

# LADR: Locality-Aware Dynamic Rescue for Efficient Text-to-Image Generation with Diffusion Large Language Models

Anonymous ACL submission

## Abstract

Discrete Diffusion Language Models have emerged as a compelling paradigm for unified multimodal generation, yet their deployment is hindered by high inference latency arising from iterative decoding. Existing acceleration strategies often require expensive re-training or fail to leverage the 2D spatial redundancy inherent in visual data. To address this, we propose **Locality-Aware Dynamic Rescue (LADR)**, a training-free method that expedites inference by exploiting the spatial Markov property of images. LADR prioritizes the recovery of tokens at the “generation frontier”, regions spatially adjacent to observed pixels, thereby maximizing information gain. Specifically, our method integrates morphological neighbor identification to locate candidate tokens, employs a risk-bounded filtering mechanism to prevent error propagation, and utilizes manifold-consistent inverse scheduling to align the diffusion trajectory with the accelerated mask density. Extensive experiments on four text-to-image generation benchmarks demonstrate that our LADR achieves an approximate  $4\times$  **speedup** over standard baselines. Remarkably, it maintains or even enhances generative fidelity, particularly in spatial reasoning tasks, offering a state-of-the-art trade-off between efficiency and quality.

## 1 Introduction

The field of generative modeling has witnessed a paradigm shift with the rapid evolution of Discrete Diffusion Language Models (DLMs) (Sahoo et al., 2024; Nie et al., 2025; Xin et al., 2025). Unlike Autoregressive (AR) models (Radford et al., 2018; Touvron et al., 2023; Achiam et al., 2023) that generate sequences strictly left-to-right, or Continuous Diffusion Models (Ho et al., 2020; Rombach et al., 2022; Liu et al., 2023) that operate in continuous latent or pixel space, DLMs formulate generation as a bidirectional masked modeling task within a discretized vector-quantized (VQ) latent space. This

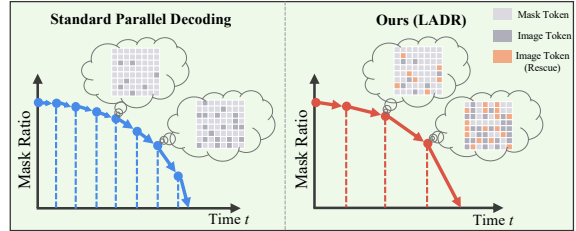


Figure 1: Comparison between Standard Parallel Decoding and our LADR method. While standard parallel decoding follows a fixed schedule, LADR accelerates decoding by exploiting spatial locality to dynamically recover neighbor tokens and keeps generation quality.

paradigm not only enables flexible, non-sequential generation orders but also facilitates unified multimodal understanding and generation within a single framework (Li et al., 2025b; You et al., 2025). By treating visual patches as discrete tokens akin to text, DLMs have achieved impressive scalability and fidelity, emerging as a powerful competitor to traditional paradigms.

However, the iterative nature of DLMs imposes a severe bottleneck on inference efficiency. High-fidelity generation typically requires 50 to 100 forward passes to progressively refine the noisy sequence. Unlike AR models that benefit from KV-caching mechanisms to reuse historical computations, masked diffusion models must re-calculate bidirectional attention interactions at every step. While acceleration techniques exist, they often fall short in practicality: distillation-based methods (Zhu et al., 2025a,b) require computationally expensive re-training and student-teacher alignment, limiting their flexibility. On the other hand, heuristic strategies borrowed from textual Masked Language Models (MLMs) (Li et al., 2025a; Ye et al., 2025) often fail to generalize to the visual domain, as they overlook the fundamental difference between 1D textual dependencies and 2D visual structures.

Our work addresses this inefficiency by exploit-

ing a property intrinsic to images but largely ignored in standard parallel decoding: *Spatial Locality*. As illustrated in Fig. 1, standard decoding schedules (Chang et al., 2022; You et al., 2025) (e.g., Cosine) assume isotropic uncertainty reduction, treating all masked tokens as independent variables. In contrast, we observe that images exhibit a strong spatial Markov property, the uncertainty of a pixel is significantly reduced if its immediate spatial neighbors are known. Based on this insight, we hypothesize that the most efficient decoding path is not random, but topological. By prioritizing the “generation frontier”, the unmasked tokens spatially adjacent to observed regions, we can accelerate the transition from noise to structure.

To materialize this insight, we propose **Locality-Aware Dynamic Rescue (LADR)**, a training-free acceleration method tailored for discrete visual generation. LADR dynamically modifies the decoding trajectory through three coupled mechanisms. First, it employs morphological operations to identify the topological neighbors of the current generation frontier. Second, to prevent the “hallucination” risks associated with aggressive acceleration, we introduce a risk-bounded filtering mechanism derived from the confidence gap of the model’s posterior. Finally, to address the distribution shift caused by rapid mask reduction, we devise a *Manifold-Consistent Inverse Scheduling* strategy that re-aligns the diffusion timesteps with the actual mask density, ensuring the denoiser operates within its trained support.

We extensively validate LADR on multiple comprehensive benchmarks, including GenEval (Ghosh et al., 2023), UniGenBench (Wang et al., 2025b), DPG-Bench (Hu et al., 2024), and T2I-CompBench (Huang et al., 2023). Experimental results demonstrate that LADR significantly outperforms standard sampling and existing acceleration baselines. Notably, our method achieves an approximate  $4\times$  **speedup** (reducing inference time from  $\sim 57s$  to  $\sim 13s$ ) without compromising generative quality. In tasks requiring spatial reasoning (e.g., object positioning and counting), LADR even surpasses the baseline performance, suggesting that enforcing spatial contiguity during decoding acts as a beneficial inductive bias.

In summary, our contributions are as follows:

- We identify *spatial locality* as a critical but underutilized source of information gain in discrete diffusion, theoretically showing that topology-

aware decoding minimizes conditional entropy more effectively than random selection.

- We propose **LADR**, a plug-and-play acceleration method that integrates morphological neighbor identification, theoretically grounded risk filtering, and inverse scheduling to safely expedite inference without re-training.
- We achieve state-of-the-art efficiency-quality trade-offs on widely adopted benchmarks, demonstrating that LADR can accelerate large-scale multimodal DLMs by  $4\times$  while maintaining robust semantic alignment and visual fidelity.

## 2 Related Work

Discrete Diffusion Language Models (DLMs) cast image generation as iterative masked token recovery in a discretized VQ space, enabling parallel decoding and substantially fewer sampling steps than continuous diffusion (Sahoo et al., 2024; Nie et al., 2025; Song et al., 2025; Arriola et al., 2025; Ho et al., 2020; Rombach et al., 2022). Initiated by MaskGIT (Chang et al., 2022), this framework has been extended by Paella (Rampas et al., 2022) and Muse (Chang et al., 2023) to improve robustness and semantic control, and more recently generalized to unified multimodal generation by modeling visual and textual tokens as a single sequence (You et al., 2025; Swerdlow et al., 2025; Xin et al., 2025; Li et al., 2025b). Despite these advances, masked discrete diffusion remains latency-bound due to its reliance on iterative refinement with bidirectional attention and dynamically changing masks, which precludes computation reuse and contrasts sharply with KV-cached autoregressive decoding (Li et al., 2024; Bai et al., 2023; Guo et al., 2025; Cai et al., 2024). Distillation-based acceleration methods compress multi-step diffusion trajectories (Hinton, 2014; Song et al., 2023; Deschenaux and Gulcehre, 2025), but adapting consistency-style objectives to discrete VQ spaces is non-trivial and typically requires expensive retraining or relaxation techniques (Zhu et al., 2025a,b). In parallel, training-free acceleration heuristics have shown promise in text diffusion and sequence models (Wu et al., 2025a; Hu et al., 2025; Wu et al., 2025b; Wang et al., 2025a; Li et al., 2025a; Israel et al., 2025), yet their direct transfer to image generation remains limited, as they fail to explicitly exploit the strong 2D spatial locality inherent in visual tokens. The extended version can be found in Appendix B.

