

# LANTERN: ACCELERATING VISUAL AUTOREGRESSIVE MODELS WITH RELAXED SPECULATIVE DECODING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Auto-Regressive (AR) models have recently gained prominence in image generation, often matching or even surpassing the performance of diffusion models. However, one major limitation of AR models is their sequential nature, which processes tokens one at a time, slowing down generation compared to models like GANs or diffusion-based methods that operate more efficiently. While speculative decoding has proven effective for accelerating LLMs by generating multiple tokens in a single forward, its application in visual AR models remains largely unexplored. In this work, we identify a challenge in this setting, which we term *token selection ambiguity*, wherein visual AR models frequently assign uniformly low probabilities to tokens, hampering the performance of speculative decoding. To overcome this challenge, we propose a relaxed acceptance condition referred to as LANTERN that leverages the interchangeability of tokens in latent space. This relaxation restores the effectiveness of speculative decoding in visual AR models by enabling more flexible use of candidate tokens that would otherwise be prematurely rejected. Furthermore, by incorporating a total variation distance bound, we ensure that these speed gains are achieved without significantly compromising image quality or semantic coherence. Experimental results demonstrate the efficacy of our method in providing a substantial speed-up over speculative decoding. In specific, compared to a naïve application of the state-of-the-art speculative decoding, LANTERN increases speed-ups by  $1.75\times$  and  $1.82\times$ , as compared to greedy decoding and random sampling, respectively, when applied to LlamaGen, a contemporary visual AR model.

## 1 INTRODUCTION

Auto-Regressive (AR) models have recently gained significant traction in image generation (Ramesh et al., 2021; Chen et al., 2020; Tian et al., 2024; Sun et al., 2024) due to their competitive performance, often matching or even surpassing diffusion models (Ho et al., 2020; Rombach et al., 2022). Notable examples include iGPT (Chen et al., 2020), DALL-E (Ramesh et al., 2021), VAR (Tian et al., 2024), and LlamaGen (Sun et al., 2024), which showcase the potential of AR models in image generation. Moreover, recent studies like Lu et al. (2023); Team (2024); Chern et al. (2024) have demonstrated that AR modeling can handle multi-modal data, including language and images, within a single unified framework. Given the remarkable success of AR models in language modeling, leading to the era of large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023; Jiang et al., 2023), it is anticipated that AR modeling will emerge as a dominant paradigm for unifying multiple modalities into a single model in the near future.

Despite the promising potential of AR models, their sequential nature poses a significant bottleneck for both efficiency and scalability since they generate a single token per forward pass. In contrast, GANs (Goodfellow et al., 2014; Karras et al., 2019), which generate images in a single forward pass, naturally avoid this issue, and diffusion models have benefited from extensive research aimed at improving their speed (Song et al., 2022; Sauer et al., 2023; Heek et al., 2024). However, transferring these acceleration techniques to visual AR models is far from straightforward due to fundamental differences in the underlying mechanisms of these models.



070 Figure 1: Images generated by vanilla decoding (*top*) and lossy speculative decoding with our relaxed  
071 acceptance condition (*bottom*) on the text-conditioned LlamaGen-XL Stage II (Sun et al., 2024). The  
072 mean accepted length for each image is displayed in white at the bottom right corner of each image.

073 One notable acceleration technique for AR models is speculative decoding (Leviathan et al., 2023;  
074 Chen et al., 2023; Cai et al., 2024; Li et al., 2024b), which has demonstrated its effectiveness in  
075 LLMs. Initially introduced by Leviathan et al. (2023), speculative decoding addresses the sequential  
076 bottleneck of AR models by introducing a *draft and verify* mechanism. In this framework, a smaller  
077 model (the *drafter*) predicts the next few tokens, which are then verified by the larger target model. If  
078 the drafter’s predictions are accurate, multiple tokens can be generated from a single forward pass,  
079 resulting in a substantial inference speed-up. This method has proven highly effective in accelerating  
080 LLM inference, making it a leading option for reducing the latency associated with AR models.

081 While speculative decoding shows great promise for accelerating AR models, its application to  
082 visual AR models remains largely unexplored. Therefore, in this paper, we take the first step toward  
083 addressing this gap by migrating speculative decoding to visual AR models. Interestingly, our findings  
084 reveal that the naïve application of existing speculative decoding methods falls short in visual AR  
085 models. Specifically, we identify a key problem, namely the *token selection ambiguity*, that hampers  
086 the effective migration of speculative decoding to visual AR models.

087 To mitigate such an obstacle in speculative decoding, we propose a solution dubbed as **LANTERN**  
088 (**L**atent **N**eighbor **T**oken **A**cceptance **R**elaxation) that leverages the interchangeability of image  
089 tokens in latent space for the relaxation of acceptance condition. By relaxing the acceptance in  
090 speculative decoding, we allow for more effective utilization of draft (candidate) tokens that would  
091 otherwise be frequently rejected despite their potential usefulness. However, our relaxation introduces  
092 some distortion to the target model’s distribution, which may cause the generated images to deviate  
093 from the original target output. To mitigate this, we further incorporate a total variation distance  
094 bound which ensures that the deviation remains controlled.

095 Our main contributions are summarized as below:

- 096 • To the best of our knowledge, we are the first to thoroughly investigate speculative decoding  
097 in visual AR models, identifying the *token selection ambiguity* problem, where near-uniform  
098 token probability distributions hinder token prioritization, causing existing methods to fail  
099 in improving speed.
- 100 • Based on our insights, we then propose LANTERN, a novel relaxation of acceptance  
101 condition for the speculative decoding that addresses the token selection ambiguity problem,  
102 successfully enabling the effective application of speculative decoding to visual AR models.
- 103 • Our experiments using LlamaGen (Sun et al., 2024) as the target and EAGLE-2 (Li et al.,  
104 2024a) as the base speculative decoding method demonstrate significant speed-ups, im-  
105 proving from 1.29× to 2.26× in greedy decoding and from 0.93× to 1.69× in random  
106 sampling, compared to the naïve application of EAGLE-2, without substantial performance  
107 drop in terms of image quality.

## 2 PRELIMINARIES (NEWLY ADDED)

**Notations** In this paper, the *target model* refers to the visual AR model we aim to accelerate. In contrast, the *drafter model* is a supplementary model used to generate draft tokens. The probability distributions modeled by the drafter and target models are represented by  $p(\cdot|\cdot)$  and  $q(\cdot|\cdot)$ , respectively. Individual tokens are denoted in lowercase  $x$ , and sequences are represented by uppercase  $X$ ; for instance,  $X_{i:j}$  represents the sequence  $(x_i, \dots, x_j)$ . The concatenation of sequences  $X_{1:N}$  and  $Y_{1:M}$  is denoted by  $(X_{1:N}, Y_{1:M}) = (x_1, \dots, x_N, y_1, \dots, y_M)$ . For simplicity, we allow certain notational liberties; for instance, expressions like  $X_{1:0}$  are treated as empty sequences to avoid unnecessary complexity in notation.

**Visual Autoregressive Modeling** In visual AR models, image generation involves two main stages: generating image tokens through auto-regression and decoding the image tokens into actual image patches. In a text-to-image generation setting, given a tokenized text prompt  $X_{1:N}$ , the model generates a sequence of image tokens  $X_{N+1:N+K}$  based on the following probability modeling:

$$P(X_{N+1:N+K} | X_{1:N}) = \prod_{\ell=1}^K P(x_{N+\ell} | X_{1:N+\ell-1}),$$

where  $K$  represents the total number of image tokens corresponding to the height and width of the image feature map. Since each token is predicted based solely on its preceding tokens, visual AR models require  $K$  sequential steps to generate all  $K$  image tokens.

Once the  $K$  image tokens are generated, they are mapped to visual representation by referring to the codebook  $\mathcal{C}$ . Specifically, the codebook  $\mathcal{C} = \{c_1, \dots, c_L\}$  consists of codes  $c_i \in \mathbb{R}^d$ , where each code is a latent feature used by the image encoder-decoder pair (such as VQVAE (Van Den Oord et al., 2017) or VQGAN (Esser et al., 2021)), and  $d$  is the dimensionality of these features. As the text tokenization, since the image tokens represent indices in the codebook  $\mathcal{C}$ , each image token  $x_{N+i}$  maps to  $c_{x_{N+i}}$ , which is then rearranged into a  $h \times w \times d$  shaped tensor in raster-scan order (from top-left to bottom-right), for  $h = H/f$  and  $w = W/f$  when  $f$  is down-sampling factor. After that, rearranged latent is fed into the image decoder  $D: \mathbb{R}^{h \times w \times d} \rightarrow \mathbb{R}^{H \times W \times C}$  to construct an actual RGB image.

**Speculative Decoding** We briefly introduce the basics of speculative decoding, which were proposed by Leviathan et al. (2023). Given a sequence of tokens  $X_{1:N} = (x_1, \dots, x_N)$ , the drafter model generates  $\gamma$  draft tokens  $\tilde{X}_{1:\gamma}$  as speculations for the next  $\gamma$  tokens following  $X_{1:N}$ . Each draft token  $\tilde{x}_i$  is sampled from the drafter distribution  $p(x|(X_{1:N}, \tilde{X}_{1:i-1}))$  for each  $i = 1, \dots, \gamma$ , thus requiring  $\gamma$  forward steps to generate  $\gamma$  draft tokens.

After obtaining the draft tokens, the concatenated sequence  $(X_{1:N}, \tilde{X}_{1:\gamma})$  is fed into the target model, which calculates the likelihood of each draft token  $q(\tilde{x}_i|(X_{1:N}, \tilde{X}_{1:i-1}))$  in parallel within a single forward pass. Each draft token  $\tilde{x}_i$  is accepted with a probability:

$$\min \left( 1, \frac{q(\tilde{x}_i|(X_{1:N}, \tilde{X}_{1:i-1}))}{p(\tilde{x}_i|(X_{1:N}, \tilde{X}_{1:i-1}))} \right).$$

If  $\tilde{x}_i$  is accepted, it is immediately set as  $x_{N+i} = \tilde{x}_i$ . Otherwise, all subsequent tokens  $\tilde{X}_{i+1:\gamma}$  are discarded, and  $x_{N+i}$  is resampled from a distribution defined by  $[q(x|(X_{1:N}, \tilde{X}_{1:i-1})) - p(x|(X_{1:N}, \tilde{X}_{1:i-1}))]_+$ , where  $[\cdot]_+$  denotes normalization over positive values only. As proven by Leviathan et al. (2023), this approach ensures that the distribution of the generated token sequence matches the distribution produced by the target model. Related works on the visual AR models and speculative decoding are presented in Appendix A.

## 3 TOKEN SELECTION AMBIGUITY LIMITS SPECULATIVE DECODING IN VISUAL AR MODELS (REVISED)

Although speculative decoding has been highly successful in LLMs, it fails to deliver comparable speed improvements when naively applied to visual AR models. As illustrated in Figure 2(a), the

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

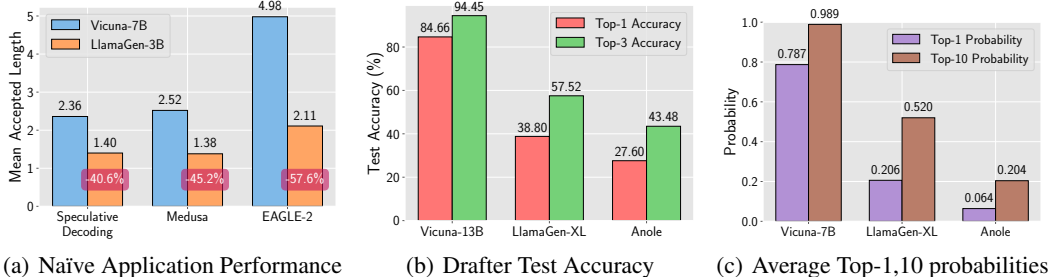


Figure 2: (a) Mean accepted length of naïve application of existing speculative decoding methods on visual AR model and LLM counterpart. (b) Top-1 and top-3 accuracy of learned drafter model for predicting the target model’s outputs. (c) An average top-1 and top-10 probabilities in the next token prediction.

visual AR model (LlamaGen-3B (Sun et al., 2024)) exhibits a 41% to 58% reduction in mean accepted length compared to their performance in LLM (Vicuna-7B (Zheng et al., 2023)), a consistent decline observed across different speculative decoding methods. Toward an in-depth exploration of this performance degradation, we train drafters for various visual AR models and conduct additional analyses on their predictions.

