

# 000 TOMA: TOKEN MERGE WITH ATTENTION FOR IMAGE 001 GENERATION WITH DIFFUSION MODELS 002 003 004

005 **Anonymous authors**

006 Paper under double-blind review  
007

## 008 009 ABSTRACT 010

011 Diffusion is one of the leading approaches for image generation. Plug-and-play  
012 token merge techniques have recently been introduced to mitigate the high com-  
013 putation cost of transformer blocks in diffusion models. However, existing meth-  
014 ods overlook two key factors: (1) the token selection process fails to account  
015 for relationships among tokens, potentially discarding important information and  
016 limiting image quality; 2) they do not take advantage of the modern, efficient  
017 implementation of attention, so that, the overhead backfires the achieved algo-  
018 rithmic efficiency. In this paper, we propose Token Merge with Attention (ToMA)  
019 with three major improvements. Firstly, we utilize a submodular-based token se-  
020 lection method to identify diverse tokens as merge destinations, representative of  
021 the entire token set. Secondly, we use efficient attention implementation for the  
022 merge operation with negligible overhead. Also, we formalize the (un-)merge  
023 as (inverse-)linear transformations, allowing shareable computation across lay-  
024 ers/iterations. Finally, we utilize the image locality to further accelerate the com-  
025 putation by performing all the operations on tokens in local tiles. ToMA achieves  
026 the best trade-offs between speed-ups and generation quality compared to the  
027 baselines.

## 028 1 INTRODUCTION 029 030



031  
032  
033  
034  
035  
036  
037  
038  
039  
040 Figure 1: Variants of ToMA generated images: ToMA\_stripe, ToMA, ToMA\_tile, ToMA

041 Diffusion Ho et al. (2020); Song et al. (2021); Dhariwal & Nichol (2021) emerges as one of  
042 the leading approaches for high-quality image generation. However, the increasing complexity  
043 of diffusion models, driven by their core transformer-based architecture, presents significant  
044 computational challenges. The design of the transformer leads to quadratic complexity with respect  
045 to the number of tokens, making them inefficient and resource-intensive as token counts increase.  
046

047 Methods with different approaches have been developed to mitigate this issue. Flash Attention Dao  
048 et al. (2022); Dao (2023) introduces a more efficient attention mechanism that reduces memory  
049 overhead, while xformers Lefaudeux et al. (2022) utilize sparse attention to lower memory usage  
050 and improve scalability. Methods like Token Pruning Kim et al. (2022) reduce computation by  
051 eliminating less relevant tokens during inference, albeit at the cost of potential quality degradation.

052 ToMeSD Bolya & Hoffman (2023) leverages token merging Bolya & Hoffman (2023); Kim et al.  
053 (2023), consolidating similar tokens during the forward pass to reduce the token count in compu-  
tational layers, thereby lowering complexity without requiring network retraining. Essentially, the

original sequence of tokens gets merged before each layer in the transformer block, including attentions and MLPs, and after the computation finishes, an unmerge is applied to transform the merged tokens back to the original sequence length. Token merge shares the spirit with token pruning by reducing the input size to transformers, and is orthogonal to other acceleration methods such as Flash Attention and xformers.

Though ToMeSD has shown considerable theoretical speedups by significantly reducing the number of tokens, it struggles to accelerate diffusion models in practice with modern attention implementation and advanced GPU architectures. This is because the merge algorithm of ToMeSD Bolya & Hoffman (2023) requires relatively costly operations on GPU (e.g., sorting). This creates significant overhead that overshadows the speedups gained from token reduction, especially with more efficient implementations of attention (e.g., Flash Attention Dao (2023)) and GPU architectures better optimized for attention-like operations.

In this paper, we propose Token Merge with Attention (ToMA) to get practical speedups for diffusion models in a plug-and-play manner. Our method first utilizes a submodular function to identify a representative subset of tokens as merge destinations with a vectorized optimization algorithm that runs efficiently on GPUs. We then perform token merge by using an attention-like operation between the destination tokens and all tokens in the sequence, resulting in a linear transformation. For unmerge, we utilize the inverse or the transpose of such a linear transformation. The design of ToMA carefully considers the advantages and limitations of GPU computations.

To further reduce the overhead of ToMA, we leverage the locality characteristics of the hidden states within the latent space, which preserves image locality, so the tokens are more likely to be similar within a local region. By partitioning the hidden states into local regions, we can run ToMA in each region independently and mitigate the overhead costs by reducing the input size to ToMA while enjoying the parallelism of computation. Moreover, we also find that the destination selection and linear transformation of merge and unmerge can be shared across network layers and diffusion steps, which further decreases the ToMA overhead costs. As a result, ToMA achieves 30%-50% speedups without noticeable sacrifice in image quality.

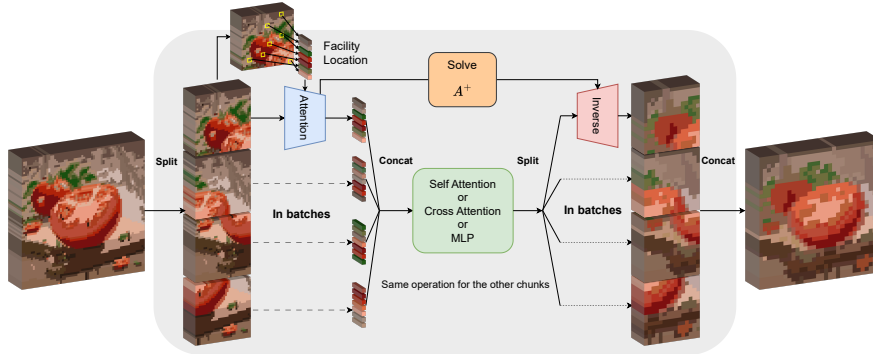
## 2 RELATED WORK

**Efficient Transformer:** The core of Transformers’ Vaswani et al. (2017) quadratic time complexity poses a bottleneck for both inference and training. Various methods have attempted to address this problem. To reduce computation complexity, ReformerKitaev et al. (2020) uses locality-sensitive hashing where LinformerWang et al. (2020), PerformerChoromanski et al. (2020), Low-Rank TransformerWinata et al. (2020) leverage low-rank approximations of the self-attention matrix to speed up computation. PerceiverJaegle et al. (2021), CharformerTay et al. (2021), and Funnel TransformersDai et al. (2020) use different ways to downsample the input to reduce the computation cost. Moreover, Sparse TransformerChild et al. (2019) and Big BirdZaheer et al. (2020) design different sparse attention patterns to let each token attend to a subset of all tokens. FRDiff So et al. (2024) accelerates diffusion inference by reusing feature maps across time steps but not merging the tokens.

**Learned Token Reduction** The majority of learned token reduction involve training auxiliary models to assess the importance of tokens in the input data. For example, DynamicViT Rao et al. (2021) employs a lightweight MLP module to generate pruning masks based on input token features. These masks are learned through a distillation process. GQA Ainslie et al. (2023) introduces an innovative mechanism that shares key and value heads across multiple query heads, balancing between the flexibility of multi-head attention and the efficiency of multi-query attention. A-ViT Yin et al. (2022) efficiently computes halting probabilities using the first channel of features, guided by auxiliary losses. Despite their effectiveness, these methods often require additional fine-tuning of auxiliary modules, which can be seen as a limitation. **Language-Vision Acceleration.** CrossGET Shi et al. (2024) combines token but on vision-language models with tasks like image captioning and image-text retrieval. TRIPS Ye et al. (2024) proposes text-relevant image patch selection but it accelerates the image-language model pertaining. DiffRate Chen et al. (2023) incorporates the compression rate and merges tokens in the vision transformers but in the training stage.

**Heuristic Token Reduction** Unlike learned token reduction techniques, some works have introduced heuristic token reduction strategies that can be directly applied to pre-trained ViTs without

108 requiring additional fine-tuning. For instance, Adaptive-Token Sampling Fayyaz et al. (2022) se-  
 109 lects tokens based on their similarity to the class token in the attention map, which outperforms the  
 110 top-k sampling. However, the requirement of the class token poses a limitation in dense prediction  
 111 tasks such as image generation. Token Pooling Marin et al. (2021) merges spatially adjacent tokens  
 112 within a local window to reduce the token count at various stages of the ViT. Token merge Bolya  
 113 et al. (2023) introduced a different pooling method that merges similar tokens based on an effi-  
 114 cient bipartite matching algorithm. ToMeSD Bolya & Hoffman (2023) randomly groups tokens into  
 115 source and destination groups and merges the source tokens with the destination tokens based on the  
 116 pair similarity score. This is one of our baselines.