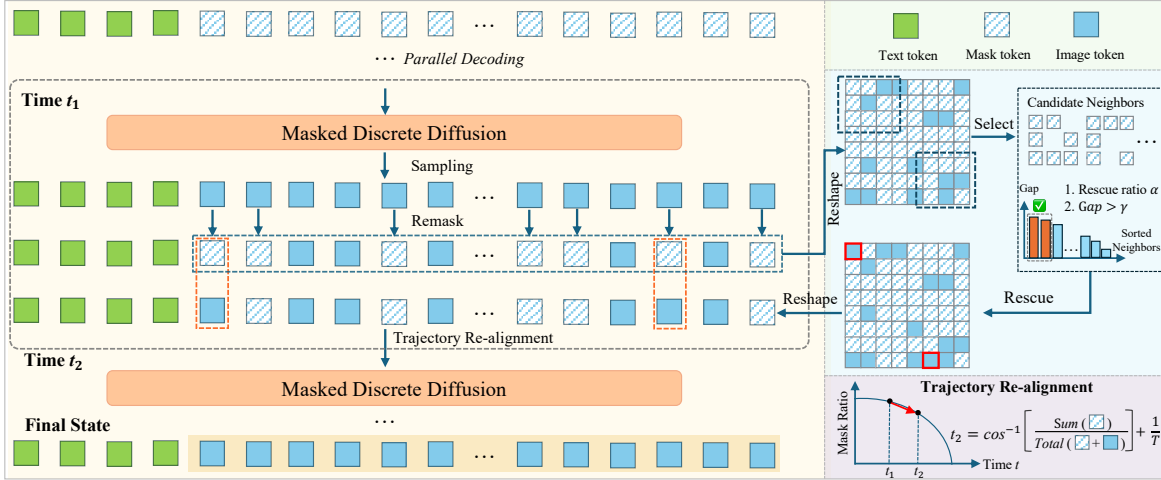


Figure 2: Overview of the LADR method. At each timestep, the flattened discrete tokens were reshaped into a 2D grid to identify candidate neighbors adjacent to resolved regions. These candidates are evaluated using the *Confidence Margin* (confidence top1-top2 gap) and are dynamically "rescued" (unmasked) based on an adaptive rescue ratio  $\alpha$  and threshold  $\gamma$ . To synchronize the generation timeline with this accelerated accumulation of tokens, the *Trajectory Re-alignment* module utilizes an inverse cosine function to re-calculate the effective timestep  $t_2$ , allowing the scheduler to skip redundant iterations while maintaining consistency.

### 3 Methodology

As illustrated in Fig. 2, we proposed **Locality-Aware Dynamic Rescue (LADR)**, a method designed to accelerate Discrete Diffusion Language Models (DLMs) while preserving generation quality. Our approach is grounded in the observation that standard parallel decoding typically treats tokens as independent variables given the global context. However, image representations derived from convolutional encoders inherently exhibit strong *spatial locality*. In this section, we first formalize the generation process and then provide the theoretical motivation grounded in information theory and risk estimation to drive our three algorithmic components: morphological neighbor identification, risk-bounded filtering, and manifold-consistent inverse scheduling.

#### 3.1 Preliminaries: Discrete Diffusion and Variational Bounds

Let  $\mathbf{z}_0 = [z_1, \dots, z_N] \in \mathcal{V}^N$  represent a discrete image sequence flattened from a  $H \times W$  feature map, where each token belongs to a codebook  $\mathcal{V}$ . The discrete diffusion process is a forward Markov chain  $q(\mathbf{z}_t | \mathbf{z}_{t-1})$  that progressively corrupts  $\mathbf{z}_0$  by replacing tokens with a special [MASK] token. The

marginal distribution at time  $t \in [0, 1]$  is given by:

$$q(\mathbf{z}_t | \mathbf{z}_0) = \prod_{i=1}^N q(z_{t,i} | z_{0,i}),$$

where  $q(z_{t,i} = [\text{M}] | z_{0,i}) = \gamma(t)$ , (1)

where  $\gamma(t)$  is a monotonic masking schedule (e.g., cosine) representing the probability of a token being masked. The reverse process  $p_\theta(\mathbf{z}_0 | \mathbf{z}_t)$  approximates the true posterior  $q(\mathbf{z}_0 | \mathbf{z}_t)$ . The training objective is to minimize the negative Evidence Lower Bound (ELBO), which simplifies to the negative log-likelihood over masked regions  $\mathcal{M}_t$ :

$$\mathcal{L} \approx \mathbb{E}_{t, \mathbf{z}_0} \left[ - \sum_{i \in \mathcal{M}_t} \log p_\theta(z_{0,i} | \mathbf{z}_{\mathcal{O}_t}) \right],$$
 (2)

where  $\mathcal{O}_t$  denotes the set of observed indices. During inference, iterative decoding approximates the joint distribution via conditional independence assumption:  $p_\theta(\mathbf{z}_0 | \mathbf{z}_t) \approx \prod_{i \in \mathcal{M}_t} p_\theta(z_i | \mathbf{z}_{\mathcal{O}_t})$ . Standard acceleration methods strictly follow  $\gamma(t)$ , discarding potentially correct predictions in early stages.

#### 3.2 Theoretical Motivation

Instead of relying on heuristic acceleration, we formulate LADR by analyzing the entropy reduction and risk bounds within the discrete latent space.

### 3.2.1 Entropy Reduction via Local Information Gain

Standard decoding assumes isotropic uncertainty reduction. However, since discrete image tokens  $\mathbf{z}$  are typically obtained via CNN-based encoders (e.g., VQGAN), the dependency between tokens decays with spatial distance due to bounded *Effective Receptive Fields (ERFs)*. We quantify the uncertainty of a masked token  $z_i$  using Conditional Entropy  $H(z_i|\mathbf{z}_O)$ . The reduction in uncertainty gained by observing an auxiliary set  $\mathcal{S}$  is quantified by the Conditional Mutual Information:

$$I(z_i; \mathbf{z}_S | \mathbf{z}_O) = H(z_i | \mathbf{z}_O) - H(z_i | \mathbf{z}_O, \mathbf{z}_S). \quad (3)$$

**Definition 1** (Generation Frontier). Given a binary mask  $\mathbf{M}$ , the generation frontier  $\mathcal{F}$  is defined as the set of masked tokens spatially adjacent to currently observed tokens:  $\mathcal{F} = \{i \mid m_i = 1 \wedge \exists j \in \mathcal{N}(i), m_j = 0\}$ , where  $\mathcal{N}(i)$  is the local spatial neighborhood.

**Proposition 1** (Locality-Driven Information Lower Bound). *Given the spatial inductive bias of the encoder, the mutual information between a token  $z_i$  and its immediate neighborhood  $\mathcal{N}(i)$  dominates that of distant context  $\mathcal{S}_{dist}$ . Formally:*

$$I(z_i; \mathbf{z}_{\mathcal{N}(i)} | \mathbf{z}_O) \gg I(z_i; \mathbf{z}_{\mathcal{S}_{dist}} | \mathbf{z}_O). \quad (4)$$

*Remark.* This proposition provides the theoretical justification for our **Morphological Neighbor Identification** strategy: by prioritizing the generation frontier  $\mathcal{F}$ , LADR maximizes the expected information gain per decoding step, guiding sampling along the local structure of the latent manifold. This locality assumption is also empirically illustrated in Fig. 3, where perturbing a small number of VQ tokens induces only spatially confined changes in the decoded image, while the global structure remains largely intact.

### 3.2.2 Safety Guarantee via Margin Bounds

Accelerating generation involves “rescuing” tokens before their scheduled timestamp. To control the quality, we must bound the probability of misclassification. We employ the *Confidence Gap*,  $\Delta_i = p_{(1)} - p_{(2)}$ , where  $p_{(1)}$  and  $p_{(2)}$  are the top-1 and top-2 probabilities.

**Theorem 1** (Margin-based Error Bound). *Consider a classification task over  $K$  classes. If the predicted distribution satisfies a confidence margin  $\Delta \geq \tau$ , the probability of error  $P(\mathcal{E})$  is strictly*

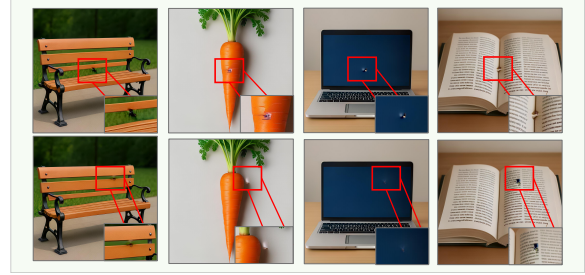


Figure 3: Visualization of localized semantic changes caused by perturbing a small set of VQ tokens. The impact remains spatially confined, supporting the locality assumption that nearby tokens dominate information gain.

