# Region-Adaptive Sampling for Diffusion Transformers

Ziming Liu[1]*, Yifan Yang[2], Chengruidong Zhang[2], Yiqi Zhang[1], Lili Qiu[2], Yang You[1]†, Yuqing Yang[2]†

[1]National University of Singapore, [2]Microsoft Research

{liuziming, youy}@comp.nus.edu.sg {yifanyang, yuqyang}@microsoft.com

https://aka.ms/ras-dit

## Abstract

*Diffusion models (DMs) have become the leading choice for generative tasks across diverse domains. However, their reliance on multiple sequential forward passes significantly limits real-time performance. Previous acceleration methods have primarily focused on reducing the number of sampling steps or reusing intermediate results, failing to leverage variations across spatial regions within the image due to the constraints of convolutional U-Net structures. By harnessing the flexibility of Diffusion Transformers (DiTs) in handling variable number of tokens, we introduce RAS, a novel, training-free sampling strategy that dynamically assigns different sampling ratios to regions within an image based on the focus of the DiT model. Our key observation is that during each sampling step, the model concentrates on semantically meaningful regions, and these areas of focus exhibit strong continuity across consecutive steps. Leveraging this insight, RAS updates only the regions currently in focus, while other regions are updated using cached noise from the previous step. The model's focus is determined based on the output from the preceding step, capitalizing on the temporal consistency we observed. We evaluate RAS on Stable Diffusion 3 and Lumina-Next-T2I, achieving speedups up to 2.36x and 2.51x, respectively, with minimal degradation in generation quality. Additionally, a user study reveals that RAS delivers comparable qualities under human evaluation while achieving a 1.6x speedup. Our approach makes a significant step towards more efficient diffusion transformers, enhancing their potential for real-time applications. Our code is available at https://github.com/microsoft/RAS.*

## 1. Introduction

Diffusion models (DMs) [8, 18, 41, 42] have proven to be highly effective probabilistic generative models, pro-

*This work is done while Ziming Liu's internship at Microsoft Research.
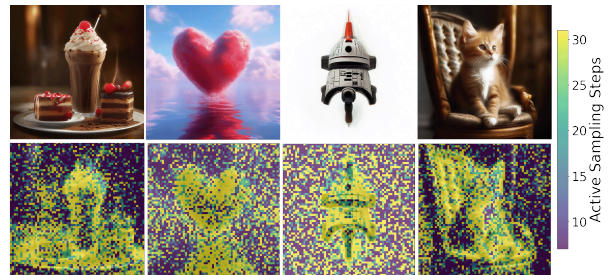
†Corresponding authors.



Figure 1. The main subject and the regions with more details are *brushed* for more steps than other regions in *RAS*. Each block represents a patchified latent token.

ducing high-quality data across various domains. Applications of DMs include image synthesis [7, 35], image super-resolution [13, 24, 51], image-to-image translation [23, 38, 47], image editing [22, 52], inpainting [29], video synthesis [3, 9], text-to-3D generation [33], and even planning tasks [20]. However, generating samples with DMs involves solving a generative Stochastic or Ordinary Differential Equation (SDE/ODE) [15, 34] in reverse time, which requires multiple sequential forward passes through a large neural network. This sequential processing limits their real-time applicability.

Considerable work has been dedicated to accelerating the sampling process in DMs by reducing the number of sampling steps. Approaches include training-based methods such as progressive distillation [39], consistency models [43], and rectified flow [1, 26, 27], and training-free methods such as DPM-solver [28], AYS [37], DeepCache [49], and Delta-DiT [5]. These methods, however, uniformly process all regions of an image during sampling, irrespective of the specific needs of different regions. Each sampling step in these methods treats every area of the image equally, predicting the noise for each region at the current time step before proceeding to the next. Intuitively, however, the complexity of different regions within an image varies: intricate foreground elements may require more sampling steps for clarity, while repetitive backgrounds could benefit from

(a) Lumina-Next-T2I

(b) Stable Diffusion 3

(c) Lumina-Next-T2I FID RAS VS Rectified Flow

(d) Lumina-Next-T2I CLIP Score RAS VS Rectified Flow

(e) Default VS RAS (1.625x throughput for Stable Diffusion 3 and 1.561x for Lumina-Next-T2I) Human Evaluation
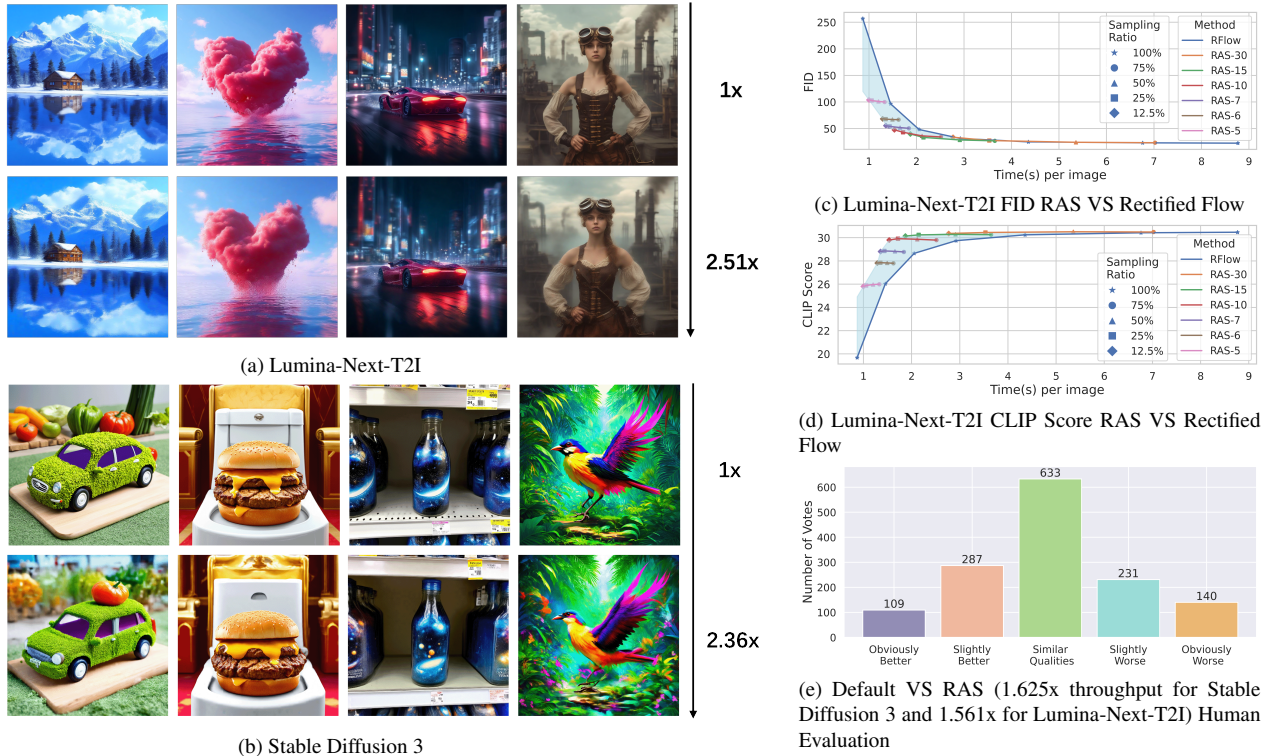
Figure 2. (a)(b) Accelerating Lumina-Next-T2I and Stable Diffusion 3, with 30 and 28 steps separately. (c)(d) Multiple configurations of RAS outperform rectified flow in both image qualities and text-following. RAS-X stands for RAS with X sampling steps in total. (e) RAS achieves comparable human-evaluation results with the default model configuration while achieving around 1.6x speedup.

more aggressive compression of sampling steps without significant loss of quality. This suggests a potential for a more flexible sampling approach that can dynamically adjust the sampling ratio across different regions, enabling faster, yet high-quality diffusion process.