Our analysis reveals that the drafter in visual AR models (LlamaGen-XL and Anole (Chern et al., 2024)) struggles to predict the target model’s outputs accurately. Specifically, as demonstrated in Figure 2(b), the trained drafters for visual AR models fail to capture the target model’s prediction precisely, whereas the drafter in LLMs exhibits considerably higher accuracy. Such a low accuracy of the drafters for visual AR models leads to a significant reduction in the speed-up achievable via speculative decoding.

To delve into the root cause of the drafter’s performance degradation, we analyze the next token probabilities of visual AR models and identify a unique problem we term *token selection ambiguity*. As shown in Figure 2(c), visual AR models present substantially lower average top-1 and top-10 probabilities in next token predictive distributions compared to language models, indicating that visual AR models are more ambiguous when selecting the next token. This lack of prioritization among tokens reflects the model’s limited confidence in any single option.

We hypothesize that this issue arises from fundamental differences between image and language data. In language models, tokens represent discrete units, such as words or subwords, that form structured and predictable sequences governed by grammar and syntax (Zipf, 1935). Consequently, the next-token probabilities are more concentrated, and the model has high confidence in the most likely token. In contrast, visual AR models treat pixels or patches as tokens, forming continuous and highly complex sequences. As a result, these models exhibit more dispersed next token probabilities and face higher uncertainty and ambiguity in predicting subsequent tokens. This *token selection ambiguity* problem in visual AR models hinders the drafter’s predictive accuracy, thereby limiting the effectiveness of speculative decoding. Further details on the setup of experiments can be found in Appendix B.1.

#### 4 DETOURING TOKEN SELECTION AMBIGUITY THROUGH LATENT SPACE

In this section, we propose LANTERN, a simple yet effective method that permits detouring the failure of speculative decoding caused by the token selection ambiguity problem by relaxing acceptance condition in Speculative Decoding (Leviathan et al., 2023). In Section 4.1, we introduce a concept of latent proximity which asserts close image tokens in the latent space are interchangeable and examines its validity on the generated images. Section 4.2 describes how we relax the acceptance condition based on the interchangeability. In Section 4.3, we present another component to ensure that the distribution of generated images does not catastrophically deviate from the original distribution.



216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269



Figure 3: Image generated by text-conditioned LlamaGen-XL Stage II model (Sun et al., 2024). The images are generated by either standard sampling method (*top*) or sampling with random replacement within 100-closest tokens in the latent space (*bottom*).

#### 4.1 LATENT PROXIMITY PERMITS TOKEN INTERCHANGEABILITY

We introduce *latent proximity*, a property in visual AR models that asserts tokens close to one another in latent space are *interchangeable* without significantly affecting the visual semantics or overall image quality. This means that replacing a token with another nearby token in latent space results in minimal changes to the generated image.

This property arises from the tokenization process unique to visual AR models. Unlike text tokenization, which is straightforward due to its discrete nature, images are spatially continuous, making tokenization more complex Esser et al. (2021). To handle this, models like VQVAEs (Van Den Oord et al., 2017) and VQGANs (Esser et al., 2021) are used to discretize the latent embeddings of images, as introduced in Section 2. These embeddings maintain a continuous mapping between changes in latent space and the visual semantics of the generated images (Kingma & Welling, 2022; Goodfellow et al., 2014; Karras et al., 2019). As a result, small shifts in latent space lead to minor shifts in the image, supporting the idea that tokens close in the latent space are effectively interchangeable.

To demonstrate this concept empirically, we perform an experiment in which, after each token is sampled, it is re-sampled uniformly from the 100 closest tokens in latent space. Figure 3 reveals that the images generated by this procedure closely resemble those produced using the original sampling method. This confirms that tokens close in latent space can be treated as interchangeable, allowing for flexible token replacement without significantly compromising the visual semantics or overall image quality. A more detailed analysis of latent proximity can be found in Appendix C.

#### 4.2 LANTERN: RELAXATION OF ACCEPTANCE CONDITION

Building on our findings about latent proximity, we introduce LANTERN, a simple yet effective solution that leverages the interchangeability of proximate tokens in latent space. By treating neighboring tokens as commutable, LANTERN effectively resolves the token selection ambiguity problem, significantly boosting the acceptance probability of candidate tokens and enabling the successful application of speculative decoding.

We start with revisiting the original acceptance condition from Leviathan et al. (2023), introduced in Section 2. The drafter model samples a draft token  $\tilde{x} \sim p(x|s)$  given a preceding sequence  $s = X_{1:N}$ . The draft token is accepted with probability

$$\min \left( 1, \frac{q(\tilde{x}|s)}{p(\tilde{x}|s)} \right) \tag{1}$$

and if rejected, the next token is re-sampled from  $[q(x|s) - p(x|s)]_+$ . Note that the acceptance depends on the alignment of probabilities between the drafter and target models.

Table 1: Average accept probabilities of LANTERN. We only use accept probability of the first draft token. An average accept probability of EAGLE-2 (Li et al., 2024a) is **0.0402**.

Average Accept Probability				
$k$	$\delta = 0.05$	$\delta = 0.1$	$\delta = 0.2$	$\delta = 0.4$
100	0.0725	0.1096	0.1703	0.2595
300	0.0759	0.1166	0.1892	0.3267
1000	0.0786	0.1186	0.2000	0.3657

However, this acceptance condition results in a sharp decline in accept probability when encountering the token selection ambiguity problem. As mentioned in Table 1, the EAGLE-2 drafter exhibits an average accept probability of 0.0402, meaning *only 4% of drafts are accepted* when applied to visual AR models. This issue arises because the target model assigns low probabilities to individual tokens and frequently misaligns with the drafter’s distribution, leading to frequent rejections of candidate tokens and reducing the overall effectiveness of speculative decoding.

To alleviate this problem, we exploit the latent proximity by aggregating the probabilities of a candidate token’s nearest neighbors, treating them as proxies. This approach effectively increases the acceptance probability by utilizing the combined likelihood of similar tokens, mitigating the impact of the token selection ambiguity problem and reducing unnecessary rejections.

Specifically, we define the neighborhood  $B_k(\tilde{x})$  as the set of  $k$ -nearest tokens to  $\tilde{x}$  in latent space, including  $\tilde{x}$  itself. The accept probability is then adjusted to

$$\min \left( 1, \frac{\sum_{x \in B_k(\tilde{x})} q(x|s)}{p(\tilde{x}|s)} \right). \quad (2)$$

Because  $\tilde{x}$  is always included in  $B_k(\tilde{x})$ , this new accept probability is guaranteed to be equal to or higher than the original acceptance condition (1). As demonstrated in Table 1, applying LANTERN significantly increases the average acceptance probability, reaching values as high as 0.37. This improvement allows us to recover candidate tokens that would have otherwise been unjustifiably rejected.

### 4.3 WITH LIMITED DISTRIBUTIONAL DIVERGENCE

Although the relaxed acceptance condition (2) effectively permits speculative decoding in visual AR models by significantly raising the accept probability, it inevitably distorts the target distribution. In particular, when we condition the target distribution on the candidate token  $\tilde{x}$ , it becomes:

$$q_k(x|s, D = \tilde{x}) = \begin{cases} \sum_{x \in B_k(\tilde{x})} q(x|s) & \text{if } x = \tilde{x} \\ 0 & \text{if } x \in B_k(\tilde{x}), x \neq \tilde{x} \\ q(x|s) & \text{otherwise} \end{cases}$$

where  $D$  is a random variable representing the candidate token and  $q_k$  denotes the distorted target distribution. In contrast, under the original acceptance condition, the target distribution remains unchanged regardless of the candidate token. For this reason, (2) may excessively distort the target distribution, leading to generating images that diverge significantly from those generated by the target model.

To mitigate this distortion, we impose an upper bound on the distributional divergence using total variation distance (TVD). Since the distortion results from redistributing probability mass, TVD effectively measures the extent of this shift, allowing us to control the magnitude of the divergence. This can be achieved by adjusting the neighborhood  $B_k(\tilde{x})$  used in the relaxation as follows.

Since the relaxation can be analogously derived using any neighborhood of  $\tilde{x}$ , we can find a neighborhood that ensures the TVD between the target distribution and the distorted target distribution induced by the neighborhood is below a specific threshold. To formulate this approach, we define the neighborhood  $A_{k,\delta}(\tilde{x})$  of  $\tilde{x}$  for a given TVD bound  $\delta > 0$  and  $k \in \mathbb{Z}^+$  as  $A_{k,\delta}(\tilde{x})$  is the largest subset of  $B_k(\tilde{x})$  such that for the total variation distance  $D_{TV}$ ,

$$D_{TV}(q_{k,\delta}(x|s, D = \tilde{x}), q(x|s, D = \tilde{x})) = D_{TV}(q_{k,\delta}(x|s, D = \tilde{x}), q(x|s)) < \delta$$

where  $q_{k,\delta}$  denotes the distorted target distribution induced by  $A_{k,\delta}(\tilde{x})$ .

We construct  $A_{k,\delta}(\tilde{x})$  by incrementally adding tokens from  $B_k(\tilde{x})$  to  $A_{k,\delta}(\tilde{x})$ , starting with the closest ones to  $\tilde{x}$ , and stopping when adding another token would exceed the TVD threshold  $\delta$ . This procedure allows us to relax the acceptance condition by incorporating probabilities of similar tokens while keeping the divergence within a predefined boundary.

By integrating the TVD constraint into the acceptance condition (2), we arrive at the final relaxed acceptance condition of LANTERN:

$$\text{Accept } \tilde{x} \text{ with probability } \min\left(1, \frac{\sum_{x \in A_{k,\delta}(\tilde{x})} q(x|s)}{p(\tilde{x}|s)}\right)$$

$$\text{Else re-sample } x \sim [q_{k,\delta}(x|s, D = \tilde{x}) - p(x|s)]_+$$

For greedy decoding, LANTERN can be simply reduced to accept  $\tilde{x}$  if  $\tilde{x} = \arg \max_x q_{k,\delta}(x|s, D = \tilde{x})$ . The algorithms for LANTERN and construction of  $A_{k,\delta}$  can be found in Appendix D.

## 5 EXPERIMENTS

In this section, we empirically demonstrate the efficacy of our method, LANTERN. In section 5.1, we report the experimental setup for our experiments. Section 5.2 evaluates LANTERN with other baselines in perspective of image quality and acceleration. In Section 5.3, we conduct an ablation study to clarify the effectiveness of each component in our method.

