\sqrt{2\pi})^n \sqrt{|\Sigma_q|}} \exp\left\{\frac{1}{2}(x - \mu_q)^T \Sigma_q^{-1}(x - \mu_q)\right\}} \tag{16}$$

$$= \frac{\sqrt{|\Sigma_q|}}{\sqrt{|\Sigma_p|}}. \tag{17}$$

Therefore:

$$\frac{p(x)}{q(x)} = \frac{p(x_T) \prod_{t=2}^{T} p(x_{t-1}|x_t)}{q(x_T) \prod_{t=2}^{T} q(x_{t-1}|x_t)} \cdot \frac{p(x_0|x_1)}{q(x_0|x_1)} \tag{18}$$

$$= \frac{\prod_{t=2}^{T} \sqrt{|\Sigma_{q,t}|}}{\prod_{t=2}^{T} \sqrt{|\Sigma_{p,t}|}} \cdot \frac{p(x_0|x_1)}{q(x_0|x_1)} \tag{19}$$

$$= \Sigma \cdot \frac{p(x_0|x_1)}{q(x_0|x_1)}, \tag{20}$$

where $\Sigma$ is the product of $\frac{\sqrt{|\Sigma_q|}}{\sqrt{|\Sigma_p|}}$ along the denoising intermediate results. Also, since $x_0 \sim q(x)$ should be verified by the target model, $p(x_0|x_1)$ is not obtained by denoising. It is obtained by substituting $x_0$ to $p(x_0|x_1)$. We keep the two terms since they should be calculated separately.

### A.2. Acceptance-Rejection Sampling

After rejection, we should resample a new token from:

$$p'(x) = \frac{max(0, p(x) - q(x))}{\int_{x'} max(0, p(x') - q(x'))dx'}. \tag{21}$$

But $p(x) - q(x)$ is hard to obtain. This term represents the subtraction between multiple Gaussian product terms. Besides, the integral $Z = \int_{x'} max(0, p(x') - q(x'))dx'$ is also hard to compute and may introduce calculation errors if we employ approximation. This integral also does not have an analytical form.

Therefore, the introduction of acceptance-rejection sampling can eliminate $Z$ by $M = \frac{1}{Z}$ to:

$$\alpha_s = \frac{max(0, p(x) - q(x))/Z}{p(x)/Z} \tag{22}$$

$$= \frac{max(0, p(x) - q(x))}{p(x)} \tag{23}$$

$$= \frac{max(0, p(x_0|x_T) - q(x_0|x_T))}{p(x_0|x_T)} \tag{24}$$

$$= \frac{max(0, p(x_T) \prod p(x_{t-1}|x_t) - q(x_T) \prod q(x_{t-1}|x_t))}{p(x_T) \prod p(x_{t-1}|x_t)} \tag{25}$$

$$= \frac{max(0, \Sigma \cdot p_\theta(x_0|x_1^p) - q_\theta(x_0|x_1^q))}{\Sigma \cdot p_\theta(x_0|x_1^p)} \tag{26}$$

Afterward, we could get the computable results as long as eliminate the intermediate denoising term by $\Sigma$. And the final expression is easy to get. The modified distribution can be sampled by this way.

## A.3. Correctness of continuous speculative decoding

We will present the correctness of continuous speculative decoding. Let $\beta$ be the acceptance probability by:

$$\beta = E_{x \sim q(x)} \min(1, \frac{p(x)}{q(x)}) = \int_x \min(p(x), q(x))dx. \tag{27}$$

Note that:

$$p'(x) = \frac{max(0, p(x) - q(x))}{\int_{x'} max(0, p(x') - q(x'))dx'} \tag{28}$$

$$= \frac{p(x) - min(q(x), p(x))}{1 - \beta}. \tag{29}$$

Now we have:

$$P(x') = p(\text{accept}, x') + p(\text{reject}, x'), \tag{30}$$

where:

$$p(\text{accept}, x') = q(x') \min(1, \frac{p(x')}{q(x')}) = \min(q(x'), p(x')), \tag{31}$$

and:

$$p(\text{reject}, x') = (1 - \beta)p'(x') = p(x') - \min(q(x'), p(x')). \tag{32}$$

Overall, the output still yields:

$$P(x') = \min(q(x'), p(x')) + p(x') - \min(q(x'), p(x')) \tag{33}$$

$$= p(x'). \tag{34}$$

Algorithm 1 shows this procedure of the speculative decoding algorithm for continuous-valued tokens with our implementation.