This concept is a natural progression in the evolution of DMs. From DDPM [18] to Stable Diffusion XL [32], diffusion models have predominantly relied on U-Nets, whose convolutional structures [36] necessitate uniform treatment of all image regions due to fixed square inputs. However, with the advent of DiTs [31] and the increasing exploration of fully transformer-based architectures [45], the research focus has shifted towards architectures that can accommodate flexible token inputs. DiTs allow any number of tokens to be processed flexibly, opening up new possibilities. This shift has inspired us to design a new sampling approach capable of assigning different sampling steps to different regions within an image.

To further explore the feasibility of this idea, we visualized several diffusion process outputs at varying sampling steps (Figure 3). We observed two key phenomena: (1) the regions of focus in adjacent steps exhibit considerable continuity during the later stages of diffusion, and (2) in each sampling step, the model tends to focus on specific

semantically meaningful areas within the image. This pattern is akin to the process of an artist completing a painting in 1000 steps on a blank canvas, where each step involves a different brush stroke that selectively refines certain areas. This observation suggests that in each adjacent step, the "brushes" used by the diffusion model are similar, and hence, the areas they refine remain consistent. Thus, areas that the model momentarily ignores could potentially be excluded from the computationally intensive DiT processing, allowing the model to focus more on regions of immediate interest.

To validate this hypothesis, we conducted an experiment where we identified the ranking of tokens at each diffusion step using the metric of output noise we introduced, representing regions of primary focus for the model. By calculating the similarity of token rankings using NDCG (Figure 4), we found a high degree of continuity in the areas of focus between adjacent steps. This continuity motivated us to design a sampling approach that assigns different sampling ratios to different regions based on their attention continuity.

As is shown in Figure 5, our method leverages the output noise from the previous step to identify the model's primary focus for the current step (fast-update regions), al-
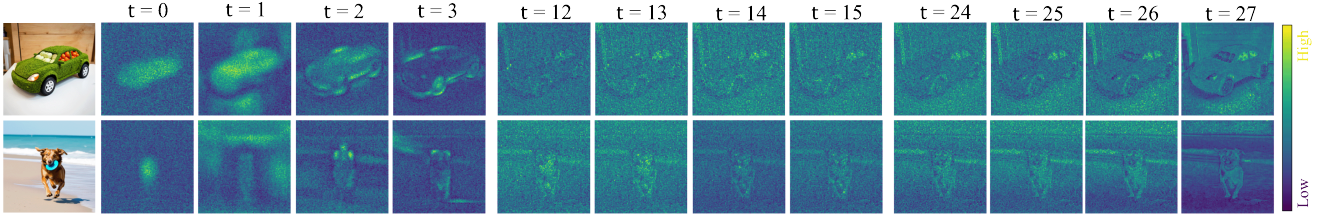
Figure 3. Visualization of predicted noise of each step. DiT model focuses on certain regions during each step and the change in focus is continuous across steps.

lowing only these regions to proceed through DiT for denoising. Conversely, for regions of less interest (slow-update regions), we reuse the previous step's noise output directly. This approach enables regional variability in sampling steps: areas of interest are updated with higher ratio, while others retain the previous noise output, thus reducing computation. After each local update, the updated regions' predicted noise values change, forming the noise map for sampling in the current step, which also serves as the selection criterion for the next fast-update regions. This operation restricts the tokens processed by DiT to those in the model's current area of focus, enhancing image refinement efficiency.
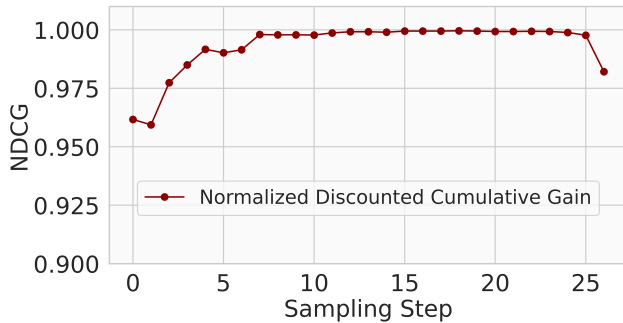


Figure 4. NDCG [21, 48] for each pair of adjacent sampling steps is high throughout the diffusion process, marking the similarities in the ranking of focused tokens ranging from 0 to 1.

For each input $X_t$, we select a fast-update rate to determine the regions needing updates in each step, while regions in the slow-update regime retain the previous noise output, which, combined with the updated fast-region noise, forms $X_{t-1}$ for the next step. To maintain global consistency, we keep features from slow-update regions as reference keys and values for subsequent steps. Although the fast-region selection is dynamic and recalculated after each update to prioritize significant areas, we periodically reset the inference for all regions to mitigate cumulative errors.

In summary, we propose *RAS*, the first diffusion sampling strategy that allows for regional variability in sampling ratios. Compared to spatially uniform samplers, this

flexibility enables our approach to allocate DiT's processing power to the model's current areas of interest, significantly improving generation quality within the same inference budget. As shown in Figure 2 (c)(d), our method achieves substantial reductions in inference cost with minimal FID increase, while outperforming the uniform sample baseline in terms of FVD within equivalent inference times. Figure 2 (a)(b) also demonstrates that with models like Lumina-Next-T2I [2] and Stable Diffusion 3 [11], our method's fast-region noise updating yields over twice the acceleration with negligible image quality loss. A user study comparing our method to uniform sampling across various generated cases further shows that our method maintains comparable generation quality at 1.6x the acceleration rate.
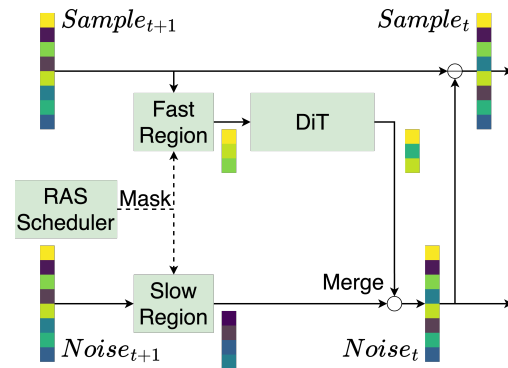


Figure 5. Overview of *RAS* design. Only current fast-update regions of each step are passed to the model.

## 2. Related Work

### 2.1. Diffusion Models: From U-Net to Transformer

Diffusion models [8, 18, 41, 42] have demonstrated significant capabilities across a range of generative tasks, often outperforming previous methods such as generative adversarial networks (GANs) [14] in many downstream applications. Historically, denoising diffusion probabilistic models (DDPMs) [18] and more recent models like Stable Diffusion XL [32] have predominantly utilized convolutional

U-Nets [36] as the backbone. Due to the structure of convolutional networks, it is necessary to maintain the original spatial resolution of the input samples to support operations like pooling, which inherently limits the ability to exploit spatial redundancy during the diffusion process, particularly when attempting to prune the latent samples used as model inputs.

Fortunately, the dilemma has been broken by the emergence of Diffusion Transformer (DiT) [31], which has also been used as the backbone of the SOTA diffusion models like Stable Diffusion 3 [10], Lumina T2X [2], Pixart-Sigma [4]. The biggest characteristic of DiT is that it completely eliminates the need for Convolutaional U-Net. DiT uses a pure Transformer [45] architecture and adds conditional information such as prompts with adaptive layer norm. In this way, the positional information is no longer provided by the convolutional operations, and the latent tokens are now positionally independent after position embedding. This enables us to utilize the redundancy we discovered in Section 1 and select the tokens that are likely to be focused in the current sampling step to compute while caching the predicted of other tokens from the previous step.