### 5.1 EXPERIMENTAL SETUP

To validate our method LANTERN, we conduct experiments on text-conditioned LlamaGen-XL Stage I model (Sun et al., 2024) as a target model that establishes the best performance among visual AR models without vision-specific modifications. We utilize the MS-COCO validation captions (Lin et al., 2014) to generate images and evaluate the image quality with the ground-truth images. For the base speculative decoding method, we employ EAGLE-2 (Li et al., 2024a), which has demonstrated state-of-the-art performance in speculative decoding in the language domain. We employ the  $\ell_2$  distance to quantify latent proximity and utilize the TVD as a metric for divergence bound.

For the assessment of speed-ups, we use 1000 MS-COCO validation captions because this sample size provides results that are nearly identical to those obtained when evaluating the entire dataset. Actual speed-up is measured by the inference time ratio between each method and standard auto-regressive decoding. The mean accepted length is determined by the average number of tokens accepted in each forward step of the target model. We evaluate each method in both the greedy decoding setting with  $\tau = 0$  and the sampling with  $\tau = 1$ . The statistical analysis on actual speed-up and number of captions can be found in Appendix F.1.

Since LANTERN can impact the quality of generated images, we evaluate image quality using FID (Heusel et al., 2017), CLIP score (Hessel et al., 2021), Precision and Recall (Kynkäänniemi et al., 2019), and HPS v2 (Wu et al., 2023) with 30K samples. For each evaluation, we do not evaluate image quality on EAGLE-2 since it theoretically guarantees exact distribution matching with the target model. Further details can be found in Appendix B.2.

### 5.2 MAIN RESULTS

In this section, we demonstrate qualitative and quantitative results of speculative decoding with our relaxed acceptance condition. First of all, Section 5.2.1 demonstrates how much speed-up can be achieved through our method. Next, Section 5.2.2 showcases that our method retains image quality within a similar level by showing qualitative samples. Afterward, Section 5.2.3 presents the trade-off between performance and efficiency in our method.

#### 5.2.1 SPEED UP COMPARISON

To confirm that LANTERN provides notable speed improvements over the baseline while maintaining image quality, we compare our method with baselines under both greedy decoding and random

Table 2: Actual speed-up, MAL (Mean Accepted Length), FID, CLIP score, Precision / Recall, and HPS v2 for each method.  $\tau = 0$  refers to the greedy decoding and  $\tau = 1$  refers to the sampling with temperature 1 for generation. The actual speed-up is measured on a single RTX 3090.

Method	$\tau = 0$						
	Acceleration ( $\uparrow$ )		Image Quality Metrics				
	Speed-up	MAL	FID ( $\downarrow$ )	CLIP score ( $\uparrow$ )	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )	HPSv2 ( $\uparrow$ )
Vanilla AR (Sun et al., 2024)	1.00 $\times$	1.00	28.63	0.3169	0.4232	0.3517	23.18
EAGLE-2 (Li et al., 2024a)	1.29 $\times$	1.60	-	-	-	-	-
LANTERN ( $\delta = 0.05, k = 1000$ )	1.56 $\times$	2.02	29.77	0.3164	0.4484	0.3158	22.62
LANTERN ( $\delta = 0.2, k = 1000$ )	<b>2.26<math>\times</math></b>	<b>2.89</b>	30.78	0.3154	0.4771	0.2773	21.69
Method	$\tau = 1$						
	Acceleration ( $\uparrow$ )		Image Quality Metrics				
	Speed-up	MAL	FID ( $\downarrow$ )	CLIP score ( $\uparrow$ )	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )	HPSv2 ( $\uparrow$ )
Vanilla AR (Sun et al., 2024)	1.00 $\times$	1.00	15.22	0.3203	0.4781	0.5633	24.11
EAGLE-2 (Li et al., 2024a)	0.93 $\times$	1.20	-	-	-	-	-
LANTERN ( $\delta = 0.1, k = 1000$ )	1.13 $\times$	1.75	16.17	0.3208	0.4869	0.5172	23.75
LANTERN ( $\delta = 0.4, k = 1000$ )	<b>1.69<math>\times</math></b>	<b>2.40</b>	18.76	0.3206	0.4909	0.4497	23.22



Figure 4: Qualitative samples generated by LlamaGen-XL Stage II model for LANTERN and standard autoregressive decoding. From top to bottom, the images are generated by standard autoregressive decoding, LANTERN ( $\delta = 0.2, \delta = 0.4$ ) where  $k$  is fixed at 1000, and images in the same column are generated using the same text prompt. Text prompts for the images are provided in Appendix I.

sampling situations. Table 2 demonstrates the speed-up and image quality across different methods. LANTERN shows a significant acceleration, even with a slight degradation in image quality, when compared to standard autoregressive decoding and EAGLE-2 (Li et al., 2024a). Our method, LANTERN demonstrates significant improvements in both mean accepted length and actual speed-up. In terms of mean accepted length, LANTERN outperforms both baselines by reaching 2.40 $\times$  and it translate to substantial actual speed-up in practice, with LANTERN achieving 1.69 $\times$  actual speed-up.

While LANTERN’s acceleration comes with a trade-off in image quality, the degradation remains minimal and well within acceptable boundaries. For instance, the HPS v2 score decreases by less than 1.5 for both greedy decoding and sampling when compared to standard autoregressive decoding. Precision remains stable or slightly improved, whereas recall shows a modest decline, indicating a slight reduction in image diversity but with individual image quality largely preserved. Other metrics, such as FID and CLIP score, further confirm that LANTERN maintains competitive image quality despite the significant acceleration.

As a result, by allowing a degree of flexibility in token selection, our method strikes a favorable balance between speed and quality, outperforming both standard autoregressive decoding and EAGLE-2 in terms of practical efficiency. Experimental results on other visual AR models including LlamaGen-XL Stage II model and Anole (Chern et al., 2024) are presented in Appendix E. Also, component-level analysis of latency is provided in Appendix F.2.

### 5.2.2 QUALITATIVE RESULTS

To confirm that image quality is preserved with LANTERN, as indicated by the various image quality metrics, we conduct a qualitative analysis. Figure 4 demonstrates that, despite the modification of



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

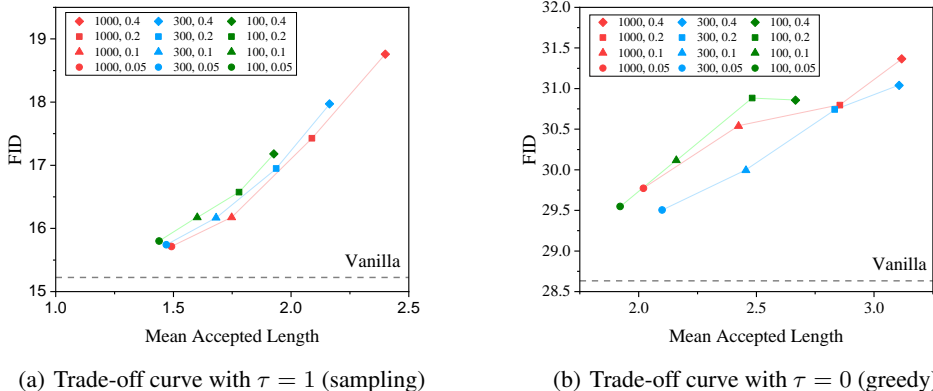


Figure 5: Trade-off curves show the relationship between performance (FID) and acceleration (mean accepted length). The results with the same  $k$  are annotated with the same color, while the same  $\delta$  values are marked with identical symbols. In the legend, the values are separated by commas, indicating  $k$  and  $\delta$ , respectively.

the target model’s probability distribution to achieve acceleration, our method effectively preserves image quality. Notably, even under the setting of  $\delta = 0.4$  and  $k = 1000$ , which achieve about  $1.64\times$  speed-up compared to the standard autoregressive decoding, generated images retain both content and style at a level comparable to standard autoregressive decoding. These qualitative results, along with the fact that LANTERN avoids significant degradation in image quality metrics as shown in the previous section, demonstrate that it effectively preserves image quality while increasing efficiency. More qualitative examples can be found in Appendix J.

### 5.2.3 TRADE-OFF BETWEEN PERFORMANCE AND EFFICIENCY

LANTERN provides various options between quality and efficiency to the end users. Therefore, we explore this trade-off across different hyperparameters by adjusting  $k$  and  $\delta$ . Figure 5 illustrates the relationship between image quality, which is measured by FID, and speed-up, which is assessed with mean accepted length for various settings of  $\delta$  and  $k$ , under both sampling ( $\tau = 1$ ) and greedy decoding ( $\tau = 0$ ). The trade-off curves highlight that increasing  $\delta$  and  $k$  generally improves speed-up, but at the expense of image quality.

The results underline that LANTERN allows for flexible tuning to balance acceleration and image quality depending on specific practical requirements. For instance, in the sampling ( $\tau = 1$ ), configurations with smaller  $\delta$  or  $k$  prioritize image quality, maintaining lower FID scores, while larger  $\delta$  or  $k$  achieve greater acceleration with a controlled loss in quality. A similar trend is observed for greedy decoding ( $\tau = 0$ ), where the trade-off is more pronounced at higher mean accepted lengths. Importantly, across all hyperparameter configurations, LANTERN demonstrates consistent and predictable trade-offs, enabling users to fine-tune the method based on their preferred balance between speed and quality. Trade-offs with other image quality metrics can be found in Appendix H.

## 5.3 ABLATION STUDY

In this section, we conduct ablation studies for LANTERN. In Section 5.3.1, we assess the impact of the metric used to measure latent proximity on performance. Then, in Section 5.3.2, we provide an ablation study on the effect of the metric for measuring the distance from the modified probability distribution.

### 5.3.1 NEAREST LATENT SELECTION

To explore the impact of various metrics for measuring latent proximity, we conduct an ablation study using representative distance metrics commonly used to measure latent proximity. Table 3 summarizes the comparisons among different strategies for selecting the nearest latent tokens, including  $\ell_2$  distance, cosine similarity, and random selection under sampling ( $\tau = 1$ ). This experiment aims to assess the role of proximity-based selection in token aggregation, with  $k = 1000$  used across all methods.

Table 3: Ablation study for latent proximity measure and probability distribution distance metrics on LlamaGen-XL Stage I model under sampling ( $\tau = 1$ ). For latent proximity measures,  $\ell_2$  distance, cosine similarity, and random selection are used, and TVD and JSD are used as probability distribution distance metrics.

Distance Metric	Latent Proximity Measure		
	Mean Accepted Length	FID	CLIP score
Cosine similarity ( $\delta = 0.2$ )	2.09	17.46	0.3206
$\ell_2$ distance ( $\delta = 0.2$ )	2.09	17.43	0.3208
$\ell_2$ distance ( $\delta = 0.05$ )	1.50	15.71	0.3203
Random ( $\delta = 0.2$ )	1.26	15.62	0.3203
Distance Metric	Probability Distribution Distance		
	Mean Accepted Length	FID	CLIP score
TVD ( $\delta = 0.3$ )	2.29	18.27	0.3206
JSD ( $\delta = 0.2$ )	2.29	18.21	0.3206
TVD ( $\delta = 0.2$ )	2.09	17.43	0.3208
JSD ( $\delta = 0.13$ )	2.09	17.48	0.3206

As shown in Table 3, the random selection significantly underperforms in terms of acceleration, achieving only a 1.26 mean accepted length, which is notably lower than both  $\ell_2$  distance and cosine similarity. Additionally, the random selection shows inferior acceleration compared to the  $\ell_2$  distance with  $\delta = 0.05$ , while maintaining a similar FID, which highlights the importance of token selection based on latent proximity. Comparing  $\ell_2$  distance and cosine similarity, both methods demonstrate nearly identical performance in terms of mean accepted length, FID, and CLIP score, suggesting robustness and flexibility in our approach to proximity measurement. These results emphasize that selecting tokens based on latent proximity is critical for achieving a balance between acceleration and image quality, offering a clear advantage over random selection. [An analysis on the size of latent proximity set is presented in Appendix G.](#)

### 5.3.2 DISTANCE BETWEEN PROBABILITY DISTRIBUTION

To ensure that the modified target distribution remains within an acceptable range of divergence from the original, we introduce  $\delta$  as an upper bound for distributional divergence. To validate the impact of the divergence metric, we evaluate two different metrics to measure the divergence: Total Variation Distance (TVD) and Jensen-Shannon Divergence (JSD). Kullback-Leibler Divergence (KLD) is not used as it is asymmetric and not a valid metric for distance in a mathematical sense.

