# Dynamic Attention-Guided Diffusion for Image Super-Resolution

Brian B. Moser[1,2,3], Stanislav Frolov[1,2,3], Federico Raue[1], Sebastian Palacio[1], and Andreas Dengel[1,2]

[1] German Research Center for Artificial Intelligence (DFKI), Germany
[2] RPTU Kaiserslautern-Landau, Germany
[3] Equal Contribution
`first.second@dfki.de`

**Abstract.** Diffusion models in image Super-Resolution (SR) treat all image regions with uniform intensity, which risks compromising the overall image quality. To address this, we introduce "You Only Diffuse Areas" (YODA), a dynamic attention-guided diffusion method for image SR. YODA selectively focuses on spatial regions using attention maps derived from the low-resolution image and the current time step in the diffusion process. This time-dependent targeting enables a more efficient conversion to high-resolution outputs by focusing on areas that benefit the most from the iterative refinement process, i.e., detail-rich objects. We empirically validate YODA by extending leading diffusion-based methods SR3 and SRDiff. Our experiments demonstrate new state-of-the-art performance in face and general SR across PSNR, SSIM, and LPIPS metrics. A notable finding is YODA's stabilization effect by reducing color shifts, especially when training with small batch sizes.

## 1 Introduction

The goal of image Super-Resolution (SR) is to enhance Low-Resolution (LR) into High-Resolution (HR) images [35]. This process has a significant impact in fields such as medical imaging, remote sensing, and consumer electronics [16, 33, 44]. Despite its long history, image SR remains a fascinating yet challenging domain due to its inherently ill-posed nature: any given LR image can lead to several valid HR images, and vice versa [2, 40]. Thanks to deep learning, image SR has made significant progress [13]. Initial regression-based methods, such as early convolutional neural networks, work great at low magnification ratios [7, 27, 43]. However, they fail to produce high-frequency details at high magnification ratios and generate over-smoothed images [34].

Recently, diffusion models have emerged with better human-rated quality compared to regression-based methods, but they also introduce new challenges [9, 21, 38, 46]. Their indiscriminate processing of image regions leads to computational redundancies and suboptimal enhancements. Some recent methods have reduced computational demands by working in latent space like LDMs [37], by

exploiting the relationship between LR and HR latent representations like Part-Diff [48], or by starting with a better-initialized forward diffusion instead of pure Gaussian noise like in CCDF [10]. However, there is still a lack of strategies that adapt the model capacity based on the spatial importance of image regions.

In this paper, we question the common approach of diffusion models for image SR: Is it necessary to update the entire image at every time step in the reverse diffusion process? We hypothesize that not all image regions require the same level of detail enhancement. For instance, a face in the foreground may need more attention than a simple, monochromatic background. Recognizing this variability in the need for detail enhancement underscores a critical inefficiency in traditional diffusion methods. Treating all image regions with uniform intensity risks compromising the overall image quality. To bridge this gap, we introduce "You Only Diffuse Areas" (YODA), an efficient diffusion mechanism focusing on detail-rich areas using time-dependent and attention-guided masking. YODA starts by obtaining an attention map, highlighting regions that need more refinement. After identification, YODA systematically replaces the highlighted regions with SR predictions during the denoising process, depending on the current time step. As a result, detail-rich areas are refined more often. Our approach is analogous to inpainting methods like RePaint [31], where only a pre-defined masked region is updated to generate complementing content. In YODA's case, however, the selection of regions to update is time-dependent. We design a dynamic approach that creates expanding masks, starting from detail-rich regions and converging towards enhancing the overall image. A key advantage of YODA is its compatibility with existing diffusion models, allowing for a plug&play application. We integrate YODA with SR3 [38] for face SR and SRDiff [25] for general SR, achieving notable improvements in image quality. Interestingly, YODA also improves the training process. For instance, our face SR experiments require a significantly reduced batch size due to limited hardware access. When training with smaller batch sizes, SR3 also suffers from color shifts [8, 42], while YODA produces color distributions faithfully. In summary, our work:

- introduces YODA, an attention-guided diffusion approach that emphasizes important image areas through masked refinement. As a result, it refines detail-rich areas more often, which leads to higher overall image quality.
- demonstrates that attention-guided diffusion results in better training conditions, accurate color predictions, and better perceptual quality.
- empirically shows that YODA outperforms leading diffusion models in face and general SR tasks.
- shows that YODA improves the training performance when using smaller batch sizes, which is crucial in limited hardware scenarios.

## 2    Background

Our method uses attention maps for attention-guided diffusion. Thus, this section introduces the main components: DDPMs [21] and the DINO framework [6].

### 2.1   DDPMs

Denoising Diffusion Probabilistic Models (DDPMs) employ two distinct Markov chains [21]: the first models the forward diffusion process $q$ transitioning from an input $\mathbf{x}$ to a pre-defined prior distribution with intermediate states $\mathbf{z}_t$, $0 < t \leq T$, while the second models the backward diffusion process $p$, reverting from the prior distribution back to the intended target distribution $p(\mathbf{z}_0 \mid \mathbf{z}_T, \mathbf{x})$. In image SR, we designate $\mathbf{x}$ as the LR image and the target $\mathbf{z}_0$ as the desired HR image. The prior distribution is generally set manually, e.g., Gaussian noise.