### 2.2. Efficient Diffusion Model Inference

To address the problem of high inference cost in diffusion models, various acceleration techniques have been proposed from different perspectives. A commonly used approach is to reduce the number of sampling steps. Some of these techniques require additional training, such as progressive distillation [39], consistency models [43], and rectified flow [1, 26, 27]. Among these methods, rectified flow has been widely used in models like Stable Diffusion 3 [10]. It learns the ODE to follow straight paths between the standard normal distribution and the distribution of the training dataset. These straight paths significantly reduce the distance between the two distributions, which in turn lowers the number of sampling steps needed.

There are also training-free methods to either reduce the number of sampling steps or decrease the computational burden within each step. DPM-solver [28], for instance, introduces a formulation that enhances the solution process of diffusion ODEs. DeepCache [49], specifically designed for U-Net-based models, leverages the U-Net architecture to cache and retrieve features across adjacent stages, allowing for the skipping of certain downsampling and upsampling operations during the diffusion process. However, these approaches treat all image regions uniformly, ignoring the varying complexity across different parts of the image. This uniform treatment can lead to significant computational inefficiency, as not all regions require the same level of processing.

As introduced in Section 1, the complexity of different regions within an image can vary substantially. To exploit the characteristics of both the diffusion process and the structure of Diffusion Transformers (DiTs), we introduce *RAS*, a novel approach designed to optimize computation by focusing on the distinctive properties of different image regions. *RAS* is also orthogonal to the methods we mentioned above, such as DiTFastAttn [50] and $\Delta$-DiT [5].

## 3. Methodology

| | |
|---|---|
| $t$ | The current timestep |
| $N$ | The noise output of the DiT model |
| $\widetilde{N}$ | The cached noise output from the previous timestep |
| $\hat{N}$ | The estimated full-length noise calculated with $N$ and $\widetilde{N}$ |
| $S$ | The unpatchified image sample |
| $x$ | The patchified input of the DiT model |
| $M$ | Mask generated to drop certain tokens in the input |
| $D$ | The number of times the tokens in a patch being dropped |

Table 1. Meanings of the symbols that are used in this paper

### 3.1. Overview

In this section, we present the *RAS* design and techniques to exploit inter-timestep token correlations and the regional token attention mechanism introduced in Section 1. (1) Based on the regional characteristics we observed in the DiT inference process, we propose an end-to-end pipeline that dynamically eliminates the computation through DiT of certain tokens at each timestep. (2) To leverage the continuity across consecutive timesteps, we propose a straightforward method to identify the fast-update regions that require refinement in upcoming timesteps, while ensuring that slow-update regions are not neglected due to insufficient diffusion steps. This approach effectively balances token focus without leading to starvation. (3) Building on our observations of continuous distribution patterns, we introduce several scheduling optimization techniques to further enhance the quality of generated content.

```python
def step(self, dit_output_fast, sample, t):
    dit_output = self.merge(dit_output_fast, self.cached_output_slow)

    # original sampling step
    diff = (self.sigma[t - 1] - self.sigma[t]) * dit_output
    prev_sample = sample - diff

    # identify and cache the slow-update region
    slow_indices, fast_indices, self.cached_output_slow = \
        self.region_identify(dit_output)

    # returning only the fast-update region to the model
    return prev_sample[fast_indices]
```

Figure 6. Sample Step with *RAS* in Python. Only two extra functions are needed to switch from the original scheduler to *RAS*.

## 3.2. Region-Adaptive Sampling

**Region-Aware DiT Inference with *RAS*.** Building on the insight that only certain regions are important at each timestep, we introduce the *RAS* pipeline for DiT inference. In U-Net-based models such as SDXL [32], tokens must remain in fixed positions to preserve positional information. However, given the structure of DiT, we can now mask and reorder elements within latent samples, as positional information is already embedded using techniques like RoPE [44]. This flexibility allows us to selectively determine which regions are processed by the model. To achieve this, some additional operations are required starting from the final step, as described in Figure 6. At the end of each timestep, the current sample is updated by combining the fresh model output for the active tokens and the cached noise for the inactive tokens. Specifically, the noise for the entire sequence is restored by integrating both the model output and the cached noise from the previous step. This mechanism enables active, important tokens to move in the new direction determined at the current timestep, while the inactive tokens retain the trajectory from the previous timestep. The next step involves updating the unpatchified sample with the scaled noise. We then compute the metric $R$, which is used to identify the fast-update regions based on the noise, update the drop count $D$ to track the frequency with which each token has been excluded, and generate the mask $M$ accordingly.

With the mask $M$, the noise for the slow-update regions is cached, while the sample for the current fast-update regions is patchified and passed through the DiT model. Since modules like Layernorm and MLP do not involve cross-token operations, the computation remains unaffected even when the sequence is incomplete. For the attention [45] module, the computation can still proceed with the query, key, and value tensors being pruned. Additionally, we introduce a caching mechanism to further enhance performance, which will be detailed later. In summary, *RAS* dynamically detects regions of focus and reduces the overall computational load of DiT by at least the same proportion as the user-defined sampling ratio.

**Region Identification.** The DiT model processes the current timestep embedding, latent sample, and prompt embedding to predict the noise that guides the current sample closer to the original image at each timestep. To quantify the refinement of tokens at each timestep, we use the model's output as a metric. Through observation, we found that the standard deviation of the noise strongly marks the regions in the images, with the main subject (fast-update regions) showing an obvious lower standard deviation than the background (slow-update region). This could be caused by the difference in the amount of information between the regions after mixing with the Gaussian noises. Utilizing the deviation as a metric achieves reasonable results of im-

age qualities and notable differences between regions, as is shown in Figure 8. Also, considering the similarities between latent samples across adjacent timesteps, we hypothesize that tokens deemed important in the current timestep are likely to remain important in the next, while the less-focused tokens can be dropped with minimal impact. Before we reach the final formulation of the metric, we need to introduce another technique to prevent starvation.

**Starvation Prevention.** During the diffusion process, the main subject regions typically require more refinement compared to the background. However, consistently dropping computations for background tokens can lead to excessive blurring or noise in the final generated image. To address this, we track how often a token is dropped and incorporate this count as a scaling factor in our metric for selecting tokens to cache or drop, ensuring less important tokens are still adequately processed.

Additionally, since DiT patchifies the latent tokens before feeding them into the model, we compute our metric at the patch level by averaging the scores of the tokens within each patch. Combining all the factors mentioned above, our metric can be written as:

$$R_t = mean_{patch}(std(\hat{N}_t)) \cdot exp(k * D_{patch}) \quad (1)$$

where $\hat{N}_t$ is the current estimated noise, $D_{patch}$ is the count of how many times the tokens in a patch have been dropped, and k is a scale factor to control the difference of sample ratios between fast-update regions and slow-update regions.
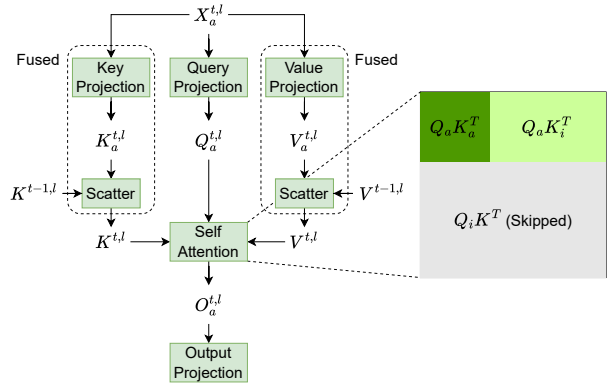