To compare the effectiveness of these distance metrics, we fix  $k = 1000$  and adjust  $\delta$  to achieve similar mean accepted lengths across the different divergence metrics. The results shown in Table 3 demonstrate that, [although the two metrics require different  \$\delta\$  for the same level of acceleration \(measured by mean accepted length\), they produce nearly identical image quality metrics at comparable acceleration levels.](#)

These results confirm that our method consistently functions as a robust trade-off controller regardless of the chosen distance metric. Since the difference in performance between two metrics is marginal, we opt to use TVD, as it is computationally lighter and thus more efficient for large-scale implementations. [The detailed analysis on the latency differences between TVD and JSD is provided in Appendix F.3.](#)

## 6 CONCLUSIONS

In this paper, we explored the application of speculative decoding to visual AR models for the first time. We revealed that the naïve application of existing methods fails due to the token selection ambiguity problem. To address this, we proposed LANTERN, a novel relaxed acceptance condition that effectively resolves this problem. Our experiments using the state-of-the-art visual AR model and speculative decoding method demonstrated that LANTERN successfully enables speculative decoding in visual AR models, achieving substantial speed-ups with minimal compromise in image generation performance. For future work, we plan to design a drafter specifically tailored to visual AR models, aiming to achieve acceleration without sacrificing the generation performance.

## REFERENCES

- 540  
541  
542 Zachary Ankner, Rishab Parthasarathy, Aniruddha Nrusimha, Christopher Rinard, Jonathan Ragan-  
543 Kelley, and William Brandon. Hydra: Sequentially-dependent draft heads for medusa decoding,  
544 2024. URL <https://arxiv.org/abs/2402.05109>.
- 545 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
546 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
547 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
548 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
549 teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
550 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot  
551 learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-  
552 vances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Asso-  
553 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
554 2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf).
- 555 Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri  
556 Dao. Medusa: Simple LLM inference acceleration framework with multiple decoding heads. In  
557 *Forty-first International Conference on Machine Learning*, 2024. URL [https://openreview.  
558 net/forum?id=PEpbUobfJv](https://openreview.net/forum?id=PEpbUobfJv).
- 559 Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John  
560 Jumper. Accelerating large language model decoding with speculative sampling. *arXiv preprint  
561 arXiv:2302.01318*, 2023.  
562
- 563 Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever.  
564 Generative pretraining from pixels. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the  
565 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine  
566 Learning Research*, pp. 1691–1703. PMLR, 13–18 Jul 2020. URL [https://proceedings.  
567 mlr.press/v119/chen20s.html](https://proceedings.mlr.press/v119/chen20s.html).
- 568 Ethan Chern, Jiadi Su, Yan Ma, and Pengfei Liu. Anole: An open, autoregressive, native large  
569 multimodal models for interleaved image-text generation, 2024. URL [https://arxiv.org/  
570 abs/2407.06135](https://arxiv.org/abs/2407.06135).
- 571 Christoph Chuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont. Laion  
572 coco: 600m synthetic captions from laion2b-en, 2022. URL [https://laion.ai/blog/  
573 laion-coco/](https://laion.ai/blog/laion-coco/). September 27th, 2024.  
574
- 575 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,  
576 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun  
577 Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin  
578 Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang,  
579 Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny  
580 Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022. URL  
581 <https://arxiv.org/abs/2210.11416>.
- 582 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-  
583 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
584 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.  
585
- 586 Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou,  
587 Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via trans-  
588 formers. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.),  
589 *Advances in Neural Information Processing Systems*, volume 34, pp. 19822–19835. Curran Asso-  
590 ciates, Inc., 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
591 2021/file/a4d92e2cd541fca87e4620aba658316d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/a4d92e2cd541fca87e4620aba658316d-Paper.pdf).
- 592 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image  
593 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
pp. 12873–12883, 2021.

- 594 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
595 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural informa-*  
596 *tion processing systems*, pp. 2672–2680, 2014.
- 597
- 598 Jonathan Heek, Emiel Hoogeboom, and Tim Salimans. Multistep consistency models, 2024. URL  
599 <https://arxiv.org/abs/2403.06807>.
- 600
- 601 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A reference-  
602 free evaluation metric for image captioning. In Marie-Francine Moens, Xuanjing Huang, Lucia  
603 Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods*  
604 *in Natural Language Processing*, pp. 7514–7528, Online and Punta Cana, Dominican Republic,  
605 November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.  
606 595. URL <https://aclanthology.org/2021.emnlp-main.595>.
- 607
- 608 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochre-  
609 iter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In  
610 I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Gar-  
611 nett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran ASSO-  
612 ciates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2017/file/8ald694707eb0fefe65871369074926d-Paper.pdf)  
613 [2017/file/8ald694707eb0fefe65871369074926d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/8ald694707eb0fefe65871369074926d-Paper.pdf).
- 614
- 615 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on*  
616 *Deep Generative Models and Downstream Applications*, 2021. URL [https://openreview.](https://openreview.net/forum?id=qw8AKxfYbI)  
617 [net/forum?id=qw8AKxfYbI](https://openreview.net/forum?id=qw8AKxfYbI).
- 618
- 619 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In  
620 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-*  
621 *ral Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc.,  
622 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf)  
623 [file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf).
- 624
- 625 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
626 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
627 L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas  
628 Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023. URL [https://arxiv.](https://arxiv.org/abs/2310.06825)  
629 [org/abs/2310.06825](https://arxiv.org/abs/2310.06825).
- 630
- 631 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
632 adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
633 *Pattern Recognition (CVPR)*, June 2019.
- 634
- 635 Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. URL [https:](https://arxiv.org/abs/1312.6114)  
636 [//arxiv.org/abs/1312.6114](https://arxiv.org/abs/1312.6114).
- 637
- 638 Tuomas Kynk nniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved  
639 precision and recall metric for assessing generative models, 2019. URL [https://arxiv.org/](https://arxiv.org/abs/1904.06991)  
640 [abs/1904.06991](https://arxiv.org/abs/1904.06991).
- 641
- 642 Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image  
643 generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer*  
644 *Vision and Pattern Recognition*, pp. 11523–11532, 2022.
- 645
- 646 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative  
647 decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.
- 648
- 649 Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle-2: Faster inference of language  
650 models with dynamic draft trees. *arXiv preprint arXiv:2406.16858*, 2024a.
- 651
- 652 Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle: Speculative sampling requires  
653 rethinking feature uncertainty, 2024b.



- 648 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
649 Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet,  
650 Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision – ECCV 2014*, pp.  
651 740–755, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10602-1.
- 652 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*  
653 *ence on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 656 Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem,  
657 and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision,  
658 language, audio, and action. *arXiv preprint arXiv:2312.17172*, 2023.
- 660 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
661 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
662 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 663 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,  
664 and Ilya Sutskever. Zero-shot text-to-image generation, 2021. URL <https://arxiv.org/abs/2102.12092>.
- 667 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
668 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-*  
669 *ence on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- 670 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
671 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
672 text-to-image diffusion models with deep language understanding. *Advances in neural information*  
673 *processing systems*, 35:36479–36494, 2022.
- 674 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion  
675 distillation, 2023. URL <https://arxiv.org/abs/2311.17042>.
- 677 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. URL  
678 <https://arxiv.org/abs/2010.02502>.
- 680 Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan.  
681 Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint*  
682 *arXiv:2406.06525*, 2024.
- 683 Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint*  
684 *arXiv:2405.09818*, 2024.
- 686 Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling:  
687 Scalable image generation via next-scale prediction. *arXiv preprint arXiv:2404.02905*, 2024.
- 688 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
689 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand  
690 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language  
691 models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- 693 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*  
694 *neural information processing systems*, 30, 2017.
- 696 Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li.  
697 Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image  
698 synthesis, 2023. URL <https://arxiv.org/abs/2306.09341>.
- 699 Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong  
700 Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan.  
701 *arXiv preprint arXiv:2110.04627*, 2021.

702 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,  
703 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-  
704 rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2(3):5, 2022.  
705  
706 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
707 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
708 chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.  
709  
710 Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat, Aditya Krishna Menon, Afshin Rostamizadeh, San-  
711 jiv Kumar, Jean-François Kagy, and Rishabh Agarwal. Distillspec: Improving speculative decoding  
712 via knowledge distillation. In *The Twelfth International Conference on Learning Representations*.  
713 GK Zipf. The psycho-biology of language: an introduction to dynamic philology. 1935.  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## APPENDIX

## A RELATED WORKS

**Visual auto-regressive models** The development of auto-regressive (AR) models in language modeling (Brown et al., 2020; Touvron et al., 2023; Jiang et al., 2023) has paved the way for AR-based approaches in generative vision tasks, where images are tokenized into discrete tokens on a 2D grid and processed as a unidirectional sequence. Early studies (Esser et al., 2021; Yu et al., 2021; Lee et al., 2022) explored different sequence orders such as row-major raster scans, spirals, and z-curves for image generation. However, these approaches often suffer from computational inefficiencies and underperform compared to diffusion models (Ho et al., 2020; Song et al., 2022; Rombach et al., 2022).

Building upon fundamental mechanism of visual AR models from DALL-E (Ramesh et al., 2021) and CogView (Ding et al., 2021), where images are discretized into tokens using vector-quantized VAEs (VQVAE) (Van Den Oord et al., 2017) or VQGANs (Esser et al., 2021) for next-token prediction, recently, visual AR models have emerged as strong competitors to diffusion models. For instance, Parti (Yu et al., 2022) and LlamaGen (Sun et al., 2024) utilize a verbose encoder-decoder architecture that incorporates frozen T5 (Raffel et al., 2020) text features via cross-attention or prefix-filling methods, drawing on insights from Imagen (Saharia et al., 2022). More recently, Chameleon (Team, 2024) has proposed a unified AR framework that enables fully token-based representations for both image and text modalities. Additionally, model like Anole (Chern et al., 2024), which build on Chameleon, further enhance image generation capabilities, demonstrating versatility across diverse tasks. In this work, we conduct experiments on state-of-the-art visual AR models—LlamaGen and Anole—for text-to-image generation to validate the effectiveness of our acceleration scheme.

**Speculative decoding for AR models** The foundational concept of speculative decoding was first introduced by Leviathan et al. (2023) to address the sequential constraints of the AR framework in language modeling. Unlike the standard AR modeling, which generates one token per forward pass, speculative decoding aims to generate multiple tokens in a single forward pass by speculating a series of tokens, referred to as *draft* tokens. To generate draft tokens efficiently, most speculative decoding methods rely on a separate, typically smaller, model called the *drafter* model. Since generating draft tokens with the drafter model is much faster than the target model, this approach can reduce the overall token generation time.