117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130 Figure 2: Overview of **ToMA**. **Facility Location** selects representative destination tokens  $D$  from the token  
 131 set  $N$  using submodular optimization. **Attention** computes a low-rank projection matrix mapping  $N$  to  $D$ ,  
 132 followed by standard transformer operations (e.g., self-attention, cross-attention, or MLP) on  $D$ . **Inverse**  
 133 reconstructs  $N$  from  $D$  via a pseudo-inverse or transpose. These steps can be applied locally to latent space  
 134 regions as batch operations for efficiency.

### 135 3 PRELIMINARIES

137 **Attention Computation and Notation.** We denote the attention computation as SDPA (scaled dot  
 138 product attention). The input query is  $Q \in \mathbb{R}^{N \times d}$ , key is  $K \in \mathbb{R}^{N \times d}$ , and value is  $V \in \mathbb{R}^{N \times d}$ .  $N$   
 139 is sequence length and  $d$  is the hidden state dimension.

$$140 \text{SDPA}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V \quad (1)$$

144 Additionally, we denote  $D$  as the size of all destination tokens and  $B$  as the batch size,  $X$  as the  
 145 latent matrix of shape  $N \times d$  that gets projected into  $Q$ ,  $K$  and  $V$  (for simplicity, we assume the same  
 146 feature dimension in latents and attention). Matrices and vectors are in **bold**, while others are not.

147 **Submodularity.** A submodular function (Fujishige, 2005) is a set function  $f : 2^V \rightarrow \mathbb{R}$  with the dimi-  
 148 nishing return property:  $f(v|A) \geq f(v|B)$  if  $v \notin B$ ,  $A \subseteq B$ , where  $f(v|A) := f(v \cup A) - f(A)$ . In-  
 149 tuitively, the property states that the gain of a smaller  
 150 subset is always greater or equal to that of a larger sub-  
 151 set. This makes submodular function  $f(A)$  very useful  
 152 in expressing the diversity of the input subset  $A$  relative  
 153 to the ground set  $V$ .

```

148 for  $i = 1 \dots k$  do
149   |  $v^* \in \arg \max_{v' \in V \setminus A} f(v'|A)$ ;
150   |  $A = A \cup \{v^*\}$ ;
151 end
152 return  $A$ 

```

154 **Algorithm 1** Greedy

156 The submodular maximization problem with a cardinality constraint is shown below 2.

$$157 \max_{A \subseteq V} f(A) \quad \text{s.t. } |A| \leq k \quad (2)$$

158  
159  
160  
161 The greedy algorithm (Alg. 1) guarantees a  $(1 - 1/e)$ -approximation of the optimal solution. It iteratively selects the element that maximizes the gain until the chosen set size reaches  $k$ .

**Conditional Diffusion Model.** The conditional diffusion model is a variation from the diffusion model. Different from the unconditional diffusion model, the conditional one estimates the data distribution with the additional information. The forward noise process is defined in 3

$$q(\mathbf{x}_t|\mathbf{x}_0) := \mathcal{N}(\mathbf{x}_t|\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad (3)$$

The model gradually adds noise to the input image over steps  $t$  to transform input  $\mathbf{x}_0$  to a latent noise representation.  $\bar{\alpha}_t := \prod_{s=0}^t \alpha_s = \prod_{s=0}^t (1 - \beta_s)$  and  $\beta_s$  represents the noise variance schedule Ho et al. (2020). The denoising process below

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\sigma}_t^2 \mathbf{I}) \quad (4)$$

is parameterized by the neural network  $\boldsymbol{\mu}_\theta$ , where  $\boldsymbol{\sigma}_t^2$  denotes the transition variance. It aims to iteratively reconstruct  $\mathbf{x}_0$  from the random noise.

## 4 TOMA

Token merge selects several destination tokens from the full set. It merges the tokens into destinations based on similarity scores, typically by assigning the merged destination as the average of the merge tokens. If the merged tokens are very similar, e.g., regions of pixels that represent the background with homogeneous colors, then the loss of information in the merge process can be minimized. Compared to token pruning, which directly throws away tokens, the token merge retains more information, thus achieving better image quality.

Token merge reduces the number of inputs processed in the transformer block, leading to significant computational savings. Thus, we can achieve a theoretical speedup based on the token merge ratio and the computational complexity of the transformer block (details in Appendix). In the unmerge process, the values of the merged tokens are redistributed back to their original tokens to reverse the merge process. This operation ensures that the information from the merged tokens is restored while maintaining the shape of the output without token merge so that later layers can process without any modifications.

ToMA consists of three key stages, where we achieve significant improvements: 1) Destination Token Selection: Efficiently selecting the most representative tokens as destinations. 2) merge Tokens: merge source tokens into their corresponding destination tokens based on similarity computed using attention. 3) Unmerge Tokens: Restoring the merged tokens to their original forms by reversing the merge process as a linear operation. We also get further speedups by a) utilizing the locality characteristics of the latent space and b) sharing destination/merge/unmerge computations across iterations and layers to reduce overhead.

### 4.1 SUBMODULAR-BASED DESTINATION SELECTION

Let  $\mathbf{S}$  be the cosine similarity matrix between all token hidden representations.  $\mathbf{S}_i$  represents the  $i$ -th row of the  $\mathbf{S}$  matrix, and  $\mathbf{S}_{i,j} := \cos(\mathbf{X}_i, \mathbf{X}_j)$ . We denote the chosen destination token set as  $T$ , and all the tokens (the ground set in submodular optimization) as  $V$ .  $L \in \mathbb{R}^{|V|}$  is the max cache vector and  $L_i$  is the max similarity score between the token  $i$  in the ground set to our chosen set  $T$ .

The submodular function we use for destination selection is the facility location function (FL)  $f_{\text{FL}}$  shown in Eq. 5. FL sums the similarity between every token  $v_i \in V$  with its most similar neighbor in the selected destination set  $v_j \in T$ . Therefore, a high value for  $f_{\text{FL}}(T)$  means every token  $v_i$  has a similar neighbor in  $T$ , and  $T$  is representative of  $V$ , which perfectly matches the goal of merge destination selection. We also note that our framework is general, and  $f_{\text{FL}}$  could be potentially replaced with any other submodular function.

$$f_{\text{FL}}(T) = \sum_{v_i \in V} \max_{v_j \in T} \mathbf{S}_{i,j} \quad (5)$$

When optimizing  $T$  using the greedy algorithm, we essentially need to identify the next best token with the largest gain  $f_{\text{FL}}(v|T')$  relative to the so-far-selected set  $T'$ . The largest gain can be

decomposed (details in Appendix) as  $\arg \max_{i \notin T'} \sum_{j=1}^N \max(0, \mathbf{S}_{i,j} - \mathbf{c}_j(T'))$  with  $\mathbf{c}$  as a vector containing the cached max values of  $T'$  and updated incrementally as we select the next destination token:  $\mathbf{c}_j(T') = \max_{v_l \in T'} \mathbf{S}_{j,l}$ . We can perform all those operations in matrix forms, which makes them a perfect fit for GPU computation. There are more efficient submodular optimization algorithms compared to greedy, such as lazier-than-lazy greedy (Mirzasoleiman et al., 2015). However, the more complicated algorithms introduce operations (e.g., random subset selection) that are not ideal for GPU implementations.

#### 4.1.1 WHY SUBMODULAR

The submodular function, particularly the facility location function, offers a theoretical guarantee for optimizing the selection of elements with the highest information gain from a set. This characteristic aligns well with our requirements for token selection, ensuring that we choose the most similar tokens for merge with minimal information loss. Furthermore, facility location is highly compatible with GPU implementations, as it takes advantage of matrix operations, significantly boosting computational efficiency. Finally, facility location is compatible with an arbitrary similarity function, where we use cosine similarity that more closely aligns with the attention computation (supposing the input tokens are properly normalized using, e.g., Layernorm (Ba, 2016)).

We have also considered clustering-based methods such as k-means for token selection. However, we opted against them for several reasons. First, k-means provides a soft target, which might introduce artifacts and achieve suboptimal performance. Second, k-means assumes that clusters have ball-like boundaries, which is not flexible and imposes extra assumptions on the latent space. Third, the k-means method requires variable iterations to converge, and it is hard to control the computational costs and trade-off with the clustering quality.

## 4.2 MERGE AND UNMERGE WITH ATTENTION

In ToMA, we achieve token merge through a linear transformation approach that uses attention weights to assign tokens. This allows for a more generalized merge process. Accordingly, the unmerge operation can be the inverse of the linear transformation.