Figure 7. A *RAS* self-attention module using Attention Recovery to enhance generation quality. $X_a^{t,l}$, $Q_a^{t,l}$, $K_a^{t,l}$, $V_a^{t,l}$ and $O_a^{t,l}$ represent the input hidden states, query, key, value and attention output of active tokens on layer $l$ during step $t$, respectively. $K^{t,l}$ and $V^{t,l}$ denote the key and value caches. The scatter operation to partially upload the key and value caches are fused into the previous projection using a PIT GeMM kernel. The keys and values of the not-focused area ($K_i^{t,l}$ and $V_i^{t,l}$) are estimated with the cache from the last sampling step ($K^{t-1,l}$ and $V^{t-1,l}$).
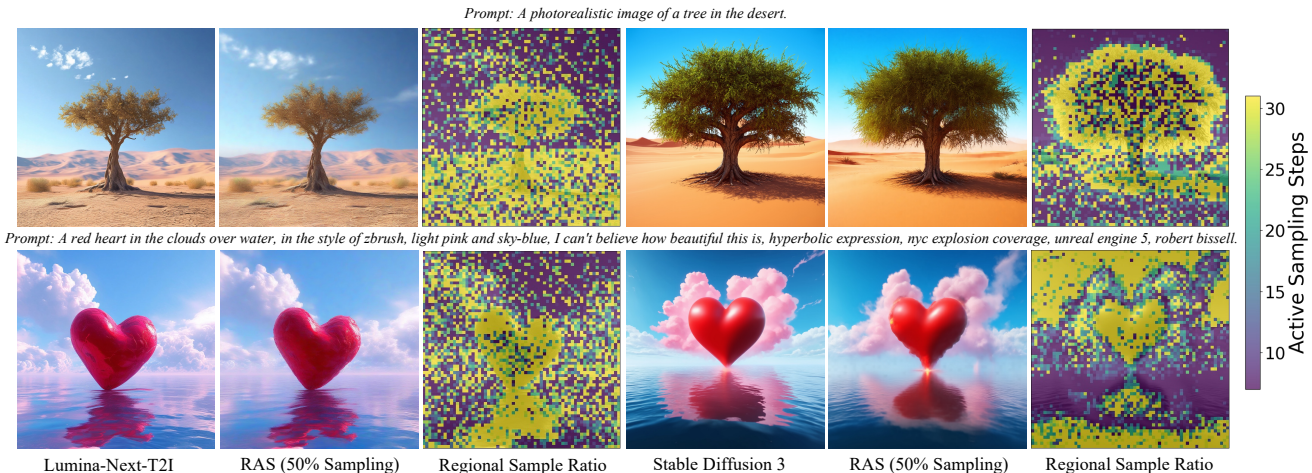
**Key and Value Caching.** As we know, the attention mech-

Figure 8. Visualization of *RAS* on Lumina-Next-T2I and Stable Diffusion 3.

anism works by using the query for each token to compute its attention score with each other tokens by querying the keys and values of the whole sequence, thus giving the relations between each two tokens. The attention of the active tokens in *RAS* can be given as:

$$O_a = softmax(\frac{Q_a K_a^T}{\sqrt{d}})V_a \qquad (2)$$

where $a$ stands for currently active tokens. However, the metric **R** we introduce to identify the current fast and slow regions does not take their contribution to the attention score into consideration. Thereby, losing these tokens during attention can cause a huge change in the final output. Our solution here is also caching. During each step, the full keys and values are cached until they are partially updated with the current active tokens. As is described in Figure 7, this solution is also based on the similarity between each two sampling steps, and now we can estimate the original attention output by:

$$O_a = softmax(\frac{Q_a [K_a, \widetilde{K}_i]^T}{\sqrt{d}})[V_a, \widetilde{V}_i] \qquad (3)$$

where $i$ stands for the inactive tokens.

### 3.3. Scheduling Optimization

**Dynamic Sampling Ratio.** As illustrated in Figure 4, the correlation between the initial timesteps is lower compared to the latter stages, where the diffusion process exhibits greater stability. This trend is also evident in Figure 3. Consequently, the strategy previously introduced is not suitable for the early stages of the diffusion process, as it could negatively impact the foundational structure of the generated image. Furthermore, we have observed that the similarity gradually increases during the stable phase of diffusion. To

address these observations, we propose a dynamic sampling ratios that maintains a 100% ratio for the initial timesteps (e.g., the first 4 out of 28 steps) to mitigate any adverse effects on the outline of the generated image. Thereafter, the sampling ratio is progressively reduced during the stable phase. This approach ensures a balance between computational efficiency and the image quality, enabling effective sampling ratios while minimizing adverse impacts on the generated output.

**Accumulated Error Resetting.** *RAS* focuses on the model's regions of interest, which tend to be similar across adjacent sampling steps. However, regions that are not prioritized for multiple steps may accumulate stale denoising directions, resulting in significant error between the original latent sample and the one generated with *RAS*. To mitigate this issue, we introduce dense steps into the *RAS* diffusion process to periodically reset accumulated errors. For instance, in a 30-step diffusion process where *RAS* is applied starting from step 4, we designate steps 12 and 20 as dense steps. During these dense steps, the entire image is processed by the model, allowing it to correct any drift that may have developed in unfocused areas. This approach ensures that the accumulated errors are reset, maintaining the denoising process in alignment with the correct direction.

### 3.4. Implementation

**Kernel Fusing.** As previously mentioned, we introduced key and value caching in the self-attention mechanism. In each attention block of the selective sampling steps, these caches are partially updated by active tokens and then used as key and value inputs for the attention functions. This partial updating operation is equivalent to a scatter operation with active token indices.

In our scenario, the source data of the scatter operation

| Method | Sample Steps | Sampling Ratio | Throughput (iter/s)↑ | FID ↓ | sFID ↓ | CLIP score ↑ |
|---|---|---|---|---|---|---|
| **Stable Diffusion 3** | | | | | | |
| RFlow | 5 | 100% | 1.43 | 39.70 | 22.34 | 29.84 |
| RAS | 7 | 25.0% | 1.45 | **31.99** | 21.70 | **30.64** |
| RAS | 7 | 12.5% | 1.48 | 32.86 | 22.10 | 30.55 |
| RAS | 6 | 25.0% | 1.52 | 33.24 | **21.51** | 30.38 |
| RAS | 6 | 12.5% | **1.57** | 33.81 | 21.62 | 30.33 |
| RFlow | 4 | 100% | 1.79 | 61.92 | 27.42 | 28.45 |
| RAS | 5 | 25.0% | 1.94 | **51.92** | **25.67** | **29.06** |
| RAS | 5 | 12.5% | **1.99** | 53.24 | 26.04 | 28.94 |
| **Lumina-Next-T2I** | | | | | | |
| RFlow | 7 | 100% | 0.49 | 48.19 | 38.60 | 28.65 |
| RAS | 10 | 25.0% | 0.59 | **45.67** | **32.36** | **29.82** |
| RAS | 10 | 12.5% | **0.65** | 47.34 | 32.69 | 29.75 |
| RFlow | 5 | 100.% | 0.69 | 96.53 | 59.26 | 26.03 |
| RAS | 7 | 25.0% | 0.70 | **53.93** | **39.80** | **28.85** |
| RAS | 7 | 12.5% | 0.74 | 54.62 | 40.23 | 28.83 |
| RAS | 6 | 25.0% | 0.75 | 67.16 | 46.46 | 27.85 |
| RAS | 6 | 12.5% | **0.78** | 67.88 | 45.88 | 27.83 |

Table 2. Pareto Improvements of rectified flow with *RAS* on COCO Val2014 1024×1024.

comprises active keys and values outputted by the previous general matrix multiplication (GeMM) kernel in the linear projection module. The extra GPU memory read/store on active keys and values can be avoided by fusing the scatter operation into the GeMM kernel, rather than launching a separate scatter kernel. Fortunately, PIT [53] demonstrates that all permutation invariant transformations, including one-dimensional scattering, can be performed in the I/O stage of GPU-efficient computation kernels (e.g. GeMM kernels) with minimal overhead. Using this method, we fused the scatter operation into the epilogue of the previous GeMM kernel.