In Leviathan et al. (2023), the drafter model is a smaller version within the same architecture family, trained on similar data and objectives as the target model. An alternative approach, DistillSpec (Zhou et al.), uses knowledge distillation to develop efficient drafter models. Recognizing that the latency of generating draft tokens is crucial for accelerating, recent works have explored lightweight designs for drafter models. Medusa (Cai et al., 2024), for instance, employs multiple separate language model heads as drafter models instead of fully independent models and leverages tree drafting and decoding to enhance the chance of accepting draft tokens. Hydra (Ankner et al., 2024) builds on Medusa by incorporating sequential dependency, yielding additional acceleration benefits. EAGLE (Li et al., 2024b;a) extends these approaches by introducing feature uncertainty in drafting and utilizing a single decoder layer as drafter models to incorporate better sequential dependency. Furthermore, they have achieved state-of-the-art performance in speculative decoding for language modeling by dynamically refining the tree structure for tree drafting.

## B EXPERIMENTAL DETAILS

### B.1 EXPERIMENTS FOR MOTIVATING EXAMPLES

For the experiment on naïve application shown in Figure 2(a), we utilize LlamaGen-L (Sun et al., 2024), an MLP with two linear layers, and a single decoder layer as the drafters for Speculative Decoding (Leviathan et al., 2023), Medusa (Cai et al., 2024), and EAGLE-2 (Li et al., 2024a), respectively. To facilitate Speculative Decoding, which requires a smaller-sized model with the same architecture as the main model and trained on the same dataset, we employ a class-conditioned LlamaGen-3B model instead of text-conditioned models. This is because text-conditioned models do not provide multiple model sizes. For each ImageNet class (Deng et al., 2009), we generate 100 images to measure the mean accepted length. Results on Vicuna-7B (Zheng et al., 2023) are taken directly from EAGLE-2.

The drafter test accuracies presented in Figure 2(b) correspond to the test accuracies of the drafters used in our main experiments (Table 2). Details on the training process for these drafters are provided in the following section. Additionally, the test accuracies of the learned drafter for Vicuna-13B are sourced from EAGLE-2.

For the experiments on average top-1 and top-10 probabilities shown in Figure 2(c), we generate 80 responses from Vicuna-7B on the MT-Bench dataset (Zheng et al., 2023) and 1000 images from LlamaGen-XL and Anole based on MS-COCO validation captions (Lin et al., 2014).

### B.2 EXPERIMENTS FOR MAIN RESULTS

To train the text-conditional model’s drafter, we sampled 100k images in LAION-COCO dataset (Chuhmann et al., 2022), which is used to train Stage I target model. We used the same amount of image sampled in ImageNet (Deng et al., 2009) dataset to train the class-conditional model’s drafter. For Anole (Chern et al., 2024), we use 118K images sampled with target model by MS-COCO (Lin et al., 2014) train caption to train drafter. We used a single-layer decoder with the same structure as the target model in the same manner as EAGLE. During training, 5% of data is set to be held out validation dataset.

Since LlamaGen (Sun et al., 2024) and Anole (Chern et al., 2024) use classifier-free guidance (Ho & Salimans, 2021) to generate images, we trained our drafter to both learn conditioned input and null-conditioned input. To do so, we dropped 10% of conditional embedding during training, as same as target model training. The batch size is 16, and the base learning rate is  $10^{-4}$ . AdamW (Loshchilov & Hutter, 2019) optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$  is used, and Linear learning rate scheduling with warm-up is used with 2000 warm-up steps. We select the best-performing model in terms of top-3 accuracy in the hold-out validation set for 20 epochs. In addition, Flan-T5 XL (Chung et al., 2022) is used to encode input text for text-conditional generation.

For LlamaGen (Sun et al., 2024) stage I and stage II, images are generated using a classifier-free guidance scale of 7.5 with top-p set to 1.0 and top-k set to 1000, which is the default generation configuration of LlamaGen official implementation for text-conditional image generation. For a class-conditional generation, the classifier-free guidance scale is set to 4.0, with the top-k sampling covering the entire vocabulary and the top-p sampling set to 1.0. For Anole (Chern et al., 2024), we use a classifier-free guidance scale of 3.0 with with top-k as 2000. For EAGLE-2 and our method, 60 candidate tokens are passed into the target model for each verification process.



## C DETAILED ANALYSIS ON LATENT PROXIMITY

### C.1 STATISTICAL ANALYSIS ON LATENT PROXIMITY

Table 4: Impact of replacing tokens with one of the  $k$ -th nearest tokens on image quality. FID and CLIP Score indicate degradation in image quality as  $k$  increases.

Randomly Replaced by one of $k$ -th nearest token	FID	CLIP Score
Vanilla AR	25.06	0.3214
$k = 50$	26.88	0.3120
$k = 100$	30.76	0.3091
$k = 1000$	88.03	0.2715

Table 4 provides statistical evidence supporting our earlier qualitative observations (Figure 3) that token replacement does not lead to significant degradation in image quality, particularly for smaller values of  $k$ . This experiment was conducted using the LlamaGen-XL Stage I model (Sun et al., 2024) and MS-COCO 2017 validation captions (Lin et al., 2014), with FID and CLIP Score used as evaluation metrics.

As  $k$  increases, the replaced token is chosen from a larger set of latent space neighbors, and its impact on image quality becomes more evident. For  $k = 50$ , the FID increases slightly from 25.06 (Vanilla AR) to 26.88, and the CLIP Score decreases marginally from 0.3214 to 0.3120, indicating minimal degradation. Similarly, for  $k = 100$ , the FID increases to 30.76, and the CLIP Score drops to 0.3091, showing that even with  $k = 100$ , the image quality remains relatively stable and acceptable. However, for  $k = 1000$ , a substantial decline is observed, with the FID increasing sharply to 88.03 and the CLIP Score dropping to 0.2715, underscoring the negative impact of replacing tokens with more distant neighbors.

These results corroborate our earlier qualitative findings, showing that token replacement up to  $k = 100$  maintains reasonable image quality, making it a viable approach in generative tasks. This demonstrates the robustness of the model under controlled token replacement scenarios.

### C.2 COMPARISON OF LATENT PROXIMITY IN VISUAL AND LANGUAGE AR MODELS

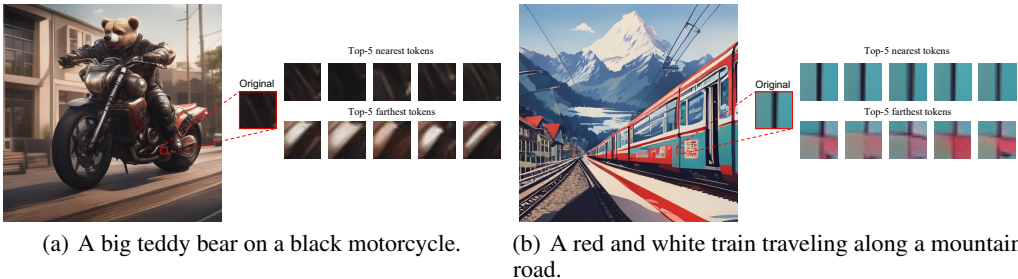


Figure 6: Visualization of top-5 nearest token and top-5 farthest token from original token in LlamaGen-XL Stage II model (Sun et al., 2024).

To achieve a deeper understanding of latent proximity, we conducted additional analyses to examine the tendencies of nearest tokens in visual AR models and language models using LlamaGen-XL Stage II (Sun et al., 2024) and Vicuna-7B (Zheng et al., 2023). For the visual AR model, we visualized the five nearest tokens and the five farthest tokens based on their L2 distances in the latent space, as shown in Figure 6. However, for the language model, since each token corresponds to subwords, it is challenging to measure proximity directly using the token itself. Therefore, we conducted the same analysis using the  $\ell_2$  distance of the input embeddings, with the results presented in Table 5.

Table 5: Top-5 nearest and farthest tokens from the original tokens based on  $\ell_2$  distance of input embeddings in Vicuna-7B (Zheng et al., 2023). [OBJ] indicates the Unicode object replacement character.

Top- $k$	Token: hi		Token: act	
	Nearest	Farthest	Nearest	Farthest
1	_Portály	\n	_Mediabestanden	[OBJ]
2	_Mediabestanden	[OBJ]	oreferrer	_Bruno
3	<0x4C>	_infinitely	<0x24>	_Ernst
4	<0x6B>	_firewall	<0x71>	_Santos
5	<0x49>	_sooner	<0x54>	_firewall

This result indicates that in the visual AR model, tokens with high latent proximity are decoded into image patches with similar visual appearances, whereas tokens with low latent proximity correspond to image patches with distinctly different appearances. On the other hand, in the language model, the similarity based on the input embeddings did not reveal a clear relationship between the nearest and farthest tokens. These results support the conclusion that latent proximity serves as an effective metric for identifying visually similar tokens in visual AR models.

## D ALGORITHMS

### D.1 SPECULATIVE DECODING WITH LANTERN

---

#### Algorithm 1 LANTERN

---

- 1: **Input:** Target model  $q(\cdot|\cdot)$ , draft model  $p(\cdot|\cdot)$ , initial sequence  $x_0, \dots, x_t$ , drafted sequence length  $L$ , minimum target sequence length  $T$ ,  $D_{TV}$  tolerance  $\delta > 0$ , and maximum cardinality of latent neighborhood  $k$ .
- 2: **Initialize:**  $n \leftarrow t$ .
- 3: **while**  $n < T$  **do**
- 4:   **for**  $t = 1, \dots, L$  **do**
- 5:     Sample draft autoregressively  $\tilde{x}_t \sim p(x|x_0, \dots, x_n, \tilde{x}_1, \dots, \tilde{x}_{t-1})$
- 6:   **end for**
- 7:   In parallel, compute  $L + 1$  sets of logits from drafts  $\tilde{x}_1, \dots, \tilde{x}_L$ :
 
$$q(x|x_0, \dots, x_n), q(x|x_0, \dots, x_n, \tilde{x}_1), \dots, q(x|x_0, \dots, x_n, \tilde{x}_1, \dots, \tilde{x}_L)$$
- 8:   **for**  $t = 1, \dots, L$  **do**
- 9:     Find the neighborhood  $A_{k,\delta}(\tilde{x}_t)$ .
- 10:     Sample  $r \sim U[0, 1]$  from an uniform distribution.
- 11:     **if**  $r < \min\left(1, \frac{\sum_{x \in A_{k,\delta}(\tilde{x}_t)} q(x|x_0, \dots, x_{n+t-1})}{p(\tilde{x}_t|x_0, \dots, x_{n+t-1})}\right)$  **then**
- 12:       Set  $x_{n+t} \leftarrow \tilde{x}_t$  and  $n \leftarrow n + 1$ .
- 13:     **else**
- 14:       Sample  $x_{n+t} \sim (q_{k,\delta}(x|x_0, \dots, x_{n+t-1}, D = \tilde{x}_t) - p(x|x_0, \dots, x_{n+t-1}))_+$  and exit the loop
- 15:     **end if**
- 16:   **end for**
- 17:   If all drafts are accepted, sample an extra token  $x_{n+L+1} \sim q(x|x_0, \dots, x_{n+L})$ .
- 18: **end while**
- 19: **Output:**  $x_{n+1}, \dots, x_{n+L}$  or  $x_{n+1}, \dots, x_{n+L+1}$