### 4.2.1 MERGE

We first compute softmax similarity weights between all destination tokens and all source tokens. This weight can be optionally sharpened or softened using the temperature parameter. Next, we normalize the matrix by counting the sum of each row and dividing the corresponding row by that value. Finally, we merge tokens together as a weighted average.

$$\mathbf{A} = \text{SDPA}(\mathbf{X}_T, \mathbf{X}, \mathbf{I}, \tau) = \text{softmax} \left( \frac{\mathbf{X}_T \mathbf{X}^\top}{\tau} \right) \mathbf{I} \tag{6}$$

$$\tilde{\mathbf{A}}_{ij} = \frac{\mathbf{A}_{ij}}{\sum_j \mathbf{A}_{ij}} \quad (\text{Normalize each row of } \mathbf{A}), \mathbf{X}_{\text{merged}} = \tilde{\mathbf{A}} \mathbf{X}$$

We describe the merge operation in Eq. 6. Here,  $\mathbf{X}_T$  is the hidden representation matrix for the set of destinations (shape  $D \times d$ ), and  $\tau$  is the temperature.  $\mathbf{X}_T$  are essentially sub-rows of  $\mathbf{X}$  so that the attention is between the destinations and all the tokens. The softmax is computed over all the destinations for every source token, where intuitively, we can think every source token gets distributed to some destinations, and the sum of the weights is 1. Because the source tokens include the destinations, in the worst case, every destination gets assigned by itself (e.g., if the destination is dissimilar to all other tokens). Note that we include the identity matrix to match the attention notation, which can be ignored in implementation.

For extremely small temperature values, the attention linear projection  $\mathbf{A}$  contains 1's and 0's, so our attention-based merge recovers the hard discrete merge by approximating the average of the merged tokens. Moreover, as we essentially compute an attention matrix and use it as a linear projection on the source tokens, our merge can be highly efficient on modern GPU architecture. Also, the linear transformation can be stored and reused later.

270 4.2.2 UNMERGE  
271

272 To unmerge tokens, we inverse the projection matrix of merge with the following two options (Eq. 7):

273 **Transpose** of the merge matrix  $A^\top$ : By multiplying the transpose of the merge matrix with the  
274 output of the transformer block, we distribute the merged token values back to their original tokens.  
275 When the temperature is extremely low, the transpose unmerge copies the computation result from  
276 the destination token to the corresponding merged tokens. This incurs very little overhead.

277 **Pseudo-inverse** of the merge matrix  $A^\dagger$ : Viewing the merge as a linear transformation, the pseudo-  
278 inverse minimizes the reconstruction error if the computation between merge and unmerge is close to  
279 linear. It is much more computationally expensive than  $A^\top$  and requires SVD or QR decomposition.

281  
282 
$$\mathbf{X}'_{\text{unmerged}} = \mathbf{A}^\top \mathbf{X}' \text{ or } \mathbf{X}'_{\text{unmerged}} = \mathbf{A}^\dagger \mathbf{X}' = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{X}' \quad (7)$$
  
283

284  $A^\top$  and  $A^\dagger$  are the same if the rows of the merge matrix  $A$  are independent, e.g., source tokens  
285 are not overlapping among different destinations. Intuitively, this means that the destinations should  
286 be as diverse as possible, which also matches the objective of the submodular optimization. Also,  
287 when the temperature is extremely low, every source token gets assigned to a single destination  
288 token, so the two options are identical. Concerning efficiency, we opt to use the transpose as the  
289 default unmerge method for ToMA.

290  
291 4.3 FURTHER SPEEDUP

292  
293 The overhead of ToMA consists of 1) the computation of destination tokens using submodular opti-  
294 mization, 2) the computation of the attention merge and unmerge matrix, and 3) applying the merge  
295 and unmerge matrices before and after layers in the transformer block. We further reduce all three  
296 overheads by considering the locality of the feature space so we can perform the computations lo-  
297 cally in every region. We also decrease the frequency to compute 1) and 2) by sharing destinations  
298 and merge matrices across iterations and layers.

299  
300 4.3.1 LOCAL REGION

301 A crucial aspect ignored in ToMe for SD is the  
302 locality of the latent space (Fig. 3). The spa-  
303 tial relationship between the generated image  
304 and the hidden states within the UNet model be-  
305 comes evident when examining this figure. We  
306 apply K-means to the tokens and recolor them  
307 based on their class affiliation. Specifically, by  
308 projecting the color from the generated image  
309 onto the hidden state feature map, we observe  
310 clear spatial coherence. The tokens in the la-  
311 tent space consistently demonstrate the greatest  
312 similarity with their neighboring tokens, creat-  
313 ing distinct localized regions.

314 This observation aligns with the intuition that  
315 images exhibit local consistency and smooth-  
316 ness. Therefore, we hypothesize that the most  
317 likely merge destination for any token resides  
318 within its local tile. This allows us to focus ex-  
319 clusively on the near tokens token while ignor-  
320 ing more distant ones.

321 To exploit this property, we limit operations to local regions, performing token selection and  
322 (un)merge within each region. We propose two region selections:

323 **Tile-shape region:** The tile region approach is particularly effective because it comprehensively  
captures the local characteristics by considering the image’s locality.

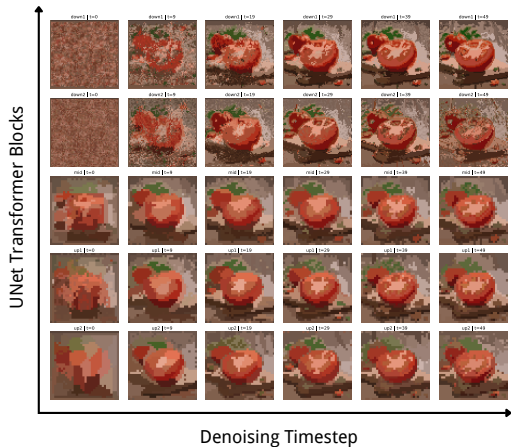


Figure 3: Recolored K-means results on UNet hidden states, across blocks/ denoising steps.

**Stripe-shape region:** The stripe region focuses on tokens on the same row, which misses the proximity in the vertical direction.

Both locality options can significantly reduce the computational overhead as we perform all operations in ToMA with a smaller number of input tokens in parallel. The tile-shape region is more coherent with the nature of the image and we find it performs better in experiments. However, turning a 2D matrix into tile-shape regions requires reshuffling, which brings additional overhead on GPUs. On the contrary, the stripe-shape option is faster as it only requires re-shaping of the 2D matrix while keeping its contiguous memory layout intact. We include both options in ToMA to provide trade-offs between speed and quality. Also, note that the tile-shape computation can be potentially accelerated as the low-level GPU operations are all in tiles, but it would require substantial re-implementation of the attention kernel. We defer this to future work.

We want to emphasize that the local region only affects components of ToMA, the computations in the transformer blocks always operate on the  $D \times d$  matrix of the merged destinations. Please refer to Appendix Alg. 3 for the detailed algorithm of ToMA with local regions.

### 4.3.2 SHARING OVERHEAD COMPUTATION

Our observations of the diffusion denoising process reveal that hidden states show substantial similarity across steps, meaning the selected destinations are also quite alike. As shown in Fig. 4, destination tokens have a significant intersection with the ones chosen in previous steps. Therefore, we can share destination selections across steps. Additionally, the linear (un)merge operations enable us to reuse the matrices across layers with minimal quality loss. By sharing both destinations and attention weights across steps and layers, we significantly reduce the number of times for destination computation and attend merge matrix computation while maintaining high output quality.

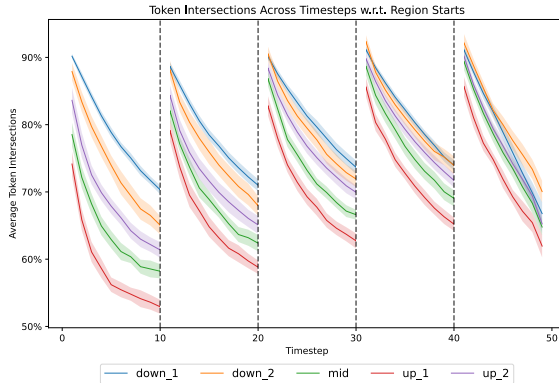


Figure 4: Intersection percentage of selected tokens between each step and the first step of its corresponding 10-step interval. Different curves refer to different layers in SDXL.

## 5 EXPERIMENTS

We evaluated ToMA on the SDXL stable-diffusion-xl-base-1.0 model using the Diffusers framework to generate  $1024 \times 1024$  images. The prompts were sourced from the GEMRec dataset (Guo et al., 2024) and ImageNet\_1K (Deng et al., 2009). To assess image quality, we used three primary metrics: CLIP, DINO, and FID (Radford et al., 2021; Caron et al., 2021; Heusel et al., 2017). For the CLIP and DINO evaluations, we generated images using 50 different prompts, each with 3 distinct seeds, and calculated the average score across all prompts and seeds. For measuring inference time, we reported the lowest wall-clock time over 100 runs.

**Diffusion Models.** We focus on the Stable Diffusion XL with the checkpoint of stable-diffusion-xl-base-1.0. SDXL is capable of generating very high-quality images and is popular in the community with abundant LoRAs available. Note that ToMA can generally apply to any transformer architecture. Thus, we can simply extend ToMA to other diffusion models like SD3 and SD2.