*upper bounded. Specifically, in the worst-case distribution scenario:*

$$P(\mathcal{E}) \leq 1 - \left(\frac{1 + \tau}{2}\right). \quad (5)$$

*Proof.* See Appendix A.1.

*Remark.* This theorem provides a controllable mechanism. By enforcing a dynamic threshold  $\tau(t)$ , we can theoretically bound the error rate of our acceleration module, ensuring that the rescued tokens satisfy a minimum reliability standard.

### 3.2.3 Manifold Consistency via Inverse Scheduling

A critical challenge in acceleration is the **Training-Inference Mismatch**. Aggressive rescue reduces the mask ratio  $\rho_{act}$  faster than the scheduler  $\gamma(t)$  expects, potentially pushing the state out-of-distribution (OOD).

**Proposition 2** (Manifold Consistency Condition). *To ensure the input state remains within the support of the trained diffusion manifold, the conditioning timestep  $t$  must be re-aligned such that the expected mask density matches the actual observation density:*

$$t_{new} = \gamma^{-1}(\rho_{act}) \quad \text{s.t.} \quad \mathbb{E}[\rho(t_{new})] \approx \rho_{act}. \quad (6)$$

*Remark.* This necessitates our **Inverse Scheduling** technique, which acts as a temporal projection operator to correct the trajectory after aggressive rescue operations.

### 3.3 The LADR Method

Guided by the theoretical motivations above, LADR dynamically updates the mask tokens  $\mathbf{M}_t$  at each timestep  $t$  through three coupled steps.

### 3.3.1 Morphological Neighbor Identification

Leveraging Proposition 1, we aim to identify the frontier  $\mathcal{F}$ . We map the 1D mask sequence to the 2D spatial grid  $\Phi : \{0, 1\}^N \rightarrow \{0, 1\}^{H \times W}$ . Using a spatial kernel  $\mathbf{K}$  (e.g.,  $3 \times 3$ ), the candidate neighbors  $\mathcal{C}_t$  are identified via morphological dilation:

$$\mathbf{M}_{\text{grid}} = \Phi(\mathbf{M}_t), \quad (7)$$

$$\mathbf{M}_{\text{frontier}} = \mathbf{M}_{\text{grid}} \wedge (\neg \mathbf{M}_{\text{grid}} \oplus \mathbf{K}), \quad (8)$$

$$\mathcal{C}_t = \{i \mid \Phi^{-1}(\mathbf{M}_{\text{frontier}})[i] = 1\}. \quad (9)$$

This explicitly selects masked tokens that share spatial connectivity with observed regions.

### 3.3.2 Phase-Aware Dynamic Filtering

For every candidate  $i \in \mathcal{C}_t$ , we compute the confidence gap  $\Delta_i$ . Guided by Theorem 1, we employ a dynamic policy  $\Pi(t) = (\alpha_t, \tau_t)$  that adapts to the entropy of the generation phase defined by the effective timestep  $t_{\text{eff}}$ :

- **Exploration Phase** ( $t_{\text{eff}} < 0.2$ ): Global entropy is high. We apply a strict threshold  $\tau = 0.05$  and limit the rescue ratio  $\alpha = 0.1$  to prevent error propagation.
- **Structure Phase** ( $0.2 \leq t_{\text{eff}} < 0.7$ ): As semantics emerge, we relax constraints ( $\tau = 0.05, \alpha = 0.3$ ).
- **Refinement Phase** ( $t_{\text{eff}} \geq 0.7$ ): We aggressively rescue neighbors ( $\tau = \emptyset, \alpha = 1.0$ ) to fill texture details.

The set of rescued tokens  $\mathcal{R}_t$  is:

$$\mathcal{R}_t = \text{TopK}_{\Delta}(\{i \in \mathcal{C}_t \mid \Delta_i > \tau_t, \lfloor |\mathcal{C}_t| \cdot \alpha_t \rfloor\}). \quad (10)$$

### 3.3.3 Trajectory Re-alignment

After unmasking  $\mathcal{R}_t$ , the sequence sparsity decreases to  $\rho_{\text{new}}$ . Crucially, continuing with the original schedule  $t$  would violate the manifold consistency (Proposition 2). We thus re-calculate the next sampling step  $t_{\text{next}}$  using the inverse schedule:

$$t_{\text{next}} = \text{clamp}\left(\gamma^{-1}(\rho_{\text{new}}) + \frac{1}{T}, 0, 1\right). \quad (11)$$

This adjustment ensures that the noise level estimates remain accurate, effectively ‘‘skipping’’ redundant diffusion steps. The complete procedure is summarized in Algorithm 1.

### Algorithm 1 Locality-Aware Dynamic Rescue (LADR)

**Require:** Pre-trained DLM  $p_{\theta}$ , Scheduler  $\gamma(\cdot)$ , Steps  $T$

**Ensure:** Discrete tokens  $\mathbf{z}$

```

1: Initialize:  $\mathbf{z} \leftarrow [\text{MASK}]^N, \mathbf{M} \leftarrow \mathbf{1}^N$ 
2:  $N_{\text{total}} \leftarrow N$ 
3: for  $\text{step} = 0$  to  $T - 1$  do
4:   if  $\sum \mathbf{M} = 0$  then
5:     break
6:   end if
7:   {Step 1: Inverse Scheduling (Prop. 1)}
8:    $\rho_{\text{curr}} \leftarrow (\sum \mathbf{M}) / N_{\text{total}}$ 
9:    $t_{\text{eff}} \leftarrow \gamma^{-1}(\rho_{\text{curr}})$  // Align time with mask density
10:   $t_{\text{next}} \leftarrow \text{clamp}(t_{\text{eff}} + 1/T, 0, 1)$ 
11:   $n_{\text{mask}} \leftarrow \lfloor N_{\text{total}} \cdot \gamma(t_{\text{next}}) \rfloor$  // Target mask count
12:  {Step 2: Parallel Prediction}
13:   $\mathbf{L} \leftarrow p_{\theta}(\mathbf{z}, t_{\text{next}})$ 
14:   $\mathbf{P} \leftarrow \text{Softmax}(\mathbf{L})$ 
15:   $\mathbf{z}_{\text{pred}} \leftarrow \text{Sampling}(\mathbf{P})$ 
16:   $\Delta \leftarrow \text{Top1}(\mathbf{P}) - \text{Top2}(\mathbf{P})$ 
17:  {Step 3: Standard Selection (Global)}
18:   $\mathcal{I}_{\text{rank}} \leftarrow \text{Argsort}(\text{Top1}(\mathbf{P}) \cdot \mathbf{M}, \text{descending})$ 
19:   $\mathbf{M}_{\text{std}} \leftarrow \mathbf{1}^N$ 
20:   $\mathbf{M}_{\text{std}}[\mathcal{I}_{\text{rank}}[n_{\text{mask}} : N]] \leftarrow 0$  // Unmask most confident
21:  {Step 4: Neighbor Rescue (Lemma 1 & Thm 1)}
22:   $\mathbf{M}_{\text{grid}} \leftarrow \text{Reshape}(\mathbf{M}_{\text{std}}, H, W)$ 
23:   $\mathbf{M}_{\text{front}} \leftarrow \mathbf{M}_{\text{std}} \wedge \text{Flatten}(\neg \mathbf{M}_{\text{grid}} \oplus \mathbf{K}_{3 \times 3})$ 
24:   $\mathcal{C}_{\text{neigh}} \leftarrow \{i \mid \mathbf{M}_{\text{front}}[i] = 1\}$ 
25:  if  $|\mathcal{C}_{\text{neigh}}| > 0$  then
26:    Get  $\alpha, \tau$  based on  $t_{\text{eff}}$  (Sec 3.3.2)
27:     $\mathcal{C}_{\text{valid}} \leftarrow \{i \in \mathcal{C}_{\text{neigh}} \mid \Delta[i] > \tau\}$ 
28:     $k_{\text{res}} \leftarrow \min(\lfloor |\mathcal{C}_{\text{neigh}}| \cdot \alpha \rfloor, |\mathcal{C}_{\text{valid}}|)$ 
29:     $\mathcal{S}_{\text{res}} \leftarrow \text{Argsort}(\Delta[\mathcal{C}_{\text{valid}}], \text{descending})[0 : k_{\text{res}}]$ 
30:     $\mathbf{M}_{\text{std}}[\mathcal{S}_{\text{res}}] \leftarrow 0$ 
31:  end if
32:  {Step 5: State Update}
33:   $\mathbf{M} \leftarrow \mathbf{M}_{\text{std}}$ 
34:   $\mathbf{z} \leftarrow \mathbf{z}_{\text{pred}} \odot (\mathbf{1} - \mathbf{M}) + [\text{MASK}] \odot \mathbf{M}$ 
35: end for
36: return  $\mathbf{z}$ 