## B. Limitations of Walltime Improvement

As described in [19], the expected walltime improvement is assumed to be:

$$\frac{1 - \alpha^{\gamma+1}}{(1 - \alpha)(\gamma c + 1)}, \tag{35}$$

where $\alpha$ is the acceptance ratio of draft tokens, $\gamma$ is the draft length, and $c$ is the inference time ratio between the draft and target models. However, due to the limited models available, the largest model we can obtain is MAR-H (943M). The inference time ratio $c$ by MAR-B (208M) is 0.38 (bs=128), which is **far more larger** than the number 0.05 or close to 0 mentioned in [19]. Also, theoretical improvement is assumed to be due to long enough generations. However, autoregressive image generation is currently limited to 256 generation steps. This further enhances the inapplicability of the theoretical results.

We anticipate that our algorithm will achieve more significant runtime improvements with larger models, like 7B, 13B, or even larger. This direction warrants further investigation in future research.

---

**Algorithm 1** ContinuousSpeculativeDecodingStep

---
**Inputs:** $M_p, M_q, prefix$.
▷ Sample $\gamma$ guesses $x_{1,\ldots,\gamma}$ from $M_q$ autoregressively.
**for** $i = 1$ **to** $\gamma$ **do**
  $q_i(x) \leftarrow M_q(prefix + [x_1, \ldots, x_{i-1}])$
  $x_i \sim q_i(x)$
**end for**
▷ Run $M_p$ in parallel, keep the $\epsilon_t$ the same in $M_q$
$p_1(x), \ldots, p_{\gamma+1}(x) \leftarrow$
    $M_p(prefix), \ldots, M_p(prefix + [x_1, \ldots, x_\gamma])$
$\Sigma \leftarrow \frac{\prod_{t=2}^T \sqrt{|\Sigma_{q,t}|}}{\prod_{t=2}^T \sqrt{|\Sigma_{p,t}|}}$
▷ Determine the number of accepted guesses $n$.
$r_1 \sim U(0,1), \ldots, r_\gamma \sim U(0,1)$
$\frac{p_i(x)}{q_i(x)} \leftarrow \Sigma \cdot \frac{p_i(x|x_1^p)}{q_i(x|x_1^q)}$
$n \leftarrow \min(\{i - 1 \mid 1 \le i \le \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})$
▷ Sample the modified distribution via
▷ acceptance-rejection sampling.
**if** $n < \gamma$ **then**
  **repeat**
    $x_t \leftarrow p_n(x|x_1^p)$
    $\alpha_s \leftarrow \frac{max(0, \Sigma \cdot p_n(x_t|x_1^p) - q_n(x_t|x_1^q))}{\Sigma \cdot p_n(x_t|x_1^p)}$
    $\epsilon \sim U(0,1)$
  **until** $\epsilon \le \alpha_s$
**end if**
▷ Return one token from $M_p$, and $n$ tokens from $M_q$.
**return** $prefix + [x_1, \ldots, x_n, x_t]$

---

## C. Implementation Details

We conduct experiments with open-sourced continuous-valued visual autoregressive model MAR [21] on ImageNet [8] $256 \times 256$ generation. The draft model is chosen from MAR-B (208M) and MAR-L (479M). The target model is chosen from MAR-L and MAR-H (943M), respectively. We use official pretrained checkpoints for all models. As commonly practiced, we set temperature and classifier-free guidance the same way as MAR does. Due to the limited open-sourced model, we couldn't implement our algorithm on larger models for further results. However, default MAR models have shown significant results for bidirectional attention in MAR. When target models verify the draft tokens, each output token can be regarded as the last since they can see every previous token. Besides, both models utilize their respective class tokens [cls], which are not shared during the speculative decoding process. Their diffusion loss is not shared either. The generation speed is measured on a single NVIDIA A100 GPU, with batch size ranging in {1, 8, 128, 256}. The FID and IS are calculated on 50k generated images. The results are averaged on ten runs of evaluations.