---

## D.2 PROXIMITY SET CONSTRUCTION

---

### Algorithm 2 Proximity Set Construction for LANTERN

---

**Require:** Latent space representation of tokens, number of neighbors  $k$ ,  $D_{TV}$  tolerance  $\delta$

- 1: **Precompute**  $B_k(\tilde{x})$  **for all tokens:**
  - 2: **for** each token  $\tilde{x}$  in the quantized latent space **do**
  - 3:   Compute distances between  $\tilde{x}$  and all other tokens.
  - 4:   Identify  $k$  nearest tokens (including  $\tilde{x}$  itself) based on  $\ell_2$  distance.
  - 5:   Store these tokens as the set  $B_k(\tilde{x})$ .
  - 6: **end for**
  - 7: **Dynamically calculate**  $A_{k,\delta}(\tilde{x})$  **during inference:**
  - 8: **for** each token  $\tilde{x}$  sampled during decoding **do**
  - 9:   Initialize  $A_{k,\delta}(\tilde{x}) \leftarrow \emptyset$ .
  - 10:   Initialize cumulative TVD,  $D_{TV} \leftarrow 0$ .
  - 11:   **for** each token  $x \in B_k(\tilde{x})$  (in increasing distance order) **do**
  - 12:     Compute potential TVD increment  $\Delta D_{TV}$  for adding  $x$  to  $A_{k,\delta}(\tilde{x})$ .
  - 13:     **if**  $D_{TV} + \Delta D_{TV} < \delta$  **then**
  - 14:       Add  $x$  to  $A_{k,\delta}(\tilde{x})$ .
  - 15:       Update  $D_{TV} \leftarrow D_{TV} + \Delta D_{TV}$ .
  - 16:     **else**
  - 17:       Break.
  - 18:     **end if**
  - 19:   **end for**
  - 20: **end for**
  - 21: **Return** proximity sets  $B_k$  (precomputed) and  $A_{k,\delta}$  (dynamically calculated).
- 

## E EXPERIMENTAL RESULTS ON OTHER VISUAL AR MODELS

Table 6: Actual speed-up, MAL (Mean Accepted Length), FID, CLIP score, Precision / Recall, and HPS v2 for each method on LlamaGen-XL Stage II and Anole.

Method	Acceleration ( $\uparrow$ )		Image Quality Metrics				
	Speed-up	MAL	FID ( $\downarrow$ )	CLIP score ( $\uparrow$ )	Prec ( $\uparrow$ )	Rec ( $\uparrow$ )	HPS v2 ( $\uparrow$ )
LlamaGen-XL Stage II	1.00 $\times$	1.00	47.60	0.2939	0.4138	0.5648	23.84
EAGLE-2 (Li et al., 2024a)	0.96 $\times$	1.22	-	-	-	-	-
LANTERN ( $\delta = 0.4, k = 1000$ )	<b>1.64<math>\times</math></b>	<b>2.24</b>	46.10	0.2925	0.4704	0.5222	23.06
Anole	1.00 $\times$	1.00	20.27	0.3215	0.6552	0.6398	23.52
EAGLE-2	0.73 $\times$	1.10	-	-	-	-	-
LANTERN ( $\delta = 0.5, k = 100$ )	<b>1.17<math>\times</math></b>	<b>1.83</b>	23.40	0.3186	0.6026	0.6178	22.92

Table 6 presents additional results comparing LANTERN to EAGLE-2 and standard decoding on LlamaGen-XL Stage II (Sun et al., 2024) and Anole (Chern et al., 2024). The table reports both actual speed-up (measured on RTX 3090 for LlamaGen and A100 80GB SXM for Anole) and image quality metrics, including FID, CLIP score, Precision/Recall, and HPS v2.

As shown in the Table 6, LANTERN consistently achieves superior acceleration compared to EAGLE-2 across multiple visual AR models, including LlamaGen-XL Stage II and Anole. While the performance of LANTERN varies slightly depending on the model, it consistently outperforms EAGLE-2 in terms of both speed-up and mean accepted length. For instance, on the LlamaGen-XL Stage II model, LANTERN achieves a speed-up of 1.64 $\times$  and a mean accepted length of 2.24, significantly higher than EAGLE-2’s 0.96 $\times$  speed-up and 1.22 mean accepted length. Similarly, for Anole, LANTERN achieves a speed-up of 1.17 $\times$  with a mean accepted length of 1.83, compared to EAGLE-2’s 0.73 $\times$  speed-up and 1.10 mean accepted length.

In terms of image quality, LANTERN demonstrates minimal degradation despite the improved acceleration. For the LlamaGen-XL Stage II model, the FID decreases slightly to 46.10 compared to the baseline of 47.60, while the CLIP score remains stable. Precision and recall metrics indicate that LANTERN maintains a balance between individual image quality (precision) and diversity (recall).

On Anole, LANTERN similarly exhibits competitive image quality metrics, with a CLIP score of 0.3186 and a moderate decrease in HPS v2.

These results highlight the versatility of LANTERN, showcasing its ability to provide significant acceleration benefits while maintaining acceptable image quality across diverse models. This reinforces LANTERN’s effectiveness as a practical and efficient approach for visual AR models.

## F DETAILED LATENCY ANALYSIS FOR LANTERN

### F.1 ANALYSIS ON THE NUMBER OF CAPTIONS

Table 7: Actual speedup results for LANTERN across different settings of  $\tau$ ,  $k$ , and  $\delta$ , with varying numbers of captions.

Num Captions	Actual Speed-up ( $\tau = 0, k = 1000, \delta = 0.05$ )	Actual Speed-up ( $\tau = 0, k = 1000, \delta = 0.2$ )	Actual Speed-up ( $\tau = 1, k = 1000, \delta = 0.1$ )	Actual Speed-up ( $\tau = 1, k = 1000, \delta = 0.4$ )
100	1.56×	2.33×	1.13×	1.73×
1000	1.56×	2.26×	1.13×	1.69×
2000	1.57×	2.27×	1.13×	1.69×
5000	1.56×	2.26×	1.13×	1.69×

Table 7 provides a statistical analysis of actual speed-up measurements conducted with different numbers of captions for MS-COCO validation captions (Lin et al., 2014) on LlamaGen-XL Stage I model (Sun et al., 2024). Across all configurations of  $\tau$ ,  $k$  and  $\delta$ , the actual speed-ups remain consistent beyond 1000 captions, with negligible differences observed. For example, for  $\tau = 0, k = 1000, \delta = 0.05$ , the speed-up values remain between 1.56× and 1.57× across 1000, 2000, and 5000 captions. Similarly, other configurations also exhibit minimal variation, confirming the reliability of using 1000 captions as a sample size.

This statistical stability justifies our choice to use 1000 captions in the main experiments, as it provides an accurate and computationally efficient estimate of actual speed-up without sacrificing reliability.

### F.2 COMPONENT-LEVEL LATENCY ANALYSIS FOR LANTERN

Table 8: Component-level latency analysis for LANTERN. Each component’s latency is averaged over a single draft-and-verify process, with measurements conducted on a single RTX 3090 GPU.

Method	Target Forward	Drafter Forward	Proximity Set $A$ Calculation
LANTERN( $\delta = 0.1, k = 1000$ )	$3.80 \times 10^{-2}$ s	$1.08 \times 10^{-2}$ s	$1.57 \times 10^{-3}$ s
LANTERN( $\delta = 0.4, k = 1000$ )			$1.19 \times 10^{-3}$ s

In this section, we conducted experiments to evaluate the contribution of each component of LANTERN to overall latency. Specifically, we identified three primary components that significantly impact latency: (1) the target model forward pass, (2) the drafter model forward pass, and (3) the computation of the proximity set  $A$ . Using the LlamaGen-XL Stage I model, we measured the average latency of each component during a single draft-and-verify process across 1,000 prompts from MS-COCO validation captions (Lin et al., 2014). Experiments is conducted under two settings,  $k = 1000, \delta = 0.1$  and  $k = 1000, \delta = 0.4$ , with the results summarized in Table 8.

Among the three key contributors to latency, the proximity set computation, which is introduced by LANTERN, required more time at lower  $\delta$  values. This can be attributed to the increased number of rejections at lower  $\delta$ , necessitating multiple proximity set computations within a single step. However, even in the  $k = 1000, \delta = 0.4$  setting, the latency of proximity set computation is 24× smaller than that of the target model forward passes and about 7× smaller than the drafter model forward pass.

Overall, while LANTERN introduces additional computational overhead compared to traditional speculative decoding methods, the trade-off results in significantly greater speed-up. Thus, our approach demonstrates superior actual speed-up, establishing its advantage over existing speculative decoding methods.



F.3 LATENCY COMPARISON ON PROBABILITY DISTANCE METRICS

Table 9: Computation time comparison between TVD and JSD in LANTERN. Each component’s latency is averaged over a single draft-and-verify process, with measurements conducted on a single RTX 3090 GPU. The hyperparameters for LANTERN are fixed to  $k = 1000$ , and  $\delta = 0.4$  for TVD and  $\delta = 0.2$  for JSD.

Distance Metric	Computation Time for Distance	Total Computation Time
TVD	$1.19 \times 10^{-3}$ s	$4.89 \times 10^{-2}$ s
JSD	$4.03 \times 10^{-3}$ s	$4.92 \times 10^{-2}$ s

Table 9 provides a comparison of computation times between TVD and JSD, supporting our decision to use TVD as the preferred distance metric. We use LlamaGen-XL Stage I with randomly sampled 1000 captions in MS-COCO validation captions (Lin et al., 2014) for this experiment. While both metrics produce nearly identical results in terms of mean accepted length and image quality, TVD is significantly more computationally efficient. Specifically, the computation time for TVD is  $1.19 \times 10^{-3}$  seconds, which is approximately three times faster than JSD’s computation time of  $4.03 \times 10^{-3}$  seconds.

Although the total computation time for each decoding step (including other processes) shows a smaller difference,  $4.89 \times 10^{-2}$  seconds for TVD and  $4.92 \times 10^{-2}$  seconds for JSD, this difference accumulates over multiple decoding steps. For large-scale applications or models generating long sequences, this efficiency advantage becomes increasingly significant. Therefore, given the negligible impact on quality metrics and the consistent computational advantage, we choose TVD over JSD as the more practical and efficient metric for our method.

Note that the hyperparameters do not affect the computation time for the distance metric itself but do impact the total computation time. Therefore, we set the values of  $\delta$  for TVD and JSD to achieve the same level of mean accepted length.

G SIZE OF LATENT PROXIMITY SETS ACROSS POSITIONS

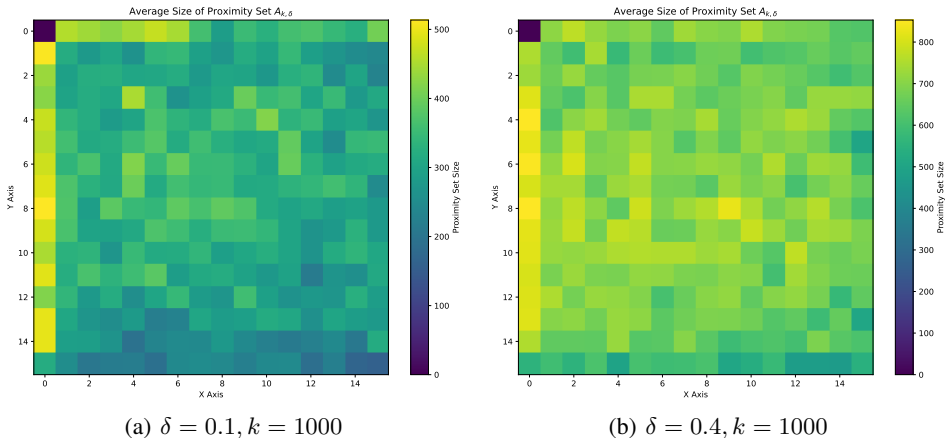


Figure 7: Average size of  $A_{k,\delta}$  on different position. Note that first token is generated right after pre-fill stage, we do not conduct speculation for first token.