# 4. Experiments

## 4.1. Experiment Setup

**Models, Datasets, Metrics and Baselines.** We evaluate *RAS* on Stable Diffusion 3 [10] and Lumina-Next-T2I [2] for text-to-image generation tasks, using 10,000 randomly selected caption-image pairs from the MS-COCO 2017 dataset [25]. To assess the quality of generated images and their compatibility with prompts, we use the Fréchet Inception Distance (FID) [17], the Sliding Fréchet Inception Distance (sFID) [17], and the CLIP score [16] as evaluation metrics. For baseline comparison, we evaluate *RAS* against widely-used Rectified-Flow-based Flow-Matching methods [1, 6, 10, 12, 26, 27], which uniformly reduce the number of timesteps in the generation process for the whole image. We implement *RAS* with varying numbers of total timesteps to assess its performance, and compare these configurations to

the original implementation under similar throughput conditions.

**Code Implementation.** We implement *RAS* using PyTorch [30], leveraging the diffusers library [46] and its Flow-MatchEulerDiscreteScheduler. The evaluation metrics are computed using public repositories available on GitHub [19, 40, 54]. Experiments are conducted on four servers, each equipped with eight NVIDIA A100 40GB GPUs, while speed tests are performed on an NVIDIA A100 80GB GPU.

## 4.2. Generation Benchmarks

We conducted a comparative evaluation of *RAS* and the rectified flow, which uniformly reduces the number of timesteps for every token during inference. To assess the performance of *RAS*, we performed experiments using various configurations of inference timesteps. The findings can be interpreted in two principal ways.

**Pushing the Efficiency Frontier.** From the first aspect, *RAS* offers a chance to further reduce the inference cost for each number of timesteps rectified flow offers. As illustrated in Figure 2 (c)(d), we generated 10,000 images using dense inference across different timesteps, ranging from 3 to 30. Subsequently, we applied *RAS* at varying average sampling ratios (75%, 50%, 25%, and 12.5%) over selective sampling timesteps, with the total number of timesteps set at 5, 6, 7, 10, 15, and 30. The results indicate that *RAS* can significantly reduce inference time while exerting only a minor effect on key evaluation metrics. For instance, employing *RAS* with 25% sampling over 30 timesteps improved

throughput by a factor of 2.25, with only a 22.12% increase in FID, a 26.22% increase in sFID, and a 0.065% decrease in CLIP score. Furthermore, the efficiency improvements achieved with *RAS* are attained at a lower cost compared to merely reducing the number of timesteps. Specifically, the rate of quality degradation observed when decreasing the sampling ratio of *RAS* is considerably lower than that observed when reducing the number of timesteps in dense inference, particularly when the number of timesteps is fewer than 10. This demonstrates that *RAS* constitutes a promising approach to enhancing efficiency while maintaining output quality and ensuring compatibility with prompts.

**Pareto Improvements of Uniform Sampling.** Through observation, we found that *RAS* can offer a Pareto improvement for rectified flow in many cases. We sorted some of the experimental results of Stable Diffusion 3 and Lumina-Next-T2I by throughput, and listed different configurations of *RAS* alongside the closest baseline in terms of throughput in Table 2 to provide a comprehensive comparison. The results clearly demonstrate that, for each instance of dense inference with rectified-flow-based flow matching in the table, there is almost consistently an option within *RAS* that offers higher throughput while delivering superior performance in terms of FID, sFID, and CLIP score. This highlights that, for achieving a given throughput level during DiT inference, *RAS* not only provides multiple configurations with both enhanced throughput and improved image quality, but also offers a broader parameter space for optimizing trade-offs between throughput, image quality, and compatibility with prompts.

(a) Drop Scheduling

| Method | FID ↓ | sFID ↓ | CLIP score ↑ |
|---|---|---|---|
| Default | 35.81 | 18.41 | 30.13 |
| Static Sampling Freq. | 37.92 | 19.11 | 29.98 |
| Random Dropping | 43.19 | 22.23 | 29.65 |
| W/O Error Reset | 46.10 | 24.85 | 30.41 |

(b) Key and Value Caching

| Method | Timesteps | FID ↓ | sFID ↓ | CLIP score ↑ |
|---|---|---|---|---|
| Default | 28 | 24.30 | 26.26 | 31.34 |
| W/O | 28 | 31.36 | 20.19 | 31.29 |
| Default | 10 | 35.81 | 18.41 | 30.13 |
| W/O | 10 | 32.33 | 20.21 | 30.27 |

Table 3. Ablation Study on Stable Diffusion 3. All techniques including dynamic sampling ratio, region identifying, error reset, and key & value recovery are necessary for high quality generation.

### 4.3. Human Evaluation

To evaluate whether *RAS* can enhance throughput while maintaining generation quality in real-world scenarios, we conducted a human evaluation. We randomly selected 14 prompts from the official research papers and blogs of Stable Diffusion 3 and Lumina, generating two images for each prompt: one using dense inference and the other using *RAS*, both with the same random seed and default number of timesteps. *RAS* was configured with 50% average sampling ratio during the selective sampling period. We invited 100 participants, comprising students and faculty members from 18 different universities and companies, to compare the generated images. Each participant was asked to determine whether one image was clearly better, slightly better, or of similar quality compared to the other. The order of the images was randomized, and participants were unaware of which image was generated with *RAS*. As shown in Figure 2 (e), 633 out of 1400 votes (45.21%) indicated that the two images were of similar quality. Additionally, 28.29% of votes favored the dense image over the *RAS* result, while 26.50% preferred *RAS* over the dense result. These results demonstrate that *RAS* achieves a significant improvement in throughput (1.625× for Stable Diffusion 3 and 1.561× for Lumina-Next-T2I) without noticeably affecting human preference.

### 4.4. Ablation Study

**Token Drop Scheduling.** As shown in Table 3 (a), we evaluate the scheduling configurations introduced in Section 3, including sampling ratio scheduling, selection of cached tokens, and the insertion of dense steps during the selective sampling period to reset accumulated errors, using 10 timesteps with an average sampling ratio of 12.5% on Stable Diffusion 3. The results indicate that each of these techniques contributes to the overall quality of *RAS*.

**Key and Value Caching.** As shown in Table 3 (b), caching keys and values from the previous step is crucial, especially when generating high-quality images with more timesteps. While dropping the keys and values of non-activated tokens during attention can improve throughput, it significantly affects the attention scores of activated tokens. A token's low ranking in the model output does not necessarily mean it has no contribution to the attention scores of other tokens.

## 5. Conclusions and Limitations

In this paper, we observed that different regions within an image require varying levels of refinement during the diffusion process, and that adjacent sampling steps exhibit significant continuity in the distribution of focused areas. Based on these observations, we proposed *RAS*, a novel diffusion sampling strategy that dynamically adjusts sampling rates according to regional attention, thereby allocating computational resources more efficiently to areas of greater importance while reusing noise predictions for less critical regions. Our approach effectively reduces computational costs while preserving high image quality. Extensive

experiments and user studies demonstrate that *RAS* achieves substantial speed-ups with minimal degradation in quality, outperforming uniform sampling baselines and paving the way for more efficient and adaptive diffusion models.