## D. Additional Experiments

### D.1. Classifier-Free Guidance

Figure 13 illustrates the relationship between draft length and the acceptance rate under different CFG scales. As the CFG scale increases, there is an overall trend of decreasing acceptance rates. This trend remains consistent mainly across each draft length. This phenomenon may indicate that as class guidance strengthens, the inconsistency between the draft model and the target model may increase, further reducing the acceptance rate.
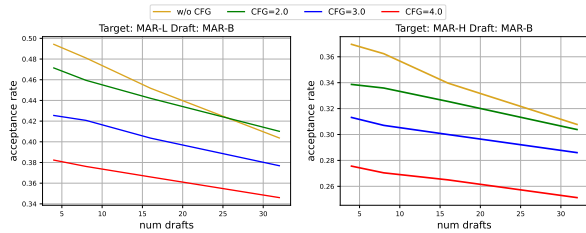


Figure 13. CFG scale has has a significant impact on the acceptance rate under different number of drafts.

### D.2. Temperature

During the denoising process of the diffusion model in MAR, temperature $\tau$ is a crucial hyperparameter. The temperature setting affects the consistency between the outputs of the draft and target models. Figure 14 illustrates the impact of the temperature $\tau$ on the acceptance rate during the generation process. The number of drafts is set to 8. The temperature influences the PDF of the final output distribution; a lower temperature may result in a sharper distribution, while a higher temperature may lead to a flatter distribution. The ratio $p(x)/q(x)$ can be influenced based on this.
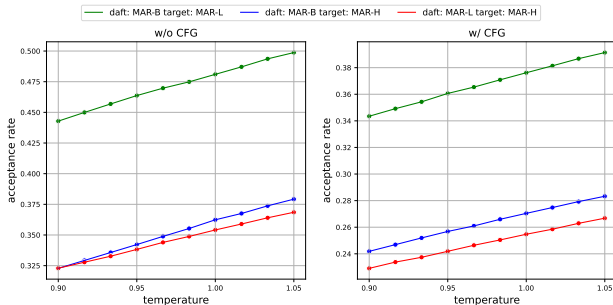


Figure 14. Temperature influence on the acceptance rate. Left: without CFG. Right: with CFG.

### D.3. Acceptance-Rejection Sampling

We observe the empirical sampling trial times during acceptance-rejection sampling. Figure 15 illustrates the re-
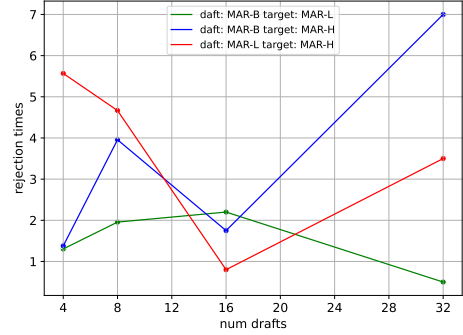


Figure 15. Empirical rejection times in acceptance-rejection sampling algorithm of the rejection phase.

lationship between the rejection times and the draft length. Empirically, acceptance-rejection sampling often requires only a few sampling steps. The runtime overhead consumed by this sampling process is negligible compared to the overall model inference time.

### D.4. Visualization of Acceptance

We visualize the acceptance and rejection of each token through a 2D heatmap. As shown in Figure 16, dark green blocks represent tokens that have been accepted, while light green blocks represent tokens that have been rejected. We observe that tokens representing backgrounds or regions with simpler textures tend to be accepted. In contrast, more detailed positions are more likely to be rejected.
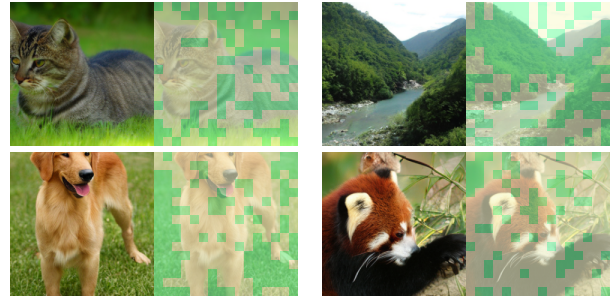


Figure 16. Visualizations of accepted token heatmap. Dark green: accepted. Light green: rejected.