We conduct experiments to better understand the behavior of LANTERN by examining how the size of the latent proximity set  $A$  varies depending on the position in images. For  $k = 1000$  and  $\delta = 0.1, 0.4$ , we generate 100 images based on MS-COCO validation captions (Lin et al., 2014) with the LlamaGen-XL Stage I model for each setting. The calculated average size of  $A$  on the different positions can be found in Figure 7.

As expected, larger  $\delta$  values generally led to a larger size of  $A$  regardless of position. Across various  $\delta$  values, the size of the proximity set  $A$  is consistently larger at the left edge of the image. We hypothesize that this phenomenon occurs because the left edge of the image corresponds to positions immediately following a line break, where uncertainty is higher compared to other positions. This increased uncertainty results in generally lower probabilities assigned to individual tokens, allowing a larger number of tokens to meet the threshold and be included in the set  $A$ .

## H TRADE-OFFS WITH ADDITIONAL METRICS

Table 10: Precision and recall (Kynkäänniemi et al., 2019) values for varying  $\delta$  and  $k$  settings. The vanilla AR decoding score is also provided for reference.

Precision / Recall	$\tau = 0$			
	$\delta = 0.05$	$\delta = 0.1$	$\delta = 0.2$	$\delta = 0.4$
<b>Vanilla AR</b>	0.4232 / 0.3517			
$k = 100$	0.4424 / 0.3212	0.4510 / 0.2989	0.4661 / 0.2754	0.4689 / 0.2813
$k = 300$	0.4488 / 0.3195	0.4619 / 0.2960	0.4696 / 0.2802	0.4750 / 0.2753
$k = 1000$	0.4484 / 0.3158	0.4659 / 0.2939	0.4771 / 0.2773	0.4854 / 0.2682
Precision / Recall	$\tau = 1$			
	$\delta = 0.05$	$\delta = 0.1$	$\delta = 0.2$	$\delta = 0.4$
<b>Vanilla AR</b>	0.4781 / 0.5633			
$k = 100$	0.4867 / 0.5389	0.4796 / 0.5303	0.4789 / 0.5140	0.4825 / 0.4946
$k = 300$	0.4856 / 0.5367	0.4834 / 0.5231	0.4894 / 0.4901	0.4895 / 0.4719
$k = 1000$	0.4865 / 0.5334	0.4869 / 0.5172	0.4880 / 0.4888	0.4909 / 0.4497

Table 11: HPS v2 (Wu et al., 2023) score for different  $\delta$  and  $k$  settings. The baseline score for Vanilla AR decoding is provided for reference.

HPS v2	$\tau = 0$			
	$\delta = 0.05$	$\delta = 0.1$	$\delta = 0.2$	$\delta = 0.4$
<b>Vanilla AR</b>	23.18			
$k = 100$	22.73	22.41	22.13	21.96
$k = 300$	22.66	22.25	21.92	21.69
$k = 1000$	22.62	22.14	21.69	21.39
HPS v2	$\tau = 1$			
	$\delta = 0.05$	$\delta = 0.1$	$\delta = 0.2$	$\delta = 0.4$
<b>Vanilla AR</b>	24.11			
$k = 100$	24.01	23.94	23.86	23.75
$k = 300$	23.97	23.85	23.70	23.55
$k = 1000$	23.91	23.75	23.47	23.22

To gain a deeper understanding of the trade-offs in LANTERN, we examined precision & recall (Kynkäänniemi et al., 2019), and HPS v2 (Wu et al., 2023) in addition to FID (Heusel et al., 2017), and the results are summarized in Table 10 and 11. The experimental settings are identical to the main experiments in Table 2 and Figure 5.

The result in Table 10 shows that LANTERN achieves comparable or slightly improved precision relative to the baseline across various settings, reflecting its ability to maintain high-quality token generation. While recall decreases marginally with increasing  $\delta$ , this is an expected trade-off due to the relaxed acceptance condition, which prioritizes sampling speed. The precision/recall results demonstrate that our method strikes a reasonable balance between quality and diversity across different hyperparameter configurations.

1188 To further quantify the aesthetic quality of generated images, we evaluate our approach using HPS v2.  
1189 Table 11 shows that while there is a slight reduction in aesthetic quality compared to the baseline, the  
1190 trade-off is well-justified by the significant improvements in generation speed. This aligns with the  
1191 intended design of LANTERN, which emphasizes efficiency while preserving acceptable quality.  
1192

1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

1242 I TEXT PROMPTS FOR QUALITATIVE SAMPLES  
1243

1244 The prompts listed below were used to generate Figure 4, with the images in the figure generated  
1245 sequentially from left to right using Prompt 1 through Prompt 9.  
1246

1247 **Prompt 1**

1248 A serene lake reflecting the colors of a sunset sky with distant mountains and a few birds  
1249 flying across the sky.  
1250

1251 **Prompt 2**

1252 A cozy bedroom with soft, warm lighting, a bed with fluffy pillows, a small bookshelf filled  
1253 with books, and a potted plant on the windowsill.  
1254

1255 **Prompt 3**

1256 A beautiful stained glass window with sunlight streaming through in a church, casting colorful  
1257 patterns on the stone floor and pews.  
1258

1259 **Prompt 4**

1260 A close-up of a single pink rose in full bloom, with soft, layered petals and gentle sunlight  
1261 illuminating its delicate curves.  
1262

1263 **Prompt 5**

1264 A classic still life of a bowl of fresh fruit, including apples, oranges, and grapes, with light  
1265 softly highlighting the textures of each fruit.  
1266

1267 **Prompt 6**

1268 A high-fashion portrait of a woman with vibrant, artistic makeup and bold accessories,  
1269 wearing a modern, avant-garde outfit against a plain studio background.  
1270

1271 **Prompt 7**

1272 A close-up of a man gazing thoughtfully out of a rain-covered window, with soft reflections  
1273 and water droplets creating a melancholic, introspective atmosphere.  
1274

1275 **Prompt 8**

1276 A close-up of a small bird perched on a snow-covered branch, with soft, fluffy feathers and a  
1277 delicate beak, illuminated by gentle winter sunlight.  
1278

1279 **Prompt 9**

1280 A close-up of a wolf with intense, focused eyes and thick gray fur, staring directly at the  
1281 camera, set against a blurred forest background.  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

1296 J QUALITATIVE RESULTS  
 1297  
 1298  
 1299



Vanilla



1330 Figure 8: Qualitative sample for the changes in the generated images according to various  
 1331  $\delta \in \{0.05, 0.1, 0.2, 0.4\}$  and  $k \in \{100, 300, 1000\}$  at  $\tau = 1$ . Input prompt is 'A kitchen with a  
 1332 refrigerator, stove and oven with cabinets'. Target model is LlamaGen-XL Stage I.  
 1333

1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349



1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

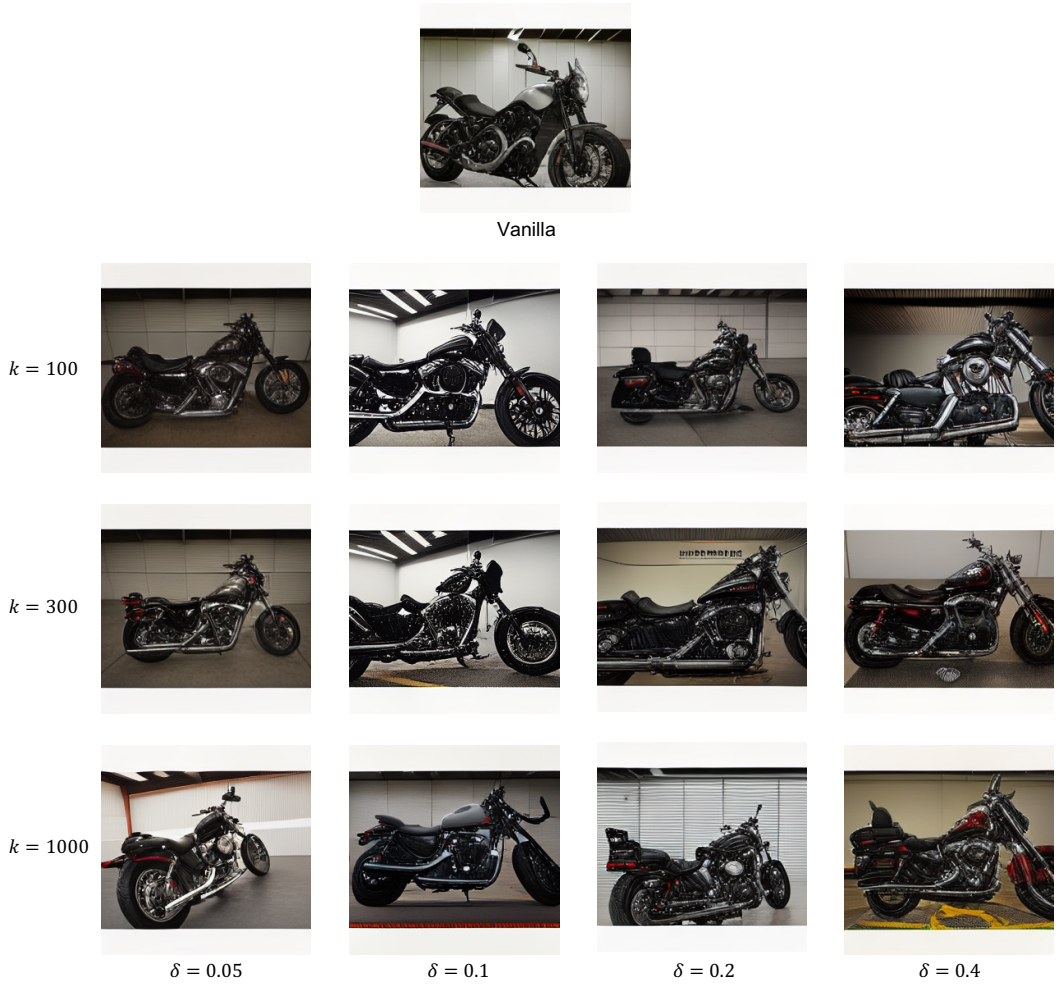
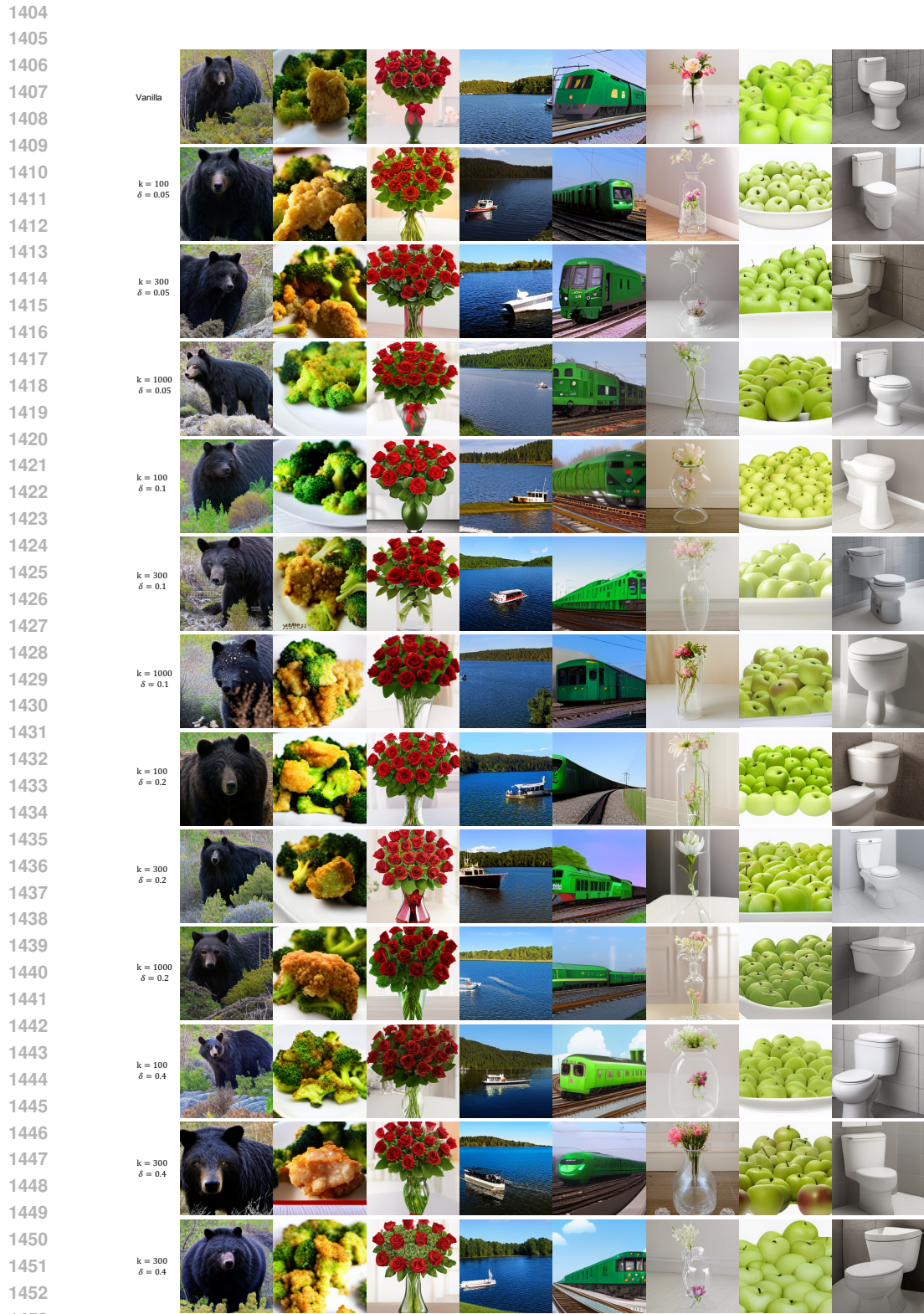