**Baseline.** ToMeSD (Bolya & Hoffman, 2023) selects tokens either in fixed or random small tiles, using a rigid approach. ToMeSD then discretely merges tokens by recording token pairs based on their computed similarity. In the unmerge phase, ToMeSD restores the original tokens by copying the values of the merged tokens back to their corresponding original tokens. We note that ToFu (Kim et al., 2023) is another relevant method that dynamically selects whether to prune or merge tokens based on the function’s linearity to accelerate diffusion models. The work done by ToFu is orthogo-

nal to our research, meaning our methods can be seamlessly integrated with ToFu to further enhance its performance. Therefore, we don't include ToFu in our comparison.

### 5.1 RESULTS ON QUALITY AND EFFICIENCY

The primary objective of this comprehensive experiment is to evaluate the trade-off between the quality of generated images and the computational efficiency across different token merge methods. Specifically, we compare **ToMA** with **ToMeSD**, using three key metrics: **CLIP**, **DINO**, and **FID**. These metrics measure visual similarity, attention mechanisms, and image fidelity, respectively. Additionally, we assess the generation time and speed-up ratio for each method, offering insights into computational gains.

We introduce three versions of ToMA to explore the impact of different attention mechanisms and facility location strategies:

**ToMA(Stripe Facility Location + Global Attention)**: Combines stripe-based facility location with global attention for token merge.

**ToMA\_stripe (Stripe Facility Location + Stripe Attention)**: Utilizes both stripe-based facility location and stripe attention for localized merge.

**ToMA\_tile (Tile Facility Location + Tile Attention)**: Applies tile-based facility location and attention for tile-wise merge.

For the Stripe Attention and Tile Attention versions, a tile size of 16 is used, while the Global Attention version uses 256 tiles. We compute facility location every 10 steps and attention weights every 5 steps over a total of 50 steps. The ToMeSD method serves as a baseline for comparison across all versions, using the same CLIP, DINO, and FID metrics.

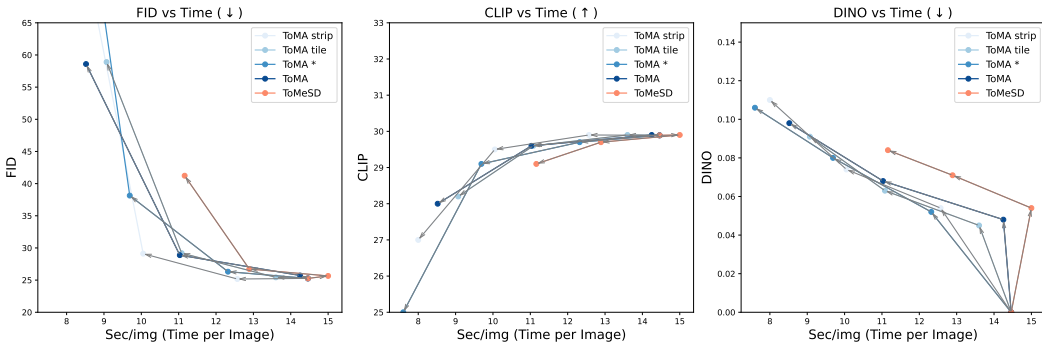


Figure 5: Quality metrics vs. generation time for SDXL-base on Nvidia V100. The merge ratio progresses from 0 to 0.75, moving from right to left following the directions of the arrows. Metrics are denoted as (↑: higher is better, ↓: lower is better).

In most cases, ToMeSD either increases computation costs or achieves minor speed-ups. Both versions of ToMA—stripe facility location with global attention and stripe facility location with stripe attention—maintain high image quality while delivering significant speed-ups. Tile facility location with tile attention achieves the best image quality, but the overhead is substantial, indicating great potential for further improvement through optimized low-level implementations.

In addition to evaluating image quality, we examine the generation time and speed-up ratio more variances. The token merge methods include ToMeSD, ToMA (with stripe and tile variations), and ToMA\*, which employs a "merge-once" strategy (merge and unmerge once for the entire transformer block instead of applying ToMA on every component of the transformer individually). Furthermore, we compare these methods to the LB (lower bound for speed-up), which is the best speedup we can get with a linear project merge and unmerge approach (apply random merge and unmerge projections without other overhead computations).



This comprehensive experiment, conducted under the same experimental setup (including dataset, tile size, and step scheduler), allows us to test the trade-off between **time** and **quality** across different token merge methods.

Token Merge Method	Generation Time / Speed Up Ratio		
	0.25	0.5	0.75
Baseline (ratio=0)	6.07s / 0.0%	6.07s / 0.0%	6.07s / 0.0%
ToMeSD	8.66s / +42.7%	8.73s / +43.8%	8.16s / +34.4%
ToMA	6.03s / -0.7%	5.04s / -17.0%	<b>4.34s / -28.5%</b>
ToMA_stripe	5.56s / -8.4%	<b>4.62s / -23.9%</b>	4.48s / -26.2%
ToMA_tile	6.20s / +2.1%	6.27s / +3.3%	6.23s / +2.6%
ToMA*	<b>5.45s / -10.2%</b>	4.91s / -19.1%	4.87s / -19.8%
LB	5.16s / -15.0%	4.01s / -33.9%	3.13s / -48.4%

Table 1: Comparison of Token Merge Methods and their Generation Time / Speed Up Ratios (RTX6000Ada). Negative percentages indicate faster times than the baseline, while positive percentages indicate slower times.

From Tab. 1, we find that ToMeSD exhibits increased generation times and negative speed-up ratios compared to the baseline across all token reduction ratios, ranging from +34.4% to +43.8%. This indicates that ToMeSD actually adds computational overhead rather than improving speed, making it less efficient than the baseline. We find that all ToMA variations, except ToMA\_tile, deliver significant speed improvements over ToMeSD ranging from -28.5% to -0.7%. Moreover, methods like ToMA\* and the ToMA stripe variants approach the theoretical speed limit, showcasing remarkable computational efficiency gains.

## 5.2 ABALATION TEST

In this section, we demonstrate on how we decide our default setting and other variances by comparing the image quality or speeds of combinations of different units in token merge.

### 5.2.1 FACILITY LOCATION & TILE FOR DESTINATION SELECTION

Destination Selection	CLIP	DINO	Time (s)	num_tiles	CLIP	DINO	MSE	sec/img
Global Facility	30.9486	0.0688	33.2	<b>4</b>	30.7747	0.0690	1564.0227	11.36
Tile Facility	<b>31.0185</b>	<b>0.0550</b>	<b>5.1</b>	<b>16</b>	30.9914	0.0566	1345.3114	6.44
Stripe Facility	30.9861	0.0740	5.16	<b>64</b>	31.0185	<b>0.0550</b>	<b>1274.1736</b>	5.04
Random	30.5527	0.0904	4.55	<b>256</b>	<b>31.0273</b>	0.0569	1296.3734	<b>5.01</b>

Table 2: Comparison of generated image metrics using different destination selection methods. Table 3: Tile facility comparison with 50 recompute steps, ratio=0.5, global attention

From Tab. 2, we find that the facility location demonstrates great performance in the generated image metric, which proves our theory that we should find the most representative tokens during selection. Also, the tile facility achieves the best **CLIP** and **DINO** score which aligns with our observation of the hidden states locality. By restricting the token merge process in a local region, we get a better image quality. Thus, we utilize the tile facility as our default setting for ToMA.

From Tab. 3, we examine the influence of different tile sizes. We find that the tile sizes of 64 achieve the best score in **DINO** and MSE while 256 shows great performance in **CLIP** and time. Generally, the metric difference is not significant between these two tile sizes. Thus, we select 256 as our default setting due to its lead in speed. We report ablation results on comparison between transpose and pseudo-inverse as well as different sharing schedules in the Appendix.

## 5.3 MERGE AND UNMERGE SPEED

In this section, we compare the speed ToMA (un)merge which generalizes this process as linear transformation, and the discrete (un)merge of ToMeSD. In this experiment, we utilize the transpose strategy as the unmerge of the merge matrix.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

$N = 4096$	0.25	0.5	0.75
ToMeSD Merge	0.2107	0.2048	0.2028
ToMA Merge	<b>0.0437</b>	<b>0.0403</b>	<b>0.0421</b>
ToMeSD Unmerge	0.1607	0.1811	0.1579
ToMA Unmerge	<b>0.0399</b>	<b>0.0483</b>	<b>0.0451</b>

Table 4: Comparison of ToMeSD and ToMA speeds for size 4096.

$N = 1024$	0.25	0.5	0.75
ToMeSD Merge	0.2022	0.2021	0.1932
ToMA Merge	<b>0.0390</b>	<b>0.0388</b>	<b>0.0388</b>
ToMeSD Unmerge	0.1605	0.1601	0.1440
ToMA Unmerge	<b>0.0402</b>	<b>0.0405</b>	<b>0.0396</b>

Table 5: Comparison of ToMeSD and ToMA speeds for size 1024.