## E. More Qualitative Results

In Figure 17, 18, 19, 20 and 21, we provide more additional images generated under our continuous speculative decoding compared with the target model only one. While the target model has achieved satisfactory quality in generating realistic and high-fidelity images, our continuous speculative decoding can show comparable performance, similar generation results, and much faster inference speed.
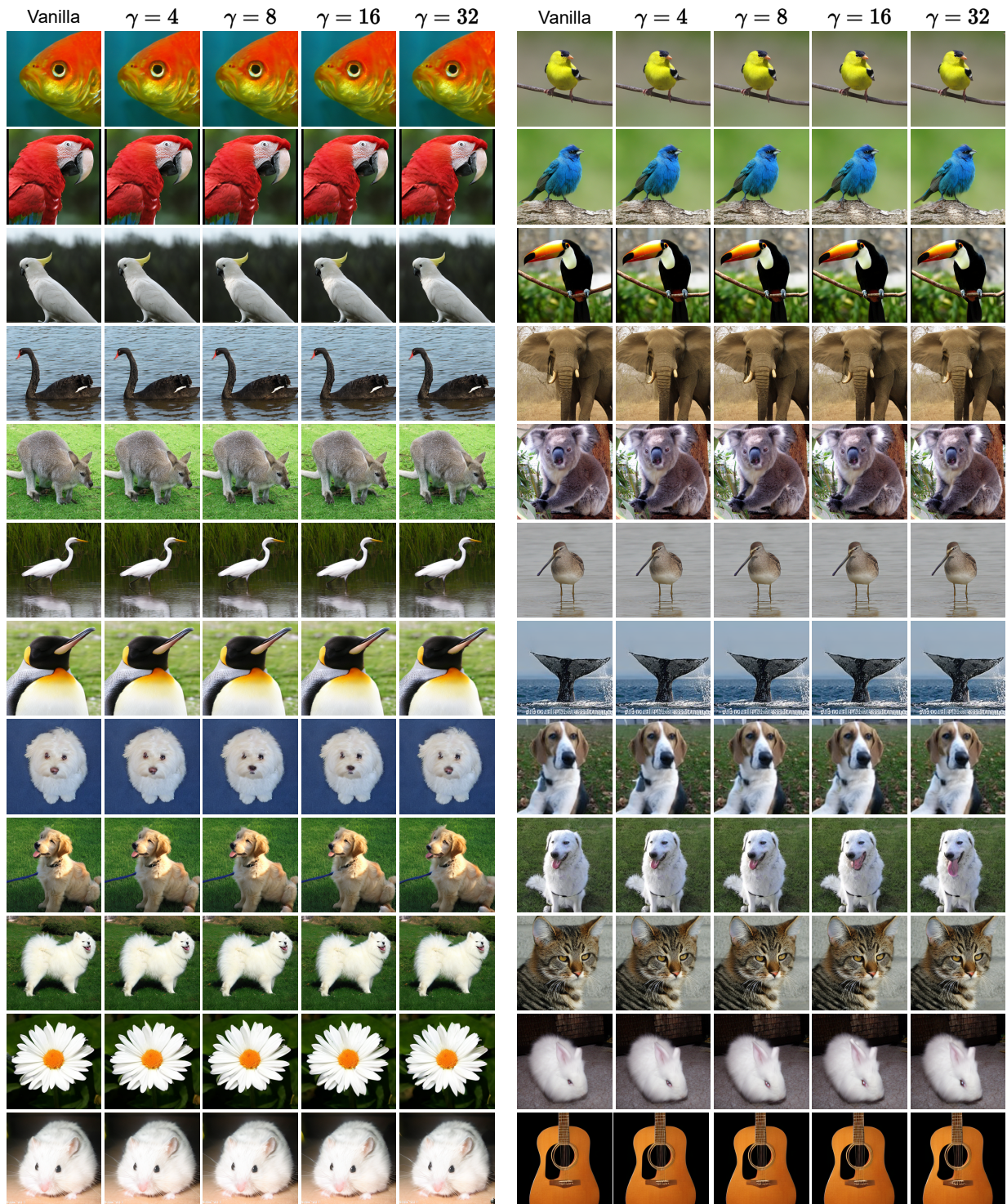
Figure 17. Visual quality with increasing draft length $\gamma$ compared with vanilla target model only generation. *Best viewed zoom-in.*

Figure 18. Visualization examples under $\gamma = 4$. Class label: arctic fox (297).