Figure 9: Qualitative sample for the changes in the generated images according to various  $\delta \in \{0.05, 0.1, 0.2, 0.4\}$  and  $k \in \{100, 300, 1000\}$  at  $\tau = 0$ . The input prompt is 'A motorcycle parked in a parking space next to another motorcycle.'. The target model is LlamaGen-XL Stage I.



1454 Figure 10: Additional qualitative samples with various  $k$  and  $\delta$  at  $\tau = 1$ . The target model is  
 1455 LlamaGen-XL Stage I, and MS-COCO validation captions are used. Images within the same column  
 1456 are generated using the same text prompt.

1457



1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

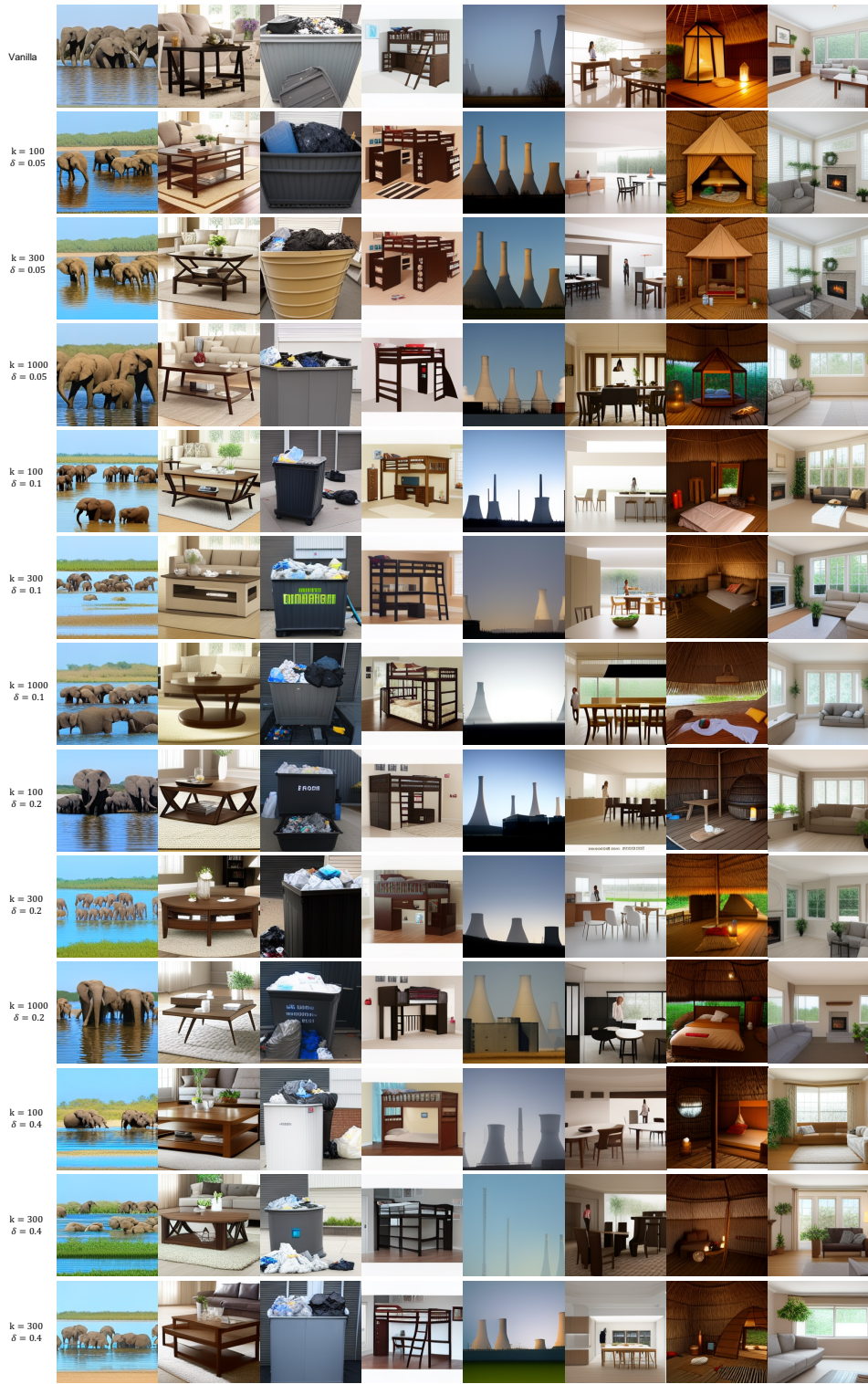


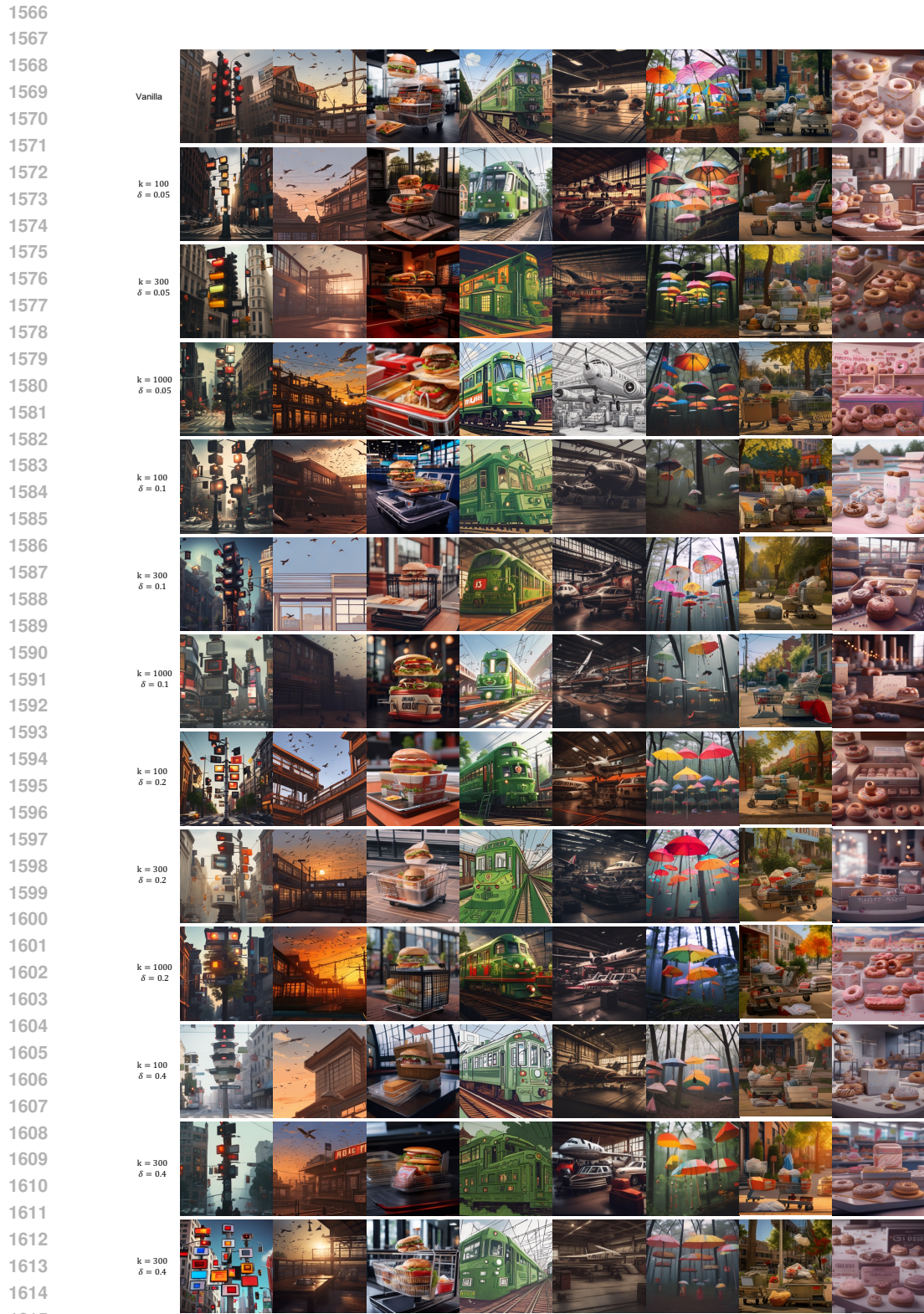
Figure 11: Additional qualitative samples with various  $k$  and  $\delta$  at  $\tau = 0$ . The target model is LlamaGen-XL Stage I, and MS-COCO validation captions are used. Images within the same column are generated using the same text prompt.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565



Figure 12: Qualitative sample for the changes in the generated images according to various  $\delta \in \{0.05, 0.1, 0.2, 0.4\}$  and  $k \in \{100, 300, 1000\}$  at  $\tau = 1$ . The input prompt is 'A pile of oranges in crates topped with yellow bananas.' The target model is LlamaGen-XL Stage II.





1616 Figure 13: Additional qualitative samples with various  $k$  and  $\delta$  at  $\tau = 1$ . The target model is  
 1617 LlamaGen-XL Stage II, and MS-COCO validation captions are used. Images within the same column  
 1618 are generated using the same text prompt.

1619