```

## 4 Experiments

### 4.1 Experimental setup

**Benchmarks and Baselines.** To strictly evaluate the effectiveness of our method on both inference efficiency and visual fidelity, we conducted evaluations across four publicly popular text-to-image generation benchmarks: **GenEval** (Ghosh et al., 2023), **UniGen-Bench** (Wang et al., 2025b), **DPG-Bench** (Hu et al., 2024), and **T2I-CompBench** (Huang et al., 2023). These benchmarks provide a comprehensive assessment spanning from basic object semantics to complex compositional generation. Furthermore, to ensure a fair and focused evaluation, we compared our method against two representative training-free acceleration methods: (1) **ML-Cache** (Xin et al., 2025), the native caching optimization strategy embedded in the Lumina-DiMOO backbone; and (2) **Prophet** (Li et al., 2025a), a heuristic-based ac-

celerated decoding method originally designed for text generation, which we adapted to the visual domain to investigate the cross-modal applicability of textual heuristics.

**Implementation Details.** For a fair and controllable comparison, we adopt the unified multimodal model Lumina-DiMOO (Xin et al., 2025) as the foundational DLM backbone for all experiments since comparable open-source models are limited. The generated image resolution is  $1024 \times 1024$ . To ensure a consistent evaluation of inference latency, all models and baselines were executed locally on a single NVIDIA A100 (80GB) GPU. We adhered to the standard inference configurations of the backbone model, reporting the performance following the evaluation scripts of each benchmark. Notably, for the UniGen-Bench, we used the version of their released scripts about the open-source vision-language model Qwen2.5-VL-72B (Team, 2025) to evaluate. Follow Prophet (Li et al., 2025a), we divided the parallel decoding process into three phases, and set phase-aware thresholds to regulate the rescued neighbors in the decoding process. Much like it established a proof-of-concept for heuristic-based acceleration in text decoding, our work aims to pioneer a similar trajectory for visual decoding. Consequently, we did not perform an exhaustive grid search to obtain these parameters for each specific dataset.

## 4.2 Main Results

We empirically investigated the effectiveness of our proposed accelerated method LADR by answering two critical research questions: (1) *Does the method deliver substantial speedups compared to existing caching and heuristic strategies?* (2) *Can it maintain or even enhance generative fidelity in some scenarios?*

**Decoding Efficiency Analysis.** The primary motivation of our approach is to alleviate the computational bottleneck of parallel iterative decoding. As presented in Tables 1 through 4, our method demonstrates a dramatic reduction in inference latency across all benchmarks. On average, our accelerated model completes inference in approximately  $\sim 13$ -14 seconds, representing a  $4\times$  speedup over the No-Cache ( $\sim 57$ s) and a  $2\times$  speedup over the optimized ML-Cache and Prophet ( $\sim 32$ s). Notably, our method outperforms the text-optimized Prophet algorithm, confirming that our locality-aware rescue strategy is inherently more suitable for the 2D

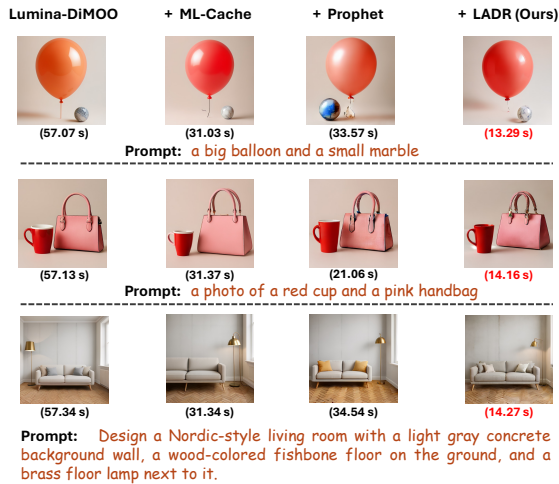


Figure 4: Qualitative comparison of different methods in terms of generation fidelity and inference time, where the corresponding text prompt is provided below each row, with the inference latency displayed in parentheses under the image.

visual domain than heuristics transplanted from 1D text generation.

**Generative Quality Analysis.** Beyond efficiency, our results indicate that the significant reduction in sampling steps does not come at the cost of visual fidelity, while it yields a highly competitive performance profile across diverse benchmarks. Conversely, we note a performance drop in fine-grained high-frequency tasks, such as the text rendering score in UniGen-Bench, suggesting that the model remains sensitive to the reduction of iterative refinement steps in some scenarios. Overall, the experiment results demonstrate that our method achieves a superior efficiency-fidelity trade-off, delivering efficient decoding speeds while preserving robust generative capabilities in image generation.

**Case Visualization.** To intuitively assess the impact of acceleration on perceptual quality, We visualized some cases generated by our proposed LADR method alongside the base model and two training-free baselines. As shown in Figure 4, our approach achieves a significant speedup ( $4\times$  faster than the backbone) while preserving intricate visual details (e.g., reflections in the "wooden boats") and correct semantic composition (e.g., "red cup and pink handbag"), demonstrating that our spatial-aware acceleration could maintain the generative quality of the underlying model.

Method	Avg. t (s)↓	Two Obj.	Colors	Attribute	Single Obj.	Position	Counting	Overall ↑
Lumina-DiMOO	57.01	<b>93.94</b>	<b>91.49</b>	73.00	97.50	79.00	<u>85.00</u>	86.66
+ ML-Cache	<u>31.95</u>	<u>93.75</u>	89.63	<b>75.50</b>	<b>100.00</b>	<u>84.50</u>	<b>85.94</b>	<b>87.83</b>
+ Prophet	32.15	91.16	86.44	70.75	96.56	74.50	84.06	83.91
+ LADR(Ours)	<b>13.22</b>	91.41	<u>91.22</u>	<u>74.75</u>	<u>99.06</u>	<b>85.50</b>	81.88	<u>87.30</u>

Table 1: Performance comparison on the **GenEval** (Ghosh et al., 2023).

Method	Avg. t (s)↓	Style	Know.	Attr.	Action	Rel.	Cmp.	Gram.	Logic.	Lay.	Text
Lumina-DiMOO	57.21	<u>91.52</u>	<b>89.87</b>	<u>79.29</u>	71.48	78.55	73.45	<b>69.79</b>	43.58	<b>85.63</b>	<b>27.87</b>
+ ML-Cache	<u>31.96</u>	91.40	<u>88.77</u>	79.17	<u>73.67</u>	<u>79.06</u>	<u>74.61</u>	<u>69.65</u>	<b>45.87</b>	83.96	26.15
+ Prophet	32.47	87.40	84.34	75.64	68.25	75.63	65.85	64.97	39.91	83.02	25.57
+ LADR(Ours)	<b>13.94</b>	<b>94.80</b>	88.61	<b>81.73</b>	<b>75.95</b>	<b>81.98</b>	<b>77.71</b>	66.31	<u>45.41</u>	<u>85.26</u>	16.38

Table 2: Performance comparison on **UniGen-Bench** (Wang et al., 2025b).