# References

[1] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022. 1, 4, 7

[2] Peng Gao andå Le Zhuo, Dongyang Liu, Ruoyi Du, Xu Luo, Longtian Qiu, Yuhang Zhang, Chen Lin, Rongjie Huang, Shijie Geng, Renrui Zhang, Junlin Xi, Wenqi Shao, Zhengkai Jiang, Tianshuo Yang, Weicai Ye, He Tong, Jingwen He, Yu Qiao, and Hongsheng Li. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers, 2024. 3, 4, 7

[3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 1

[4] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-$\sigma$: Weak-to-strong training of diffusion transformer for 4k text-to-image generation, 2024. 4

[5] Pengtao Chen, Mingzhu Shen, Peng Ye, Jianjian Cao, Chongjun Tu, Christos-Savvas Bouganis, Yiren Zhao, and Tao Chen. $\Delta$-dit: A training-free acceleration method tailored for diffusion transformers. *arXiv preprint arXiv:2406.01125*, 2024. 1, 4

[6] Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space, 2023. 7

[7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1

[8] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2024. Curran Associates Inc. 1, 3

[9] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7346–7356, 2023. 1

[10] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. 4, 7

[11] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. 3

[12] Johannes S Fischer, Ming Gui, Pingchuan Ma, Nick Stracke, Stefan A Baumann, and Björn Ommer. Boosting latent diffusion with flow matching. *arXiv preprint arXiv:2312.07360*, 2023. 7

[13] Sicheng Gao, Xuhui Liu, Bohan Zeng, Sheng Xu, Yanjing Li, Xiaoyan Luo, Jianzhuang Liu, Xiantong Zhen, and Baochang Zhang. Implicit diffusion models for continuous super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10021–10030, 2023. 1

[14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 3

[15] Philip Hartman. *Ordinary differential equations*. SIAM, 2002. 1

[16] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 7

[17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 7

[18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2020. Curran Associates Inc. 1, 2, 3

[19] Tao Hu. pytorch-fid-with-sfid. https://github.com/dongzhuoyao/pytorch-fid-with-sfid, 2022. 7

[20] Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. *arXiv preprint arXiv:2205.09991*, 2022. 1

[21] Kalervo Järvelin and Jaana Kekäläinen. Ir evaluation methods for retrieving highly relevant documents. In *Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 41–48, New York, NY, USA, 2000. Association for Computing Machinery. 3

[22] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023. 1

[23] Bo Li, Kaitao Xue, Bin Liu, and Yu-Kun Lai. Bbdm: Image-to-image translation with brownian bridge diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern Recognition*, pages 1952–1961, 2023. 1

[24] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022. 1

[25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence

Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 7

[26] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 1, 4, 7

[27] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 1, 4, 7

[28] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022. 1, 4

[29] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022. 1

[30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. *PyTorch: an imperative style, high-performance deep learning library*. Curran Associates Inc., Red Hook, NY, USA, 2019. 7

[31] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4172–4182, 2023. 2, 4

[32] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. 2, 3, 5

[33] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1

[34] Philip E Protter and Philip E Protter. *Stochastic differential equations*. Springer, 2005. 1

[35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1

[36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 2, 4

[37] Amirmojtaba Sabour, Sanja Fidler, and Karsten Kreis. Align your steps: Optimizing sampling schedules in diffusion models. In *Forty-first International Conference on Machine Learning*, 2023. 1

[38] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 conference proceedings*, pages 1–10, 2022. 1

[39] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*. 1, 4

[40] Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. https://github.com/mseitzer/pytorch-fid, 2020. Version 0.3.0. 7

[41] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, page 2256–2265. JMLR.org, 2015. 1, 3

[42] Yang Song and Stefano Ermon. *Generative modeling by estimating gradients of the data distribution*. Curran Associates Inc., Red Hook, NY, USA, 2019. 1, 3

[43] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023. 1, 4

[44] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 5

[45] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 2, 4, 5

[46] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers, 2022. 7

[47] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*, 2022. 1

[48] Yining Wang, Liwei Wang, Yuanzhi Li, Di He, Tie-Yan Liu, and Wei Chen. A theoretical analysis of ndcg type ranking measures, 2013. 3

[49] Mengwei Xu, Mengze Zhu, Yunxin Liu, Felix Xiaozhu Lin, and Xuanzhe Liu. Deepcache: Principled cache for mobile deep vision. In *Proceedings of the 24th annual international conference on mobile computing and networking*, pages 129–144, 2018. 1, 4

[50] Zhihang Yuan, Hanling Zhang, Pu Lu, Xuefei Ning, Linfeng Zhang, Tianchen Zhao, Shengen Yan, Guohao Dai, and Yu Wang. Ditfastattn: Attention compression for diffusion transformer models, 2024. 4

[51] Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. *Advances in Neural Information Processing Systems*, 36, 2024. 1

[52] Zhixing Zhang, Ligong Han, Arnab Ghosh, Dimitris N Metaxas, and Jian Ren. Sine: Single image editing with text-to-image diffusion models. In *Proceedings of the IEEE/CVF*

*Conference on Computer Vision and Pattern Recognition*, pages 6027–6037, 2023. 1

[53] Ningxin Zheng, Huiqiang Jiang, Quanlu Zhang, Zhenhua Han, Lingxiao Ma, Yuqing Yang, Fan Yang, Chengruidong Zhang, Lili Qiu, Mao Yang, et al. Pit: Optimization of dynamic sparse deep learning models via permutation invariant transformation. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 331–347, 2023. 7

[54] SUN Zhengwentai. clip-score: CLIP Score for PyTorch. https://github.com/taited/clip-score, 2023. Version 0.1.1. 7

# Region-Adaptive Sampling for Diffusion Transformers

## Supplementary Material

## 6. More Visualization of *RAS*

This section presents *RAS* accelerating Lumina-Next-T2I and Stable Diffusion 3 with a 50% sampling ratio. As illustrated in Figure 10, the main object receives more sampling steps compared to the background, demonstrating the significance of our region-adaptive sampling strategy. This approach ensures that the primary subject in the generated image consistently undergoes more sampling, while relatively smooth regions receive fewer sampling steps. For instance, in the example shown in Figure 10 with the prompt "hare in snow," the weeds in the snow are sampled more frequently, while the smooth snow receives fewer sampling steps.

In Figure 11, we visualize the standard deviation of the noise across dimensions, as well as the decoded images derived from the noise. This stems from our observation that the noise's standard deviation is consistently smaller in the main subject areas. A preliminary hypothesis is that this occurs because the main subject contains more information. When mixed with a certain proportion of noise at each diffusion step, the foreground tends to retain more deterministic information compared to the background. This allows the model to predict more consistent denoising directions. We acknowledge that further study is needed to fully understand this phenomenon.

The primary contribution of this work is to highlight that employing different sampling steps for different regions can significantly enhance the efficiency of diffusion model sampling. The method for selecting these regions is not limited to the aforementioned approach based on the noise standard deviation across dimensions. For example, we also experimented with using the $l-2$ norm of the noise output by the network as a criterion for selection. By targeting regions with larger noise norms, which indicate areas the network deems requiring more refinement, we observed a preference for more complex regions in the frequency domain as in Figure 9. This approach also achieves high-quality imaging results, as shown in Table 4. It can be seen that the methods using the $l_2$ norm and standard deviation (std) yield relatively similar results, and both significantly outperform random selection, particularly when the cache ratio is higher.