Figure 19. Visualization examples under $\gamma = 8$. Class label: balloon (417).

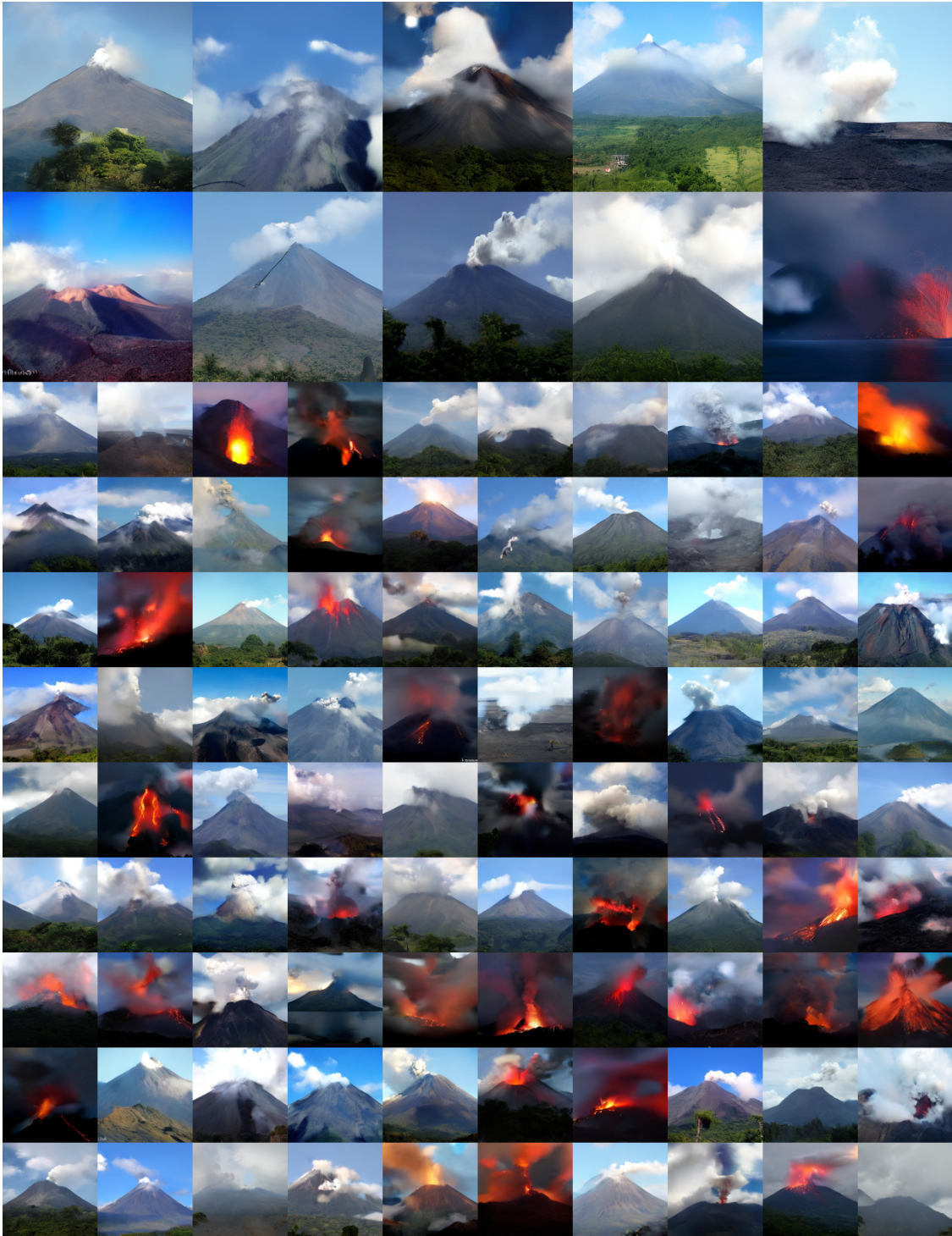Figure 20. Visualization examples under $\gamma = 16$. Class label: ice cream (928).

Figure 21. Visualization examples under $\gamma = 32$. Class label: volcano (980).