Method	Avg. t (s) ↓	Color	Shape	Texture	Spatial	Non-spatial	Complex
Lumina-DiMOO	56.92	81.07	57.02	69.30	<u>46.70</u>	<u>31.70</u>	34.98
+ ML-Cache	<u>31.78</u>	<u>81.52</u>	<u>57.75</u>	70.28	<b>46.79</b>	<b>31.83</b>	<u>35.23</u>
+ Prophet	32.43	<u>80.92</u>	55.48	<u>70.34</u>	42.37	31.46	34.92
+ LADR(Ours)	<b>13.41</b>	<b>82.25</b>	<b>58.94</b>	<b>72.34</b>	46.35	31.60	<b>36.19</b>

Table 3: Performance Comparison on **T2I-CompBench** (Huang et al., 2023).

Setting	Avg. t (s) ↓	Global	Entity	Attribute	Relation	Other	Overall ↑
Lumina-DiMOO	58.16	77.20	90.36	87.93	93.04	82.80	83.61
+ ML-Cache	<u>32.01</u>	81.46	<u>90.37</u>	<u>88.16</u>	<u>93.27</u>	<b>84.40</b>	<u>84.02</u>
+ Prophet	34.63	<u>81.76</u>	89.56	87.33	92.76	<u>83.20</u>	82.91
+ LADR(Ours)	<b>14.52</b>	<b>84.19</b>	<b>91.47</b>	<b>89.12</b>	<b>94.20</b>	81.20	<b>85.42</b>

Table 4: Performance evaluation on **DPG-Bench** (Hu et al., 2024).

### 4.3 Ablation Studies and Analysis

**Impact of Spatial Selection Strategy.** To strictly validate our hypothesis that spatial locality is the critical factor for acceleration, we conducted an ablation study on the token selection criteria. Specifically, we first determined the counts  $k$  of rescued neighbor tokens via LADR at each timestep, and then enforced this exact budget on two strategies:

- **Non-Neighbor Prioritization.** This strategy prioritizes isolated tokens with the top- $k$  confidence gaps, only reverting to neighbors if the non-neighbor set is exhausted.
- **Random Selection.** This strategy randomly sampled from the remask token, no better

neighbor or non-neighbor tokens.

Figure 5 presents the quantitative comparison on the GenEval benchmark. We observe that the **Non-Neighbor** strategy yields the lowest performance (Overall 84.05), significantly lagging behind our method (87.30). This confirms that forcing the model to resolve isolated tokens early—even those with high confidence gaps—leads to error propagation, as these predictions may lack sufficient spatial grounding. Interestingly, the random strategy (86.49) outperforms the Non-Neighbor variant but it is still lagging behind the proposed accelerated method. This indicates that our Spatial-Aware strategy provides the optimal balance, ensuring that the accelerated decoding trajectory respects the struc-

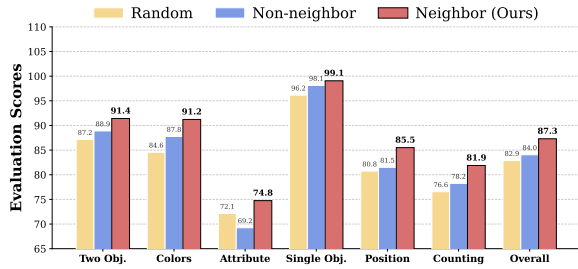


Figure 5: Ablation study of spatial selection strategies on the GenEval benchmark. “Random” means *Random Selection*, “Non-neighbor” represents *Non-Neighbor Prioritization*, and “Neighbor (Ours)” is our strategy that rescues neighbor tokens.

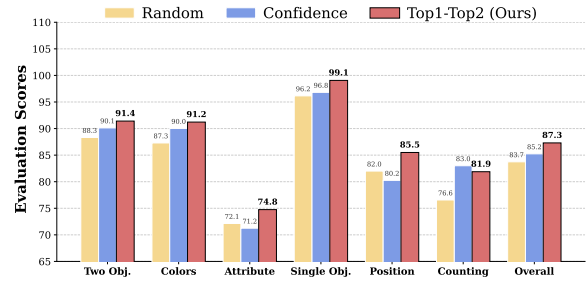


Figure 6: Ablation study about different ranking metrics of neighbor tokens on the GenEval benchmark. “Random” means *Random Neighbor Selection*, “Confidence” denotes *Standard Confidence*, and “Top1-Top2 (Ours)” means confidence margin we employed.

tural dependencies of the image.

**Effectiveness of Confidence Margin.** Besides the necessity of spatial locality, we further investigate the optimality of the ranking metric used to filter these spatial candidates. To verify whether the *Confidence Margin* (Top1-Top2 gap) provides a superior signal compared to standard confidence scores, we evaluate distinct prioritization criteria for selecting the rescued tokens  $\mathcal{R}_t$  in eq (10):

- **Standard Confidence (Top-1 Probability).** This variant ranks neighbors solely by the probability of the most likely token.
- **Random Neighbor Selection.** This baseline selects tokens stochastically from the neighborhood  $\mathcal{C}_t$  in eq (9), ignoring predictive certainty entirely.

Figure 6 reports the performance on the GenEval benchmark. The results demonstrate that our confidence margin strategy achieves the highest overall accuracy (87.30), surpassing the standard confidence baseline (85.24). While the standard confidence approach performs strongly in object-centric metrics like *Counting* (83.00), it underperforms in structural categories such as *Position* (80.25 vs. 85.50) and *Attribute* (71.25 vs. 74.75). The Random neighbor selection yields the lowest overall performance (83.74). These findings suggest that the Top1-Top2 gap is a more robust discriminator to recover tokens. It effectively penalizes ambiguously high predictions, where the model is confident in the top choice but equally confident in a competing alternative, thereby preventing the premature fixation of semantically unstable tokens.

## 5 Conclusion and Future Work

In this paper, we propose an accelerated parallel decoding strategy called LADR, which is a training-free method designed to unlock the inference efficiency of DLMs. By challenging the standard schedule generation difficulty, LADR exploits the intrinsic spatial locality of visual data. It dynamically rescues high-confidence tokens within resolved neighborhoods using a lightweight confidence margin, employing an inverse scheduling mechanism to adaptively re-align the generation timeline. Extensive evaluations across four publicly popular text-to-image generation benchmarks demonstrate that our method achieves a superior efficiency-fidelity trade-off. It delivers a significant  $4\times$  speedup over non-cached baselines and  $2\times$  speedup over heuristic-based methods, without model re-training or architectural modifications. Our findings underscore that while text-optimized heuristics provide a foundation, optimal acceleration in the visual domain requires strategies that explicitly respect the 2D spatial structure of the modality, paving the way for plug-and-play DLMs.

**Future Work.** While LADR demonstrates strong efficiency-quality trade-offs for text-to-image diffusion, several directions remain open for future exploration. First, extending locality-aware rescue to temporally structured modalities such as video generation may require jointly modeling spatial and temporal frontiers, where locality spans both space and time. Second, we anticipate that integrating LADR with emerging architectural optimizations (e.g., sparse attention or lightweight distillation) may yield complementary gains, pushing DLMs closer to real-time multimodal generation.

## 530 Limitations

531 While LADR demonstrates the potential of exploit-  
532 ing image spatial locality for acceleration of par-  
533 allel decoding, our current method still has some  
534 limitations. The implementation relies on empiri-  
535 cally determined hyperparameters, such as the con-  
536 fidence threshold  $\tau$  and rescue ratios  $\alpha$ . These  
537 values were selected to validate the core hypothesis  
538 that spatial neighbors facilitate faster convergence  
539 rather than to locate the global optimum.

## 540 References

541 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
542 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
543 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
544 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
545 cal report. *arXiv preprint arXiv:2303.08774*.

546 Marianne Arriola, Subham Sekhar Sahoo, Aaron  
547 Gokaslan, Zhihan Yang, Zhixuan Qi, Jiaqi Han,  
548 Justin T Chiu, and Volodymyr Kuleshov. 2025. [Block  
549 diffusion: Interpolating between autoregressive and  
550 diffusion language models](#). In *Proceedings of Inter-  
551 national Conference on Learning Representations*.

552 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
553 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
554 Huang, and 1 others. 2023. Qwen technical report.  
555 *arXiv preprint arXiv:2309.16609*.