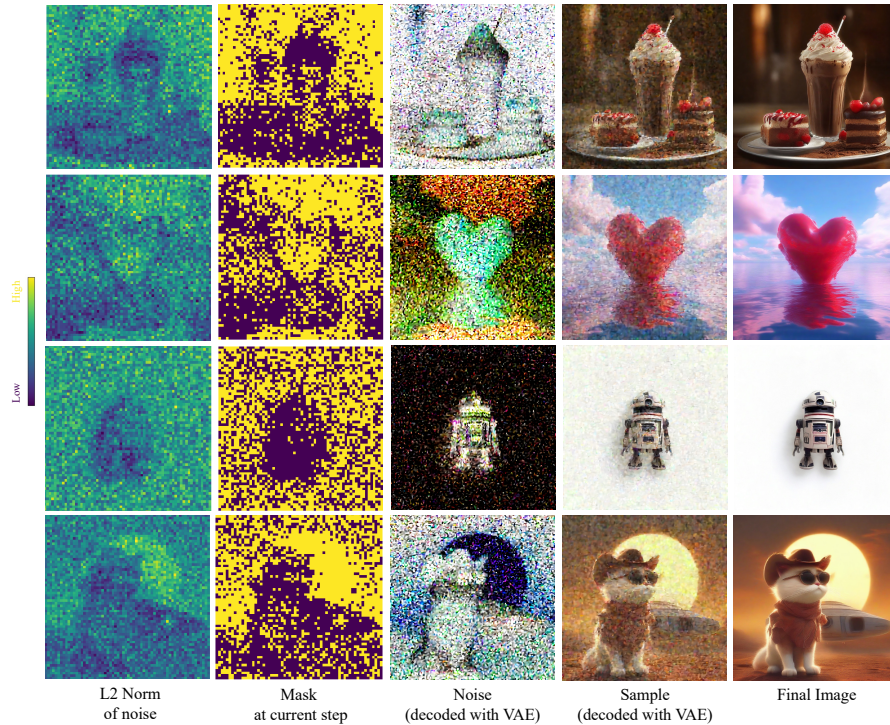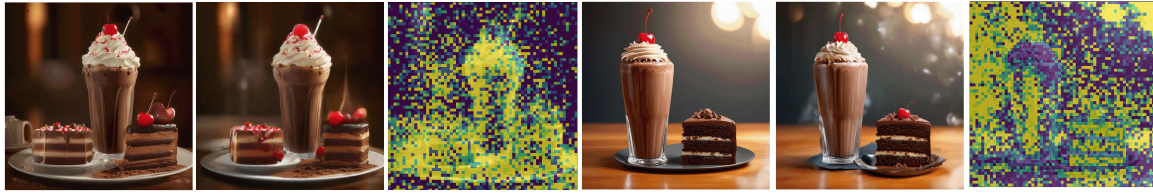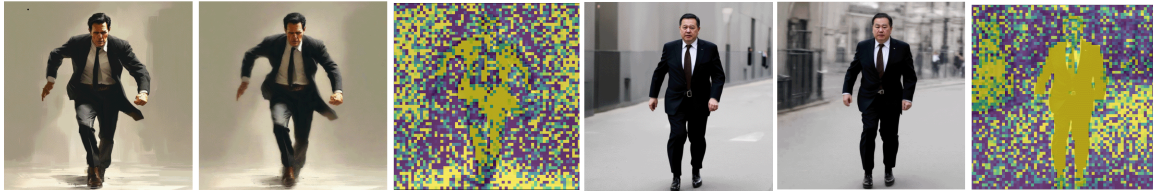| L2 Norm of noise | Mask at current step | Noise (decoded with VAE) | Sample (decoded with VAE) | Final Image |

Figure 9. *RAS* using norm as the metric, accelerating Lumina-Next-T2I with 50% sample ratio and 30 total steps. The noise, masks and samples are from the 20th step.
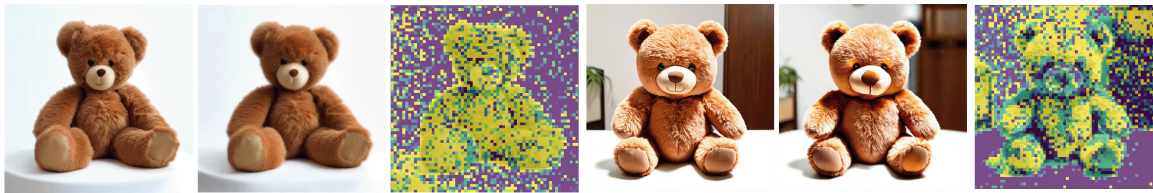
Figure 10. *RAS* VS default sampling and the active sampling step for each latent token.

| Method | Sample Steps | Sampling Ratio | Throughput (iter/s)↑ | FID ↓ | sFID ↓ | CLIP score ↑ |
|--------|-------------|----------------|---------------------|-------|--------|--------------|
| RFlow | 7 | 100.0% | 1.01 | 27.23 | 17.76 | 30.87 |
| RAS-Std | 7 | 25.0% | 1.45 | 31.99 | 21.7 | 30.64 |
| RAS-Norm | 7 | 25.0% | 1.45 | 31.65 | 21.24 | 30.59 |
| Random | 7 | 25.0% | 1.45 | 33.26 | 22.10 | 30.67 |

Table 4. Experiments on using L2 Norm as the metric for *RAS* on Stable Diffusion 3. The sample ratio of the first 4 steps is 100% to guarantee generation qualities.



Figure 11. The 20th sampling step (out of 30) of Lumina-Next-T2I using *RAS*.

# 7. Full Experiment Results of *RAS*

In this section, we present the full experiment results of *RAS* against rectified flow, with the same settings as is described in the experiment section. Both Table 5 and 6 are ordered by the throughputs.