Tab. 4 and Tab. 5 clearly demonstrate that the ToMA method significantly outperforms ToMeSD in both merge and unmerge speeds across all token reduction ratios (0.25, 0.5, and 0.75). For size 4096, ToMA achieves merge speeds approximately 80% faster than ToMeSD, with the lowest recorded merge time being 0.0403s compared to 0.2048s for ToMeSD. Similarly, the unmerge times for ToMA are consistently lower, with improvements ranging from 72% to 75% across the different token reduction ratios. This trend is mirrored in the 1024 size table, where ToMA again demonstrates its advantage, with merge and unmerge times consistently around 80% faster than ToMeSD. These results highlight the clear efficiency gains of the ToMA method in terms of both merge and unmerge processes, making it a more computationally efficient solution.



Figure 6: Image generated from left to right: original, ToMA\_stripe, ToMA\*, ToMA\_tile, ToMA

#### 5.4 DISCUSSION

Although the combination of tile facility location and tile merge produces high-quality images, it still falls short in speed. Optimizing this operation at a lower level could significantly reduce computational costs. Additionally, improving the implementation of the pseudo-inverse API would allow us to apply it to larger matrices, potentially enhancing image quality. Moreover, we utilize linear transformation, specifically SDPA, for token merge, where the parameter  $V$  is currently set as an identity matrix and ignored during inference. This  $V$  matrix holds potential for future training, which could further boost image quality.

**Broader impact.** ToMA enables speed improvements across a wide range of GPU architectures. On one hand, it accelerates image generation without compromising quality, making the process more efficient. On the other hand, it broadens the accessibility of diffusion models, allowing even those with less powerful or outdated GPUs to benefit from advanced techniques. ToMA reduces computational demands, making high-quality image generation feasible on a wider variety of hardware, thus making diffusion models more accessible to a larger audience.

## 6 CONCLUSION

In this work, we propose ToMA to enhance the existing token merge method in three key areas: 1) more representative token selection, 2) a more flexible and efficient merge and unmerge operation, and 3) the introduction of locality and sharing strategies. As a result, we achieve significant speedup while maintaining high image quality. For future work, we aim to further speed up the low-level tile region computation as well as fine-tune the merge attention for better generation quality.

## REFERENCES

- 540  
541  
542 Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit  
543 Sanghai. GQA: Training generalized Multi-Query transformer models from Multi-Head check-  
544 points. In *arXiv*, May 2023.
- 545 Jimmy Lei Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- 546  
547 Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion, 2023. URL <https://arxiv.org/abs/2303.17604>.
- 548  
549 Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy  
550 Hoffman. Token merging: Your ViT but faster. In *ICLR*, 2023.
- 551  
552 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Ar-  
553 mand Joulin. Emerging properties in self-supervised vision transformers. *CoRR*, abs/2104.14294,  
554 2021. URL <https://arxiv.org/abs/2104.14294>.
- 555  
556 Mengzhao Chen, Wenqi Shao, Peng Xu, Mingbao Lin, Kaipeng Zhang, Fei Chao, Rongrong Ji,  
557 Yu Qiao, and Ping Luo. Diffrate : Differentiable compression rate for efficient vision transform-  
558 ers, 2023. URL <https://arxiv.org/abs/2305.17997>.
- 559  
560 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
561 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 562  
563 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Jared Davis, Tamas  
564 Sarlos, David Belanger, Lucy Colwell, and Adrian Weller. Masked language modeling for pro-  
565 teins via linearly scalable long-context transformers. *Proceedings of ICLR*, 2020.
- 566  
567 Zihang Dai, Guokun Lai, Yiming Yang, and Quoc V Le. Funnel-transformer: Filtering out sequen-  
568 tial redundancy for efficient language processing. *Proceedings of NeurIPS*, 2020.
- 569  
570 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv*  
571 *preprint arXiv:2307.08691*, 2023.
- 572  
573 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-  
574 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,  
575 35:16344–16359, 2022.
- 576  
577 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-  
578 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
579 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 580  
581 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
582 *in neural information processing systems (NIPS/NeurIPS)*, 34:8780–8794, 2021.
- 583  
584 Mohsen Fayyaz, Soroush Abbasi Koohpayegani, Farnoush Rezaei Jafari, Sunando Sengupta, Hamid  
585 Reza Vaezi Joze, Eric Sommerlade, Hamed Pirsiavash, and Jürgen Gall. Adaptive token sampling  
586 for efficient vision transformers. In *ECCV*, pp. 396–414, 2022.
- 587  
588 Satoru Fujishige. *Submodular functions and optimization*. Elsevier, 2005.
- 589  
590 Yuanhe Guo, Haoming Liu, and Hongyi Wen. Gemrec: Towards generative model recommendation.  
591 In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*,  
592 volume 9 of *WSDM '24*, pp. 1054–1057. ACM, March 2024. doi: 10.1145/3616855.3635700.  
593 URL <http://dx.doi.org/10.1145/3616855.3635700>.
- 594  
595 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
596 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30,  
597 2017.
- 598  
599 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
600 *neural information processing systems (NIPS/NeurIPS)*, 33:6840–6851, 2020.

- 594 Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David  
595 Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A  
596 general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021.  
597
- 598 Minchul Kim, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. Token fusion: Bridg-  
599 ing the gap between token pruning and token merging, 2023. URL [https://arxiv.org/  
600 abs/2312.01026](https://arxiv.org/abs/2312.01026).
- 601 Sehoon Kim, Sheng Shen, David Thorsley, Amir Gholami, Woosuk Kwon, Joseph Hassoun, and  
602 Kurt Keutzer. Learned token pruning for transformers, 2022. URL [https://arxiv.org/  
603 abs/2107.00910](https://arxiv.org/abs/2107.00910).
- 604  
605 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In  
606 *International Conference on Learning Representations*, 2020. URL [https://openreview.  
607 net/forum?id=rkgNKkHtvB](https://openreview.net/forum?id=rkgNKkHtvB).
- 608 Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano, Sean  
609 Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, Daniel Haziza, Luca  
610 Wehrstedt, Jeremy Reizenstein, and Grigory Sizov. xformers: A modular and hackable trans-  
611 former modelling library. <https://github.com/facebookresearch/xformers>,  
612 2022.
- 613  
614 Dmitrii Marin, Jen-Hao Rick Chang, Anurag Ranjan, Anish Prabhu, Mohammad Rastegari, and  
615 Oncel Tuzel. Token pooling in vision transformers. In *arxiv*, October 2021.
- 616 Baharan Mirzasoleiman, Ashwinkumar Badanidiyuru, Amin Karbasi, Jan Vondrák, and Andreas  
617 Krause. Lazier than lazy greedy. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
618 volume 29, 2015.
- 619  
620 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
621 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
622 Sutskever. Learning transferable visual models from natural language supervision, 2021. URL  
623 <https://arxiv.org/abs/2103.00020>.
- 624 Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. DynamicViT:  
625 Efficient vision transformers with dynamic token sparsification. In *NeurIPS*, November 2021.  
626
- 627 Dachuan Shi, Chaofan Tao, Anyi Rao, Zhendong Yang, Chun Yuan, and Jiaqi Wang. Crossget:  
628 Cross-guided ensemble of tokens for accelerating vision-language transformers, 2024. URL  
629 <https://arxiv.org/abs/2305.17455>.
- 630  
631 Ethan Smith, Nayan Saxena, and Aninda Saha. Todo: Token downsampling for efficient generation  
632 of high-resolution images. *arXiv preprint arXiv:2402.13573*, 2024.
- 633  
634 Junhyuk So, Jungwon Lee, and Eunhyeok Park. Frdiff : Feature reuse for universal training-free  
635 acceleration of diffusion models, 2024. URL <https://arxiv.org/abs/2312.03517>.
- 636  
637 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
638 Poole. Score-based generative modeling through stochastic differential equations. In *International  
639 Conference on Learning Representations (ICLR)*, 2021. URL [https://openreview.  
640 net/forum?id=PXTIG12RRHS](https://openreview.net/forum?id=PXTIG12RRHS).
- 641  
642 Yi Tay, Vinh Q Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Si-  
643 mon Baumgartner, Cong Yu, and Donald Metzler. Charformer: Fast character transformers via  
644 gradient-based subword tokenization. *arXiv preprint arXiv:2106.12672*, 2021.
- 645  
646 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
647 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
648 tion processing systems (NIPS/NeurIPS)*, 30, 2017.
- 649  
650 Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with  
651 linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.

648 Genta Indra Winata, Samuel Cahyawijaya, Zhaojiang Lin, Zihan Liu, and Pascale Fung. Lightweight  
649 and efficient end-to-end speech recognition using low-rank transformer. In *ICASSP 2020-2020*  
650 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6144–  
651 6148. IEEE, 2020.