556 Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu  
557 Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao  
558 Chang, Junjie Hu, and Xiao Wen. 2024. Pyra-  
559 midkv: Dynamic kv cache compression based on  
560 pyramidal information funneling. *arXiv preprint  
561 arXiv:2406.02069*.

562 Huiwen Chang, Han Zhang, Jarred Barber,  
563 AJ Maschinot, José Lezama, Lu Jiang, Ming-  
564 Hsuan Yang, Kevin Murphy, William T Freeman,  
565 Michael Rubinstein, and 1 others. 2023. Muse:  
566 Text-to-image generation via masked generative  
567 transformers. In *Proceedings of International  
568 Conference on Machine Learning*, pages 4055–4075.

569 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and  
570 William T Freeman. 2022. Maskgit: Masked gen-  
571 erative image transformer. In *Proceedings of the  
572 IEEE/CVF Conference on Computer Vision and Pat-  
573 tern Recognition*, pages 11315–11325.

574 Justin Deschenaux and Caglar Gulcehre. 2025. [Be-  
575 yond autoregression: Fast LLMs via self-distillation  
576 through time](#). In *Proceedings of International Con-  
577 ference on Learning Representations*.

578 Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig  
579 Schmidt. 2023. Geneval: An object-focused frame-  
580 work for evaluating text-to-image alignment. In *Pro-  
581 ceedings of Advances in Neural Information Process-  
582 ing Systems*, pages 52132–52152.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao  
583 Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-  
584 rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.  
585 Deepseek-r1: Incentivizing reasoning capability in  
586 llms via reinforcement learning. *arXiv preprint  
587 arXiv:2501.12948*.

Satoshi Hayakawa, Yuhta Takida, Masaaki Imaizumi,  
589 Hiromi Wakaki, and Yuki Mitsufuji. 2025. [Distilla-  
590 tion of discrete diffusion through dimensional corre-  
591 lations](#). In *Proceedings of International Conference  
592 on Machine Learning*.

G Hinton. 2014. Distilling the knowledge in a neu-  
594 ral network. In *Deep Learning and Representation  
595 Learning Workshop in Conjunction with NIPS*.  
596

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. De-  
597 noising diffusion probabilistic models. In *Proceed-  
598 ings of Advances in Neural Information Processing  
599 Systems*, pages 6840–6851.  
600

Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng,  
601 and Gang Yu. 2024. Ella: Equip diffusion models  
602 with llm for enhanced semantic alignment. *arXiv  
603 preprint arXiv:2403.05135*.  
604

Zhanqiu Hu, Jian Meng, Yash Akhauri, Mohamed S Ab-  
605 delfattah, Jae-sun Seo, Zhiru Zhang, and Udit Gupta.  
606 2025. Accelerating diffusion language model infer-  
607 ence via efficient kv caching and guided diffusion.  
608 *arXiv preprint arXiv:2505.21467*.  
609

Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and  
610 Xihui Liu. 2023. T2i-compbench: A compre-  
611 hensive benchmark for open-world compositional text-  
612 to-image generation. In *Proceedings of Advances  
613 in Neural Information Processing Systems*, pages  
614 78723–78747.  
615

Daniel Mingyi Israel, Guy Van den Broeck, and Aditya  
616 Grover. 2025. [Accelerating diffusion LLMs via adap-  
617 tive parallel decoding](#). In *Proceedings of Annual  
618 Conference on Neural Information Processing Sys-  
619 tems*.  
620

Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki  
621 Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong  
622 He, Yuki Mitsufuji, and Stefano Ermon. 2024. [Con-  
623 sistency trajectory models: Learning probability flow  
624 ODE trajectory of diffusion](#). In *Proceedings of Inter-  
625 national Conference on Learning Representations*.  
626

Pengxiang Li, Yefan Zhou, Dilxat Muhtar, Lu Yin,  
627 Shilin Yan, Li Shen, Yi Liang, Soroush Vosoughi,  
628 and Shiwei Liu. 2025a. Diffusion language models  
629 know the answer before decoding. *arXiv preprint  
630 arXiv:2508.19982*.  
631

Shufan Li, Jiuxiang Gu, Kangning Liu, Zhe Lin, Zijun  
632 Wei, Aditya Grover, and Jason Kuen. 2025b. Lavidao:  
633 Elastic large masked diffusion models for unified  
634 multimodal understanding and generation. *arXiv  
635 preprint arXiv:2509.19244*.  
636

637	Shufan Li, Konstantinos Kallidromitis, Hritik Bansal,	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	691
638	Akash Gokul, Yusuke Kato, Kazuki Kozuka, Ja-	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	692
639	son Kuen, Zhe Lin, Kai-Wei Chang, and Aditya	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	693
640	Grover. 2025c. Lavidia: A large diffusion language	Azhar, and 1 others. 2023. Llama: Open and effi-	694
641	model for multimodal understanding. <i>arXiv preprint</i>	cient foundation language models. <i>arXiv preprint</i>	695
642	<i>arXiv:2505.16839</i> .	<i>arXiv:2302.13971</i> .	696
643	Yuhong Li, Yingbing Huang, Bowen Yang, Bharat	Xu Wang, Chenkai Xu, Yijie Jin, Jiachun Jin, Hao	697
644	Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai,	Zhang, and Zhijie Deng. 2025a. Diffusion llms can	698
645	Patrick Lewis, and Deming Chen. 2024. Snapkv:	do faster-than-ar inference via discrete diffusion forc-	699
646	Llm knows what you are looking for before genera-	ing. <i>arXiv preprint arXiv:2508.09192</i> .	700
647	tion. In <i>Proceedings of Advances in Neural Informa-</i>	Yibin Wang, Zhimin Li, Yuhang Zang, Yujie Zhou, Jiazi	701
648	<i>tion Processing Systems</i> , pages 22947–22970.	Bu, Chunyu Wang, Qinglin Lu, Cheng Jin, and Jiaqi	702
649	Xingchao Liu, Chengyue Gong, and 1 others. 2023.	Wang. 2025b. Pref-grpo: Pairwise preference reward-	703
650	Flow straight and fast: Learning to generate and	based grpo for stable text-to-image reinforcement	704
651	transfer data with rectified flow. In <i>The Eleventh</i>	learning. <i>arXiv preprint arXiv:2508.20751</i> .	705
652	<i>International Conference on Learning Representa-</i>	Chengyue Wu, Hao Zhang, Shuchen Xue, Shizhe Diao,	706
653	<i>tions</i> .	Yonggan Fu, Zhijian Liu, Pavlo Molchanov, Ping	707
654	Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang,	Luo, Song Han, and Enze Xie. 2025a. Fast-dllm	708
655	Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong	v2: Efficient block-diffusion llm. <i>arXiv preprint</i>	709
656	Wen, and Chongxuan Li. 2025. Large language dif-	<i>arXiv:2509.26328</i> .	710
657	fusion models. <i>arXiv preprint arXiv:2502.09992</i> .	Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu,	711
658	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya	Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han, and	712
659	Sutskever, and 1 others. 2018. Improving language	Enze Xie. 2025b. Fast-dllm: Training-free accel-	713
660	understanding by generative pre-training.	eration of diffusion llm by enabling kv cache and	714
661	Dominic Rampas, Pablo Pernias, and Marc Aubre-	parallel decoding. <i>arXiv preprint arXiv:2505.22618</i> .	715
662	ville. 2022. A novel sampling scheme for text-and	Yi Xin, Qi Qin, Siqi Luo, Kaiwen Zhu, Juncheng	716
663	image-conditional image synthesis in quantized la-	Yan, Yan Tai, Jiayi Lei, Yuewen Cao, Keqi Wang,	717
664	tent spaces. <i>arXiv preprint arXiv:2211.07292</i> .	Yibin Wang, and 1 others. 2025. Lumina-dimoo:	718
665	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	An omni diffusion large language model for multi-	719
666	Patrick Esser, and Björn Ommer. 2022. High-	modal generation and understanding. <i>arXiv preprint</i>	720
667	resolution image synthesis with latent diffusion mod-	<i>arXiv:2510.06308</i> .	721
668	els. In <i>Proceedings of the IEEE/CVF Conference</i>	Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui	722
669	<i>on Computer Vision and Pattern Recognition</i> , pages	Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong.	723
670	10684–10695.	2025. Dream 7b: Diffusion large language models.	724
671	Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron	<i>arXiv preprint arXiv:2508.15487</i> .	725
672	Gokaslan, Edgar Marroquin, Justin Chiu, Alexan-	Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli	726
673	der Rush, and Volodymyr Kuleshov. 2024. Simple	Shechtman, Fredo Durand, William T Freeman, and	727
674	and effective masked diffusion language models. In	Taesung Park. 2024. One-step diffusion with dis-	728
675	<i>Proceedings of Advances in Neural Information Pro-</i>	tribution matching distillation. In <i>Proceedings of</i>	729
676	<i>cessing Systems</i> , pages 130136–130184.	<i>the IEEE/CVF Conference on Computer Vision and</i>	730
677	Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya	<i>Pattern Recognition</i> , pages 6613–6623.	731
678	Sutskever. 2023. <b>Consistency models</b> . In <i>Proceed-</i>	Zebin You, Shen Nie, Xiaolu Zhang, Jun Hu, Jun Zhou,	732
679	<i>ings of International Conference on Machine Learn-</i>	Zhiwu Lu, Ji-Rong Wen, and Chongxuan Li. 2025.	733
680	<i>ing</i> , pages 32211–32252.	Llada-v: Large language diffusion models with visual	734
681	Yuxuan Song, Zheng Zhang, Cheng Luo, Pengyang Gao,	instruction tuning. <i>arXiv preprint arXiv:2505.16933</i> .	735
682	Fan Xia, Hao Luo, Zheng Li, Yuehang Yang, Hongli	Yuanzhi Zhu, Xi Wang, Stéphane Lathuilière, and Vicky	736
683	Yu, Xingwei Qu, and 1 others. 2025. Seed diffusion:	Kalogeiton. 2025a. Di [m] o: Distilling masked diffu-	737
684	A large-scale diffusion language model with high-	sion models into one-step generator. In <i>Proceedings</i>	738
685	speed inference. <i>arXiv preprint arXiv:2508.02193</i> .	<i>of the IEEE/CVF International Conference on Com-</i>	739
686	Alexander Swerdlow, Mihir Prabhudesai, Siddharth	<i>puter Vision</i> , pages 18606–18618.	740
687	Gandhi, Deepak Pathak, and Katerina Fragkiadaki.	Yuanzhi Zhu, Xi Wang, Stéphane Lathuilière, and Vicky	741
688	2025. Unified multimodal discrete diffusion. <i>arXiv</i>	Kalogeiton. 2025b. Soft-di [m] o: Improving one-	742
689	<i>preprint arXiv:2503.20853</i> .	step discrete image generation with soft embeddings.	743
690	Qwen Team. 2025. <b>Qwen2.5-vl</b> .	<i>arXiv preprint arXiv:2509.22925</i> .	744