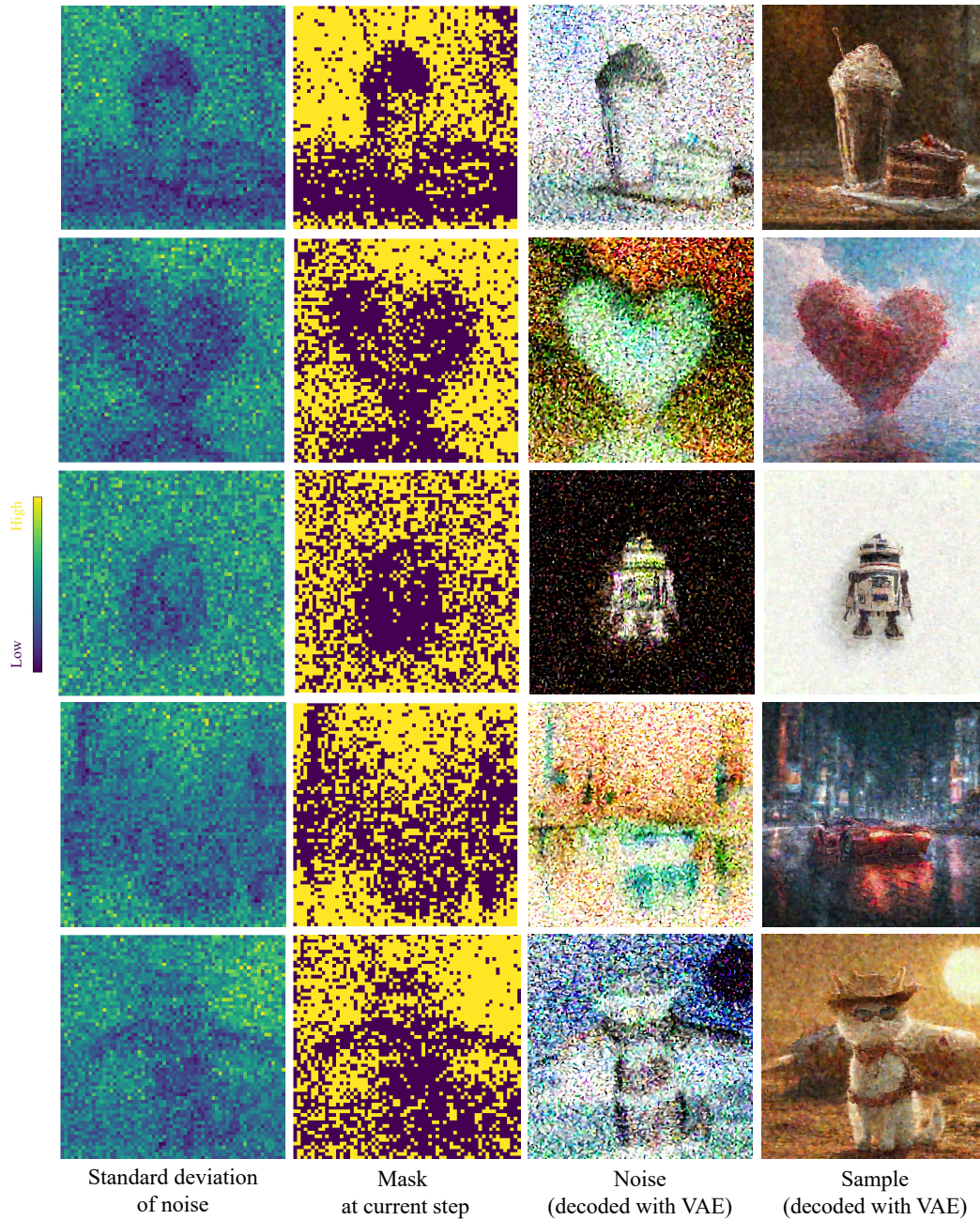| Method | Sample Steps | Sampling Ratio | Throughput (iter/s)↑ | FID ↓ | sFID ↓ | CLIP score ↑ |
|--------|--------------|----------------|----------------------|-------|--------|--------------|
| RFlow | 30 | 100.0% | 0.11 | 22.46 | 16.59 | 30.47 |
| RAS | 30 | 75.0% | 0.14 | 23.31 | 17.73 | 30.49 |
| RFlow | 23 | 100.0% | 0.15 | 23.10 | 17.91 | 30.42 |
| RAS | 30 | 50.0% | 0.18 | 24.10 | 18.83 | 30.51 |
| RFlow | 15 | 100.0% | 0.23 | 24.88 | 21.02 | 30.25 |
| RAS | 30 | 25.0% | 0.26 | 27.44 | 20.95 | 30.45 |
| RAS | 15 | 75.0% | 0.27 | 26.82 | 23.33 | 30.26 |
| RAS | 30 | 12.5% | 0.31 | 33.64 | 23.44 | 30.36 |
| RAS | 15 | 50.0% | 0.33 | 28.48 | 25.17 | 30.29 |
| RFlow | 10 | 100.0% | 0.34 | 31.35 | 27.84 | 29.74 |
| RAS | 10 | 75.0% | 0.40 | 34.19 | 30.57 | 29.79 |
| RAS | 15 | 25.0% | 0.43 | 33.28 | 27.41 | 30.24 |
| RAS | 15 | 12.5% | 0.48 | 39.75 | 28.88 | 30.14 |
| RAS | 10 | 50.0% | 0.48 | 36.18 | 32.36 | 29.86 |
| RFlow | 7 | 100.0% | 0.49 | 48.19 | 38.60 | 28.65 |
| RAS | 7 | 75.0% | 0.54 | 50.45 | 40.19 | 28.78 |
| RAS | 10 | 25.0% | 0.59 | 42.96 | 33.51 | 29.91 |
| RAS | 7 | 50.0% | 0.61 | 51.78 | 40.51 | 28.82 |
| RAS | 6 | 75.0% | 0.62 | 66.12 | 46.58 | 27.80 |
| RAS | 10 | 12.5% | 0.65 | 47.34 | 32.70 | 29.75 |
| RAS | 6 | 50.0% | 0.67 | 66.54 | 46.71 | 27.83 |
| RAS | 7 | 25.0% | 0.70 | 53.93 | 39.80 | 28.85 |
| RAS | 7 | 12.5% | 0.74 | 54.62 | 40.23 | 28.83 |
| RAS | 6 | 25.0% | 0.74 | 67.16 | 46.46 | 27.85 |
| RAS | 5 | 75.0% | 0.75 | 99.01 | 56.26 | 26.02 |
| RAS | 6 | 12.5% | 0.78 | 67.88 | 45.89 | 27.83 |
| RFlow | 5 | 100.0% | 0.69 | 96.53 | 59.26 | 26.03 |
| RAS | 5 | 50.0% | 0.83 | 99.81 | 56.57 | 26.01 |
| RAS | 5 | 25.0% | 0.95 | 101.50 | 56.40 | 25.93 |
| RAS | 5 | 12.5% | 1.00 | 102.90 | 55.25 | 25.84 |
| RFlow | 3 | 100.0% | 1.15 | 256.90 | 94.80 | 19.67 |

Table 5. Full experiment results of *RAS* and rectified flow on Lumina-Next-T2I and COCO Val2014 1024×1024.

| Method | Sample Steps | Sampling Ratio | Throughput (iter/s)↑ | FID ↓ | sFID ↓ | CLIP score ↑ |
|---|---|---|---|---|---|---|
| RFlow | 28 | 100% | 0.26 | 25.8 | 15.32 | 31.4 |
| RAS | 28 | 75.0% | 0.33 | 24.43 | 15.94 | 31.39 |
| RAS | 28 | 50.0% | 0.42 | 24.86 | 16.88 | 31.36 |
| RFlow | 14 | 100% | 0.51 | 24.49 | 14.78 | 31.34 |
| RAS | 28 | 25.0% | 0.55 | 25.16 | 17.11 | 31.29 |
| RFlow | 12 | 100% | 0.59 | 24.36 | 14.89 | 31.3 |
| RAS | 14 | 75.0% | 0.62 | 23.61 | 15.92 | 31.35 |
| RAS | 28 | 12.5% | 0.63 | 25.72 | 17.3 | 31.22 |
| RFlow | 10 | 100% | 0.71 | 24.17 | 15.39 | 31.22 |
| RAS | 14 | 50.0% | 0.74 | 24.6 | 17.24 | 31.32 |
| RAS | 14 | 25.0% | 0.91 | 25.88 | 17.97 | 31.24 |
| RAS | 10 | 75.0% | 0.91 | 24.39 | 16.29 | 31.12 |
| RAS | 14 | 12.5% | 0.98 | 26.48 | 18.14 | 31.18 |
| RAS | 10 | 50.0% | 1.0 | 27.1 | 17.5 | 30.93 |
| RFlow | 7 | 100% | 1.01 | 27.23 | 17.76 | 30.87 |
| RAS | 7 | 75.0% | 1.16 | 27.57 | 18.76 | 30.81 |
| RAS | 10 | 25.0% | 1.2 | 30.97 | 18.36 | 30.67 |
| RAS | 10 | 12.5% | 1.3 | 35.81 | 18.41 | 30.13 |
| RAS | 7 | 50.0% | 1.3 | 30.04 | 20.34 | 30.73 |
| RAS | 6 | 75.0% | 1.3 | 31.23 | 19.98 | 30.48 |
| RAS | 6 | 50.0% | 1.41 | 32.21 | 20.86 | 30.43 |
| RFlow | 5 | 100% | 1.43 | 39.7 | 22.34 | 29.84 |
| RAS | 7 | 25.0% | 1.45 | 31.99 | 21.7 | 30.64 |
| RAS | 7 | 12.5% | 1.48 | 32.86 | 22.1 | 30.55 |
| RAS | 6 | 25.0% | 1.52 | 33.24 | 21.51 | 30.36 |
| RAS | 6 | 12.5% | 1.57 | 33.81 | 21.62 | 30.33 |
| RAS | 5 | 75.0% | 1.59 | 44.02 | 23.14 | 29.53 |
| RAS | 5 | 50.0% | 1.75 | 48.65 | 24.51 | 29.29 |
| RFlow | 4 | 100% | 1.79 | 61.92 | 27.42 | 28.45 |
| RAS | 5 | 25.0% | 1.94 | 51.92 | 25.67 | 29.06 |
| RAS | 5 | 12.5% | 1.99 | 53.24 | 26.04 | 28.94 |
| RFlow | 3 | 100% | 2.38 | 121.61 | 36.92 | 25.32 |

Table 6. Full experiment results of *RAS* and rectified flow on Stable Diffusion 3 and COCO Val2014 1024×1024.