652 Wei Ye, Chaoya Jiang, Haiyang Xu, Chenhao Ye, Chenliang Li, Ming Yan, Shikun Zhang, Songhang  
653 Huang, and Fei Huang. Efficient vision-and-language pre-training with text-relevant image patch  
654 selection, 2024. URL <https://arxiv.org/abs/2403.07883>.

655 Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-ViT:  
656 Adaptive tokens for efficient vision transformer. In *CVPR*, pp. 10809–10818, June 2022.

657  
658 Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon,  
659 Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer  
660 sequences. *Proceedings of NeurIPS*, 2020.

661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## 702 A DETAILS ABOUT FACILITY LOCATION OPTIMIZATION

### 703 A.1 GAIN FUNCTION COMPUTATION IN FL

704 In the greedy algorithm we select the token  $v^*$  from  $V \setminus A$  that maximizes the following gain  
705 function:

$$706 v^* = \arg \max_{v \in (V-A)} f(v|A)$$

707 where  $f(v|A)$  is defined as:

$$708 f(v^*|A) = f(\{v^*\} \cup A) - f(A)$$

$$709 = \sum_{v \in V} \max_{v' \in (\{v^*\} \cup A)} \text{sim}(v, v') - \sum_{v \in V} \max_{v'' \in A} \text{sim}(v, v'')$$

710 Since for each  $v \in V$ , find the maximum corresponding  $v'$  in the updated representatative set  $\{v^*\} \cup$   
711  $A$  is equivalent to compare  $v'$  in  $A$  and  $v^*$ , namely:

$$712 \sum_{v \in V} \max_{v' \in (\{v^*\} \cup A)} \text{sim}(v, v') = \sum_{v \in V} \max \left[ \left( \max_{v' \in A} \text{sim}(v, v'), \text{sim}(v, v^*) \right) \right]$$

713 Therefore

$$714 f(v^*|A) = \sum_{v \in V} \max \left[ \left( \max_{v' \in A} \text{sim}(v, v'), \text{sim}(v, v^*) \right) \right] - \sum_{v \in V} \max_{v'' \in A} \text{sim}(v, v'')$$

$$715 = \sum_{v \in V} \max \left[ \left( \max_{v' \in A} \text{sim}(v, v'), \text{sim}(v, v^*) \right) - \max_{v'' \in A} \text{sim}(v, v'') \right]$$

$$716 = \sum_{v \in V} \max \left( 0, \text{sim}(v, v^*) - \max_{v'' \in A} \text{sim}(v, v'') \right)$$

717 Eventually,

$$718 v^* = \arg \max_{v' \in (V-A)} \sum_{v \in V} \max \left( 0, \text{sim}(v', v^*) - \max_{v'' \in A} \text{sim}(v, v'') \right)$$

## A.2 FACILITY LOCATION OPTIMIZATION ALGORITHM

**Algorithm 2:** Facility Location Token Selection Algorithm**Input:** Similarity matrix  $S \in \mathbb{R}^{N \times N}$ , number of tokens to select  $D$ **Output:** Selected token indices  $T$ **Initialize:**  $T \leftarrow \{\}$ ;**for**  $i = 1$  **to**  $N$  **do**    | Compute row sums:  $s_i = \sum_{j=1}^N S_{i,j}$ ;**end**Select the first token:  $t_1 \leftarrow \arg \max_i s_i$ ;Add  $t_1$  to  $T$ :  $T \leftarrow T \cup \{t_1\}$ ;Initialize the largest row:  $c \leftarrow S_{t_1}$ ;Set  $S_{t_1} \leftarrow 0$ ;**for**  $k = 2$  **to**  $D$  **do**    **for each token**  $i$  **not in**  $T$  **do**        | Compute gain:  $g_i = \sum_{j=1}^N \max(0, S_{i,j} - c_j)$ ;    **end**    Select next token:  $t_k \leftarrow \arg \max_{i \notin T} g_i$ ;    Add  $t_k$  to  $T$ :  $T \leftarrow T \cup \{t_k\}$ ;    Update largest row:  $c_j \leftarrow \max(c_j, S_{t_k,j})$  for all  $j = 1$  to  $N$ ;    Set  $S_{t_k} \leftarrow 0$ ;**end****return**  $T$

---

## B OVERALL DETAILED ALGORITHM OF TOMA

---

### Algorithm 3: ToMA with local regions

---

**Input:** Tensor  $\mathbf{X} \in \mathbb{R}^{B \times N \times d}$  (input tensor),  $D$  (number of destinations),  $\tau$  (attention temperature),  $F$  (computational layer)

- 1  $\mathbf{X} \leftarrow (\mathbf{X}_1, \dots, \mathbf{X}_P)$ ; /\* Reorganize  $\mathbf{X}$  as local regions \*/
- where  $\mathbf{X}_p \in \mathbb{R}^{B \times N_{local} \times d}$  for  $p = 1 \dots P$  and  $N_{local} \times P = N$ ;
- 2  $D_{local} \leftarrow D/P$ ;  $\mathbf{X} \leftarrow \mathbf{X}.reshape(B \times P, N_{local}, d)$ ;
- Step 1: Facility Location**
- 3 GPU Greedy to get:  $(T_1, T_2, \dots, T_{B \times P}) \leftarrow \text{Greedy}(f_{FL}, D_{local}, \mathbf{X})$ ;
- 4 Gather  $\mathbf{X}_T \leftarrow (\mathbf{X}_{1,T_1}, \mathbf{X}_{2,T_2}, \dots, \mathbf{X}_{B \times P, T_{B \times P}})$ ; /\* Shape:  $B \times P, D_{local}, d$  \*/
- Step 2: Merge**
- 5  $\mathbf{A} \leftarrow \text{SDPA}(\mathbf{X}_T, \mathbf{X}, \mathbf{I}, \tau)$ ; /\* Shape:  $B \times P, D_{local}, N_{local}$  \*/
- 6  $\tilde{\mathbf{A}} \leftarrow \mathbf{A} / \mathbf{A}.sum(-1)$ ; /\* Normalize each row \*/
- 7  $\mathbf{X}_{merged} \leftarrow \tilde{\mathbf{A}} \mathbf{X}$ ; /\* Apply Merge, Shape:  $B \times P, D_{local}, d$  \*/
- Computational Layer:**
- 8  $\mathbf{X}' \leftarrow F(\mathbf{X}_{merged}.reshape(B, D, d))$ ;
- Step 3: Unmerge**
- 9  $\mathbf{X}'_{unmerged} \leftarrow \tilde{\mathbf{A}}^T \mathbf{X}'$ ;
- 10 Group  $\mathbf{X}'_{unmerged}$  back to reverse the local region split;
- 11 **return**  $\mathbf{X}'_{unmerged}$

---



## C MORE ABLATION RESULTS

### C.1 UNMERGE COMPARISON(TRANPOSE VS PSEUDO-INVERSE)

Unmerge method	CLIP	DINO	MSE	Time (s)
<b>transpose</b>	<b>31.0273</b>	<b>0.0569</b>	1296.3734	<b>4.75</b>
<b>pseudo-inverse</b>	30.9972	0.0571	<b>1288.2609</b>	10.07

Table 6: Comparison of unmerge methods (transpose vs. pseudo-inverse) under the condition: 50 recompute steps, ratio=0.5, global attention (globalAttn)

From Tab. 6, we find the difference between transpose and pseudo-inverse method of unmerge shows similar outcome in scores while transpose is significantly faster than pseudo-inverse, we then set transpose as our default setting.

### C.2 SCHEDULE

In this section, we compare the metric of different sharing schedules of destination selections and attention weights computation

Dst steps	Attn steps	CLIP	DINO	MSE
first step	first step	30.0429	0.0773	2488.5682
every 10 steps	every 10 steps	30.8171	0.0729	1735.4429
every 10 steps	every 5 steps	30.8652	0.0699	1632.3441
every 10 steps	every 1 step	<b>30.9971</b>	<b>0.0668</b>	<b>1524.6436</b>
every 5 steps	every 5 steps	30.8923	0.0692	1608.6099
every 1 step	every 1 step	30.9196	0.0668	1551.5814

Table 7: Recompute schedule comparison with CLIP, DINO, and MSE metrics

From Tab. 7, we observe that recomputing attention and destination (Dst) steps more frequently generally results in slightly improved performance across the CLIP, DINO, and MSE metrics. Specifically, recomputing attention every step yields the best scores in all metrics, while less frequent recomputation (e.g., every 10 steps) results in slightly worse scores but still competitive results. The difference between recomputation frequencies becomes more noticeable in the MSE metric, where recomputing more frequently leads to lower error values. We select a recomputation schedule of computing attention every 5 steps and destinations every 10 steps because this provides nearly similar performance to the most frequent recomputation (every step) while likely being faster due to reduced computation overhead. This approach strikes a good balance between performance and efficiency.