## A Theoretical Proofs and Derivations

In this section, we provide the detailed mathematical derivations for the propositions and theorems presented in the main methodology.

### A.1 Proof of Theorem 1 (Margin-based Error Bound)

**Problem Statement:** Let  $\mathbf{p} = [p_1, p_2, \dots, p_K]$  be the probability distribution over  $K$  classes, sorted such that  $p_{(1)} \geq p_{(2)} \geq \dots \geq p_{(K)}$ . The predicted class is  $\hat{y} = \operatorname{argmax}_k p_k$ . The probability of error is  $P(\mathcal{E}) = 1 - p_{(1)}$ . Given the margin constraint  $p_{(1)} - p_{(2)} \geq \tau$ , we seek the upper bound of  $P(\mathcal{E})$ .

*Proof.* We aim to find the maximum possible error, which corresponds to minimizing the confidence  $p_{(1)}$  subject to the constraints.

1. **Define Constraints:** The probability distribution must sum to 1, and the margin condition must hold:

$$\sum_{k=1}^K p_{(k)} = 1, \quad (12)$$

$$p_{(1)} - p_{(2)} \geq \tau. \quad (13)$$

2. **Worst-Case Analysis:** To maximize error ( $1 - p_{(1)}$ ), we must minimize  $p_{(1)}$ . This occurs when the remaining probability mass  $1 - p_{(1)}$  is concentrated as much as possible in the ‘‘competitor’’ classes to satisfy the constraints tightest.

From Eq. (13), we express  $p_{(2)}$  in terms of  $p_{(1)}$ :

$$p_{(2)} \leq p_{(1)} - \tau. \quad (14)$$

3. **Bounding the Sum:** Since  $p_{(1)}$  is the maximum, all other probabilities  $p_{(k)}$  for  $k > 2$  are also bounded by  $p_{(2)}$  (and thus  $p_{(1)} - \tau$ ), but a stricter bound for the binary case (or focusing on the top-2 dominating case) gives the tightest constraint on  $p_{(1)}$ .

Consider the relation between the top-2 probabilities and the total sum:

$$p_{(1)} + p_{(2)} \leq \sum_{k=1}^K p_{(k)} = 1. \quad (15)$$

4. **Deriving the Inequality:** Substitute the margin constraint  $p_{(2)} \leq p_{(1)} - \tau$  into the sum inequality:

ity:

$$p_{(1)} + (p_{(1)} - \tau) \leq p_{(1)} + p_{(2)} \leq 1 \quad (16)$$

$$2p_{(1)} - \tau \leq 1 \quad (17)$$

$$2p_{(1)} \leq 1 + \tau \quad (18)$$

$$p_{(1)} \leq \frac{1 + \tau}{2}. \quad (19)$$

*Correction for Lower Bound of Correctness:* Wait, we are looking for the *minimum* valid  $p_{(1)}$  to find the *maximum* error. The constraint is  $p_{(1)} + p_{(2)} + \sum_{k=3}^K p_{(k)} = 1$ . To minimize  $p_{(1)}$ , we need to maximize the other terms. The largest possible value for  $p_{(2)}$  is  $p_{(1)} - \tau$ .

Thus, in the worst case (where probability is split maximally between top-2):

$$p_{(1)} + (p_{(1)} - \tau) \approx 1 \quad (\text{assuming } \sum_{k=3}^K p_{(k)} \approx 0) \quad (20)$$

$$2p_{(1)} \approx 1 + \tau \quad (21)$$

$$p_{(1)} \geq \frac{1 + \tau}{2}. \quad (22)$$

Any  $p_{(1)}$  smaller than  $\frac{1+\tau}{2}$  would require  $p_{(2)} > p_{(1)} - \tau$  to sum to 1, violating the margin condition.

Therefore, the lower bound on confidence is:

$$p_{(1)} \geq \frac{1 + \tau}{2}. \quad (23)$$

5. **Calculating the Error Bound:** The probability of error is the complement of the top-1 probability:

$$P(\mathcal{E}) = 1 - p_{(1)} \quad (24)$$

$$\leq 1 - \frac{1 + \tau}{2} \quad (25)$$

$$= \frac{1 - \tau}{2}. \quad (26)$$

This concludes the proof. The error is strictly bounded by a linear function of the threshold  $\tau$ .  $\square$

### A.2 Justification of Proposition 1 (Locality-Induced Information Gain)

**Proposition Statement:** The mutual information  $I(z_i; \mathbf{z}_{\mathcal{N}(i)})$  dominates  $I(z_i; \mathbf{z}_{\mathcal{S}_{dist}})$ .

*Proof.* We derive this property from the definition of Conditional Mutual Information and the Markov property of Convolutional Neural Networks (CNNs).

820 **1. Definition of Information Gain:** Let  $\Omega$  be  
 821 the set of all tokens. The information gain for a  
 822 target token  $z_i$  given a subset  $\mathcal{S}$  is:

$$823 \begin{aligned} IG(\mathcal{S}) &= I(z_i; \mathcal{S} | \Omega \setminus (\{z_i\} \cup \mathcal{S})) \\ 824 &= H(z_i | \Omega \setminus (\{z_i\} \cup \mathcal{S})) - H(z_i | \Omega \setminus \{z_i\}). \end{aligned} \quad (27)$$

825 **2. Spatial Correlation Decay:** For latent codes  
 826  $z$  derived from a VQ-GAN encoder  $E$ , the covari-  
 827 ance between features at spatial locations  $u$  and  $v$   
 828 generally follows a decay function dependent on  
 829 Euclidean distance  $d(u, v)$ :

$$830 \text{Cov}(z_u, z_v) \propto \exp\left(-\frac{d(u, v)^2}{2\sigma^2}\right), \quad (28)$$

831 where  $\sigma$  corresponds to the Effective Receptive  
 832 Field (ERF).

833 **3. Entropy and Correlation:** For Gaussian-like  
 834 distributions, the conditional entropy is related to  
 835 the correlation coefficient  $\rho$ :

$$836 H(z_i | z_j) \approx \frac{1}{2} \log(1 - \rho_{ij}^2) + \text{const.} \quad (29)$$

837 Higher correlation  $\rho_{ij}$  leads to lower conditional  
 838 entropy  $H(z_i | z_j)$ .

839 **4. Comparing Neighbors vs. Distant Tokens:**  
 840 Let  $j \in \mathcal{N}(i)$  be a spatial neighbor and  $k \in \mathcal{S}_{dist}$   
 841 be a distant token.

$$842 d(i, j) \ll d(i, k) \quad (30)$$

$$843 \Rightarrow \rho_{ij} > \rho_{ik} \quad (31)$$

$$844 \Rightarrow H(z_i | z_j) < H(z_i | z_k). \quad (32)$$

845 **5. Conclusion:** Since observing neighbors re-  
 846 duces the conditional entropy more than observing  
 847 distant tokens:

$$848 I(z_i; \mathbf{z}_{\mathcal{N}(i)}) > I(z_i; \mathbf{z}_{\mathcal{S}_{dist}}). \quad (33)$$

849 Thus, the optimal strategy for variance reduction is  
 850 to prioritize  $\mathcal{N}(i)$ .  $\square$

### 851 A.3 Derivation of Inverse Scheduling 852 (Proposition 2)

853 **Objective:** Find the effective timestep  $t_{new}$  such  
 854 that the model’s training distribution matches the  
 855 current observation density.

856 *Proof.* **1. Forward Process Definition:** The mask-  
 857 ing probability at time  $t$  is given by the schedule  
 858 function  $\gamma(t)$ :

$$859 p(z_{t,i} = [\text{MASK}]) = \gamma(t). \quad (34)$$

860 **2. Expected Mask Ratio:** For a sequence of  
 861 length  $N$ , the number of masked tokens  $M_t$  follows  
 862 a Binomial distribution. The expected mask ratio  
 863 is:

$$864 \mathbb{E}\left[\frac{|M_t|}{N}\right] = \gamma(t). \quad (35)$$

865 **3. Perturbation via Rescue:** The LADR al-  
 866 gorithm unmask a set of tokens  $\mathcal{R}$ , changing the  
 867 actual mask ratio to  $\rho_{act}$ :

$$868 \rho_{act} = \frac{|\mathcal{M}_{prev}| - |\mathcal{R}|}{N}. \quad (36)$$

869 **4. Manifold Alignment:** To ensure the input to  
 870 the denoiser  $p_\theta(\mathbf{z}_0 | \mathbf{z}_{t_{new}})$  is In-Distribution (ID),  
 871 we require the expected mask ratio at  $t_{new}$  to equal  
 872 the actual current ratio:

$$873 \gamma(t_{new}) = \rho_{act}. \quad (37)$$

874 **5. Solving for Timestep:** Assuming  $\gamma(t)$  is  
 875 monotonic and invertible (e.g., cosine schedule),  
 876 we apply the inverse function:

$$877 t_{new} = \gamma^{-1}(\rho_{act}). \quad (38)$$

878 This creates the mapping required for the Manifold  
 879 Consistent Inverse Scheduling.  $\square$

## 880 B Extended Related Work

881 In this section, we provide a detailed elaboration  
 882 on the development of Masked Discrete Diffusion  
 883 for image generation and the current landscape of  
 884 acceleration strategies.

### 885 B.1 Masked Discrete Diffusion for Image 886 Generation

887 Discrete Diffusion Language Models (DLMs) (Sa-  
 888 hoo et al., 2024; Nie et al., 2025; Song et al., 2025;  
 889 Arriola et al., 2025) have reformulated the gener-  
 890 ation process as masked modeling within a dis-  
 891 cretized vector-quantized (VQ) space. Pioneered  
 892 by MaskGIT (Chang et al., 2022), this paradigm  
 893 utilizes a bidirectional Transformer coupled with a  
 894 mask-scheduling strategy to enable image synthesis  
 895 via iterative parallel decoding. Compared to stan-  
 896 dard continuous diffusion models (Ho et al., 2020;  
 897 Rombach et al., 2022), this formulation signifi-  
 898 cantly curtails the required sampling steps. Build-  
 899 ing upon this foundation, subsequent architectures  
 900 have rapidly expanded the field: Paella (Rampas  
 901 et al., 2022) optimized U-Net backbones with noise-  
 902 robust objectives, while Muse (Chang et al., 2023)

903 demonstrated scalability by integrating pre-trained  
904 LLMs for enhanced semantic control.

905 More recently, the field has witnessed a shift  
906 towards unified multimodal understanding and gener-  
907 ation (You et al., 2025; Swerdlow et al., 2025;  
908 Xin et al., 2025; Li et al., 2025b,c). Notably, frame-  
909 works like Lumina-DiMOO (Xin et al., 2025) and  
910 LaVida-O (Li et al., 2025b) adopt a generalized dis-  
911 crete diffusion approach that treats visual and tex-  
912 tual tokens as a shared sequence. While facilitating  
913 versatile generative capabilities across modalities,  
914 the iterative mask recovery process still imposes a  
915 non-negligible computational overhead. This un-  
916 derscores the necessity for efficient acceleration  
917 strategies that can expedite inference without com-  
918 promising generative integrity.

## 919 **B.2 Acceleration of Masked Discrete** 920 **Diffusion**

921 The efficacy of discrete diffusion hinges on iterative  
922 refinement, where multiple forward passes resolve  
923 the joint distribution of tokens. Unlike autoregres-  
924 sive models that benefit from causal masking and  
925 KV-caching (Li et al., 2024; Bai et al., 2023; Guo  
926 et al., 2025; Cai et al., 2024), masked diffusion  
927 relies on bidirectional attention with dynamically  
928 shifting mask states (Sahoo et al., 2024; Xin et al.,  
929 2025). This characteristic precludes the reuse of  
930 historical computations, creating a distinct latency  
931 bottleneck.

932 **Distillation-Based Approaches.** To alleviate la-  
933 tency, research has gravitated towards model distil-  
934 lation (Hinton, 2014; Song et al., 2023; Deschenaux  
935 and Gulcehre, 2025; Yin et al., 2024; Hayakawa  
936 et al., 2025). While Consistency Models (Song  
937 et al., 2023; Kim et al., 2024) are effective in con-  
938 tinuous pixel space, adapting them to discrete VQ  
939 space requires specialized formulations due to the  
940 absence of explicit ODE trajectories (Zhu et al.,  
941 2025a). Works such as DiMO (Zhu et al., 2025a)  
942 and Soft-DiMO (Zhu et al., 2025b) address this us-  
943 ing policy gradients and soft embedding relaxations  
944 to compress multi-step trajectories. However, these  
945 methods necessitate computationally expensive re-  
946 training and student-teacher alignment, limiting  
947 their plug-and-play applicability.

948 **Training-Free Heuristics.** Advancements in  
949 DLMS have also explored architectural optimiza-  
950 tions (Wu et al., 2025a; Hu et al., 2025; Wu et al.,  
951 2025b; Wang et al., 2025a) and adaptive sam-  
952 pling (Li et al., 2025a; Israel et al., 2025) primar-

953 ily for text generation. While effective for 1D se-  
954 quences, their direct adaptation to the visual do-  
955 main is non-trivial. The inherent gap between the  
956 sequential dependencies of text and the 2D spa-  
957 tial correlations of images renders text-optimized  
958 heuristics suboptimal for visual generation. This  
959 discrepancy highlights the need for acceleration  
960 strategies explicitly tailored to the spatial redun-  
961 dancy and structural properties of images.