## D THEORETICAL COMPLEXITY

We keep necessary constants for the complexity estimate as they are essential factors in practical speedup calculation. Also, we count the total number of multiplications by treating matrix multiplication as multiple dot products, ignoring algorithms with better theoretical complexity. The complexity is  $7d^2N + 2dN^2$  for a self-attention block. After using token merge, the complexity is:  $(7d^2D + 2dD^2)$  as we reduce the input size from  $N$  to  $D$ . We also define  $r := D/N$  as the reduction ratio. Thus, the speedup in terms of the reduction ratio is  $\text{Speedup} = \frac{7d+2N}{7d \cdot r + 2N \cdot r^2}$ .

The overhead of submodular optimization is:  $N^2d$

The overhead of computing merge attention projection is:  $NDd + Nd$

The overhead of merge is:  $NDd$

The overhead of unmerge with transpose is:  $NDd$

## E DETAILED RESULTS ON GEMREC & IMAGENET1K

Method	Ratio	FID	CLIP	DINO	MSE	RTX6000	V100	RTX8000
baseline_SDXL	0	25.265	29.889	0.000	0.000	6.1	14.5	16.1
ToMe	0.25	25.650	29.861	0.054	1716.131	8.7	15.0	16.9
	0.5	26.726	29.712	0.071	2279.389	8.7	12.9	14.6
	0.75	41.227	29.091	0.084	2344.868	8.2	11.2	12.4
ToMA strip	0.25	25.168	29.903	0.054	1604.185	5.6	12.6	14.5
	0.5	29.110	29.524	0.074	2199.760	4.6	10.1	12.0
	0.75	89.932	26.973	0.110	3185.344	4.5	8.0	9.5
ToMA tile	0.25	25.432	29.856	0.045	1348.644	6.2	13.6	15.7
	0.5	29.192	29.629	0.063	1912.216	6.3	11.1	13.2
	0.75	58.896	28.174	0.091	2802.324	6.2	9.1	10.7
ToMA *	0.25	26.311	29.696	0.052	1866.684	5.5	12.3	13.5
	0.5	38.138	29.061	0.080	3451.150	4.9	9.7	11.5
	0.75	123.366	24.963	0.106	5440.233	4.9	7.6	8.9
ToMA	0.25	25.718	29.858	0.048	1432.562	6.0	14.3	15.9
	0.5	28.875	29.640	0.068	2012.134	5.0	11.0	12.8
	0.75	58.592	27.961	0.098	2785.680	4.3	8.5	9.8
LTB	0.25	–	–	–	–	5.2	12.1	3.1
	0.5	–	–	–	–	4.0	9.9	7.8
	0.75	–	–	–	–	3.1	8.3	6.5

Table 8: Comparison of different methods with respect to FID, CLIP, DINO, MSE, and various GPU performance metrics (RTX6000Ada, V100, RTX8000).

972 F QUALITATIVE RESULT

973

974 F.1 MORE ON TOMA

975

976 Please refer to Fig. 7 for more qualitative result of ToMA.

977

978 F.2 COMPARISON WITH OTHER BASELINE MODELS

979 Please refer to Fig. 8 for more qualitative result of ToMA and other baseline models.

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

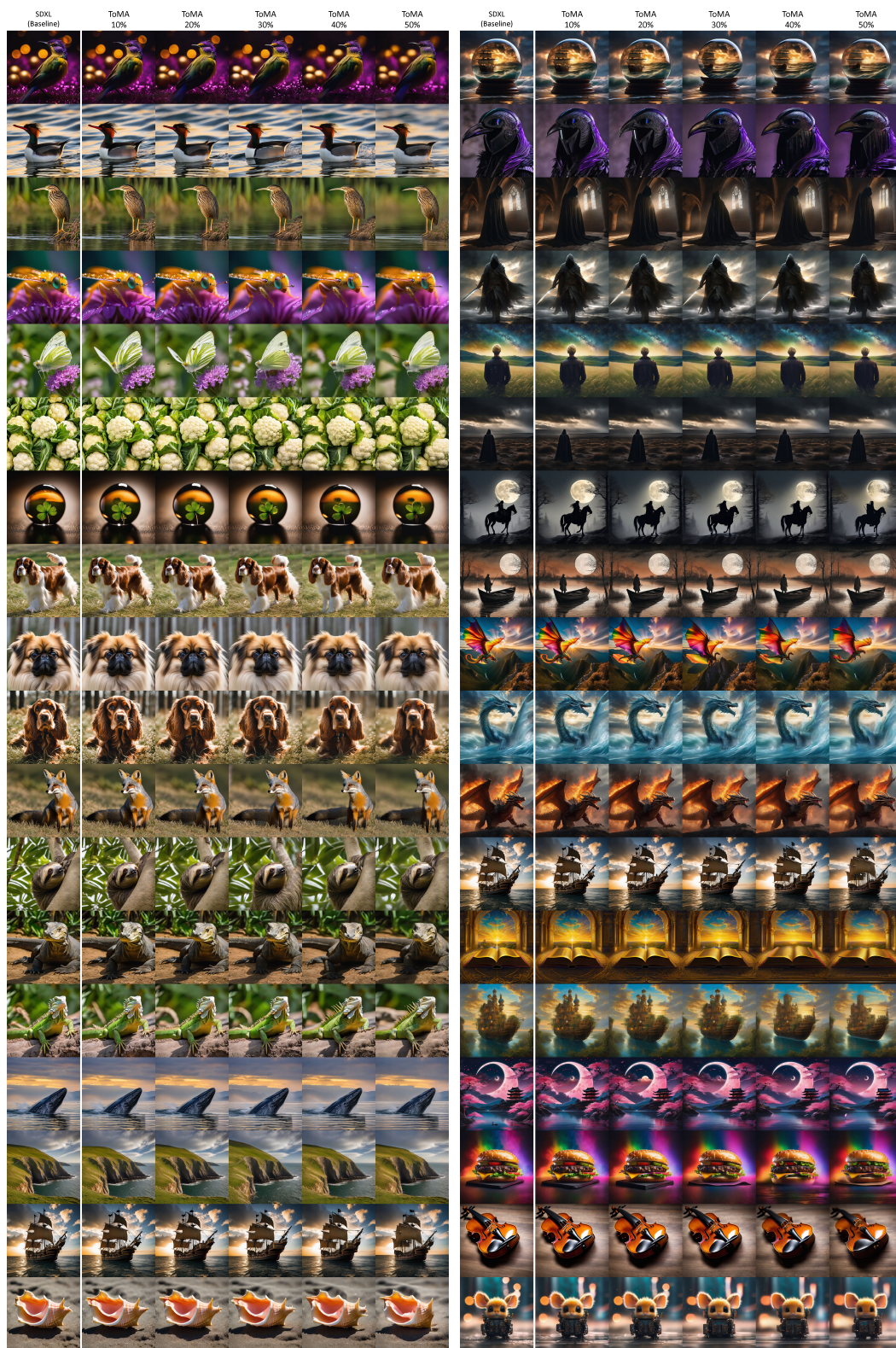


Figure 7: Visual examples of ToMA. Even with half of the tokens merged, ToMA consistently preserves image quality and often demonstrates greater robustness compared to other methods (ToDo, ToMeSD).

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

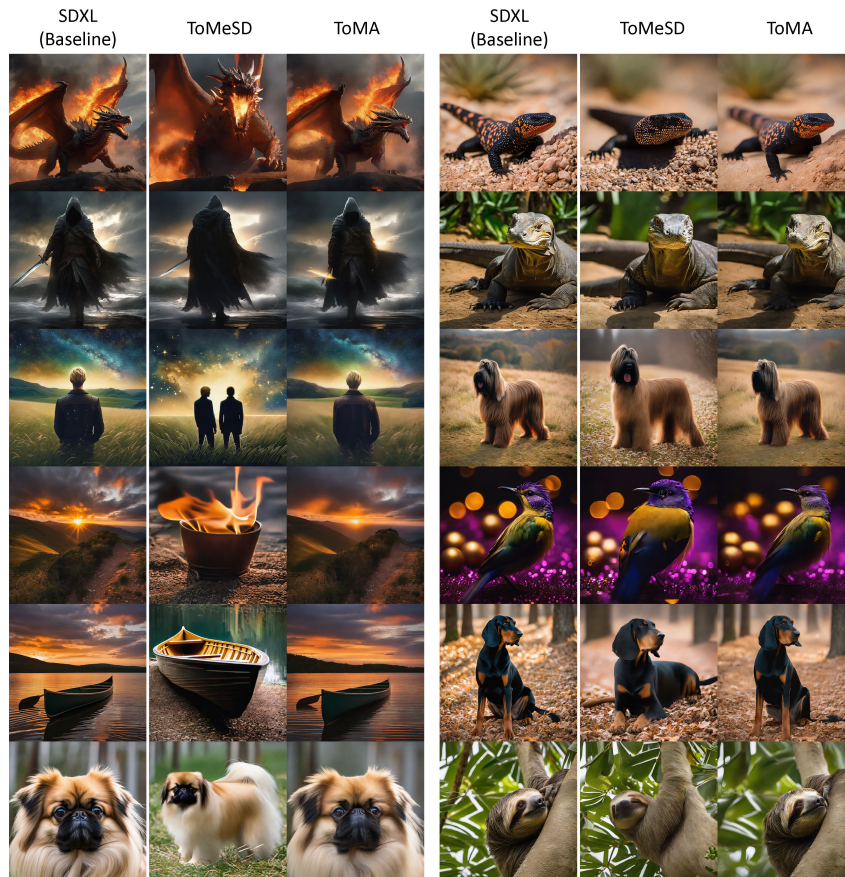


Figure 8: Qualitative comparison between Baseline SDXL-base, ToMeSD, and ToMA.

## G RESULT ON MORE BASELINE MODELS

We have compared ToMA with other baselines(eg. ToMeSD, ToFu, ToDoSmith et al. (2024)) and from the result we find that across all the ratios the ToMA achieve the most speedup and get better image quality compared to ToFu and ToDo.

Ratio	Method	FID↓	CLIP↑	DINO↓	MSE↓	Sec/img↓
<b>Baseline</b>	SDXL-base	25.27	29.889	0	0	6.07
0.25	ToMA	25.72	29.858	0.048	1,433	6.03
	ToMeSD	25.65	29.861	0.054	1,716	8.66
	ToFu	35.15	29.340	0.072	2,639	6.92
0.50	ToMA	28.88	29.640	0.068	2,012	5.04
	ToMeSD	26.73	29.712	0.071	2,279	8.73
	ToFu	142.08	25.039	0.134	7,408	6.83
0.75	ToMA	58.59	27.961	0.098	2,786	4.34
	ToMeSD	41.23	29.091	0.084	2,345	8.16
	ToFu	161.47	24.126	0.148	5,318	6.76
	ToDo	68.59	27.635	0.093	3,694	5.67

Table 9: Comparison of SDXL-base and various methods for generating 1024x1024 images for 50 denoising steps. ToDo is given a consistent ratio of 0.75 since it applies a 4x downsample for KV. Metrics are denoted as (↑: higher is better, ↓: lower is better), with the best performance highlighted.

## H DIFFUSION TRANSFORMERS (DiTs)

### H.1 DiT LOCALITY

We examined the hidden states of DiT models, focusing specifically on the `FLUX.1-dev` setting. Using visualization techniques, we analyzed the hidden states at the start of each block and across the denoising timesteps. As shown in Figure 9, the hidden states, despite the lack of convolutional layers, appear to closely represent the true image. Our analysis indicates that this locality is introduced apart from the VAE through the positional embeddings incorporated in DiT models, such as rotary embeddings in Flux and sin/cos embeddings in SD3 and SD3.5. Practically, through our experiments, we applied submodular-based token selection within local regions, which resulted in high-quality images.

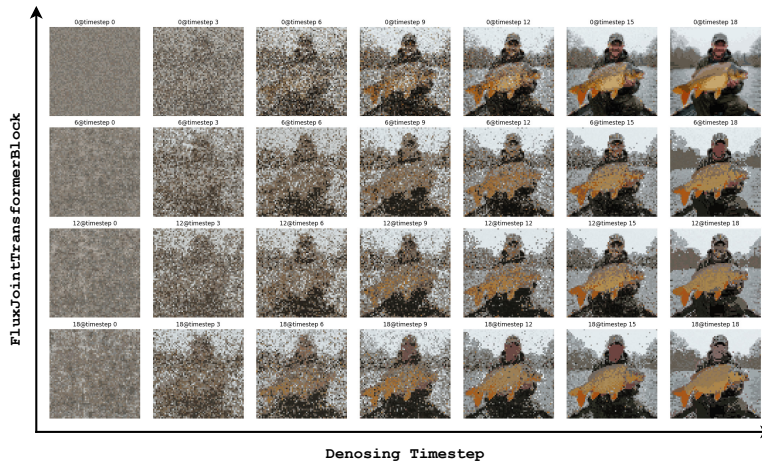


Figure 9: Recolored K-means results on hidden states of Flux.1-dev, across blocks & denoising steps.

### H.2 SPECIAL DESIGN OF TRANSFORMER BLOCKS AND POSITIONAL EMBEDDING

Due to the unique design of the transformer blocks in DiT models, which combine attention blocks and MLPs differently compared to the traditional setup of self-attention, cross-attention, and MLP, existing token merging methods such as ToMeSD, ToFu, and ToDo cannot be directly applied which would lead to all black or pure white noise. Additionally, the influence of positional embeddings further complicates their applicability since the naive application of token merging can lead to the selection of tokens that are not the most similar, which significantly degrades performance.

To address these issues, we implemented specific adaptations to the transformer blocks and positional embeddings, allowing our approach to successfully generate correct images with minimal quality loss which is shown in Tab. FIXME. Our method was selectively skip the first 10 transformer blocks in FLUX.1 to enable better the blend of text and image.

### H.3 RESULTS ON DiT

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

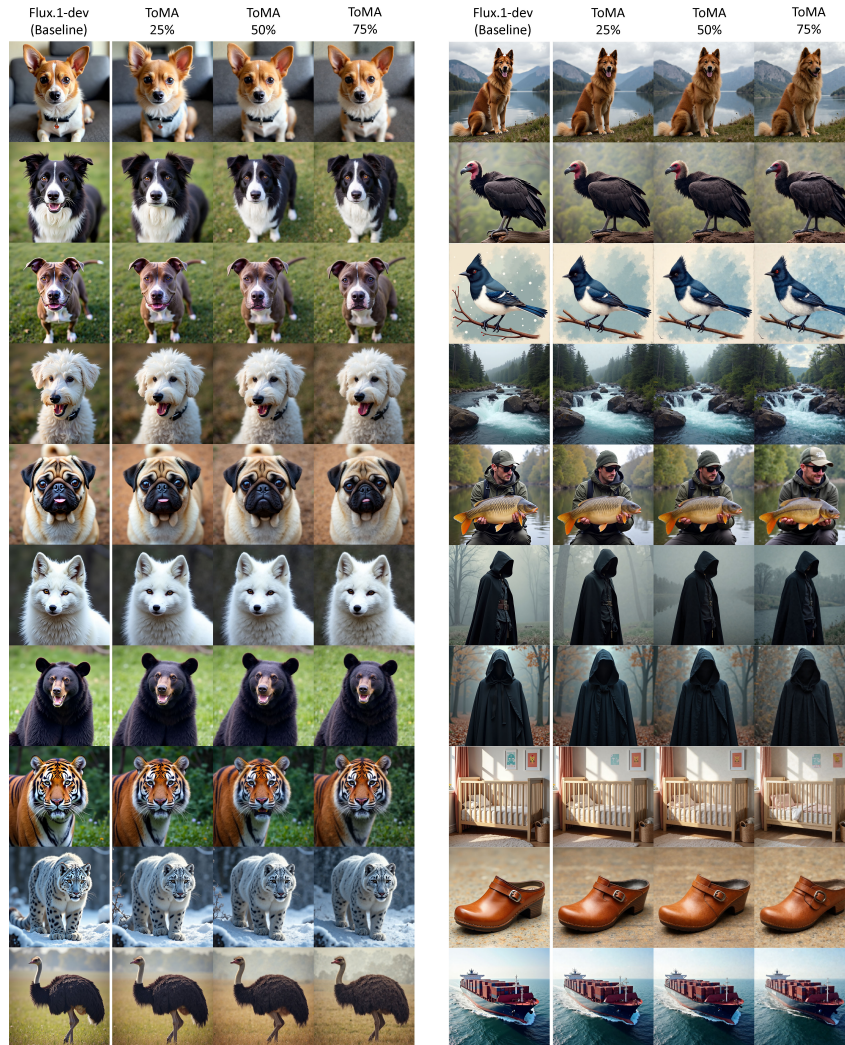


Figure 10: Qualitative comparison between Baseline Flux1.0-dev and ToMA.



1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

Method	Ratio	FID↓	CLIP↑	DINO↓	MSE↓	Sec/img↓
<b>Baseline</b>	Flux.1-dev	31.56	29.026	0	0	21.03
ToMA	0.25	30.80	29.068	0.043	1,340	20.14
	0.50	31.70	29.091	0.051	1,579	18.58
	0.75	33.39	28.976	0.064	2,041	16.12

Table 10: Performance of Flux.1-dev and various methods for generating 1024x1024 images for 35 denoising steps. Metrics are denoted as (↑: higher is better, ↓: lower is better). No other model works on DiT models.