**Forward Diffusion**  In the forward diffusion phase, an HR image $\mathbf{z}_0$ is incrementally modified by adding Gaussian noise over a series of time steps. This process can be mathematically represented as:

$$q(\mathbf{z}_t \mid \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t \mid \sqrt{1 - \alpha_t}\, \mathbf{z}_{t-1}, \alpha_t \mathbf{I}) \tag{1}$$

The hyperparameters $0 < \alpha_{1:T} < 1$ represent the noise variance injected at each time step. It is possible to sample from any point in the noise sequence without needing to generate all previous steps through the following simplification [39]:

$$q(\mathbf{z}_t \mid \mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t \mid \sqrt{\gamma_t}\, \mathbf{z}_0, (1 - \gamma_t)\mathbf{I}), \tag{2}$$

where $\gamma_t = \prod_{i=1}^{t}(1 - \alpha_i)$ . The intermediate step $\mathbf{z}_t$ is derived by:

$$\mathbf{z}_t = \sqrt{\gamma_t} \cdot \mathbf{z}_0 + \sqrt{1 - \gamma_t} \cdot \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{3}$$

**Backward Diffusion**  The backward diffusion process is where the model learns to denoise, effectively reversing the forward diffusion to recover the HR image. In image SR, the reverse process is conditioned on the LR image to guide the generation of the HR image:

$$p_\theta(\mathbf{z}_{t-1} \mid \mathbf{z}_t, \mathbf{x}) = \mathcal{N}(\mathbf{z}_{t-1} \mid \mu_\theta(\mathbf{z}_t, \mathbf{x}, \gamma_t), \Sigma_\theta(\mathbf{z}_t, \mathbf{x}, \gamma_t)) \tag{4}$$

The mean $\mu_\theta$ depends on a parameterized denoising function $f_\theta$, which can either predict the added noise $\varepsilon_t$ or the underlying HR image $\mathbf{z}_0$. Following the standard approach of Ho et al. [21], we focus on predicting the noise. Hence, the mean is:

$$\mu_\theta(\mathbf{x}, \mathbf{z}_t, \gamma_t) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{z}_t - \frac{1 - \alpha_t}{\sqrt{1 - \gamma_t}} f_\theta(\mathbf{x}, \mathbf{z}_t, \gamma_t)\right) \tag{5}$$

Following Saharia et al. [38], setting the variance of $p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})$ to $(1 - \alpha_t)$ yields the subsequent refining step with $\varepsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$:

$$\mathbf{z}_{t-1} \leftarrow \mu_\theta(\mathbf{x}, \mathbf{z}_t, \gamma_t) + \sqrt{1 - \alpha_t}\,\varepsilon_t \tag{6}$$

**Optimization**  The optimization goal for DDPMs is to train the parameterized model to accurately predict the noise added during the diffusion process. The loss function used to measure the accuracy of the noise prediction is:

$$\mathcal{L}(\theta) = \underset{(\mathbf{x}, \mathbf{z}_0)}{\mathbb{E}}\, \underset{t}{\mathbb{E}} \left\| \varepsilon_t - f_\theta(\mathbf{x}, \mathbf{z}_t, \gamma_t) \right\|_1 \tag{7}$$

## 2.2   DINO

DINO is a self-supervised learning approach for feature extractors on unlabeled data [6]. It employs a teacher and a student network, where the student learns to imitate the features learned by the teacher. The student gets only local views of the image (i.e., $96 \times 96$), whereas the teacher receives global views (i.e., $224 \times 224$). This setup encourages the student to learn "local-to-global" correspondences. The features learned through self-supervision are directly accessible in the self-attention modules. These self-attention maps provide information on the scene layout and important object boundaries. We leverage the generality, availability, and robustness of these attention maps as a measure of an image's saliency to guide the diffusion process. In another context, a similar approach has been applied to image compression, demonstrating its ability to capture essential image content in the attention maps [3]. More details can be found in the supplementary material.

## 3   Methodology

This section describes the components of our proposed method, "You Only Diffuse Areas" (YODA). YODA is divided into three main phases:

– **Identifying Key Regions:** Estimate the importance of pixel positions in a LR image $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$ with an attention map $\mathbf{A} \in \mathbb{R}^{H \times W}$.
– **Time-Dependent Masking:** Using $\mathbf{A}$, define time-dependent, binary masks $\mathbf{M} : \mathbb{N}_0 \to \{0, 1\}^{H \times W}$ to identify salient areas at time step $t \in \mathbb{N}_0$ of the diffusion process.
– **Guided Backward Diffusion:** Condition the diffusion process on the regions identified by the time-dependent masking $\mathbf{M}(t)$.

## 3.1   Identifying Key Regions

Super-resolving with YODA begins by identifying areas of an image with greater details. This is achieved by generating an attention map $\mathbf{A}$ with $0 \leq \mathbf{A}_{i,j} \leq 1$ from the LR image $\mathbf{x}$. The greater the value of $\mathbf{A}_{i,j}$, the more important it is. Note that $\mathbf{A}$ has to be generated only once for each image.

For generating $\mathbf{A}$, we evaluated several approaches, including innate methods (i.e., not-learnable) and learnable methods, i.e., ResNet and Transformer architectures [14, 20]. For the latter, we leverage the DINO framework for its robustness in self-supervised learning, extracting refined attention maps directly from LR images without necessitating extra annotated data [6]. This choice is motivated by DINO's demonstrated efficacy in highlighting essential features within images using pre-existing models, e.g., for image compression [3]. An additional overview of how DINO and the attention maps are used is shown in the supplementary materials. Next, we describe the process of creating time-dependent masks for backward diffusion, utilizing the attention map $\mathbf{A}$.

## 3.2 Time-Dependent Masking

Given the LR input image $\mathbf{x}$ and the attention map $\mathbf{A}$, we introduce a strategy to focus the diffusion process to detail-rich areas. Each position in $\mathbf{A}$ indicates its semantic importance, influencing the number of refinement steps it receives. Thus, for two positions $(i, j)$ and $(i', j')$ with $\mathbf{A}_{i,j} > \mathbf{A}_{i',j'}$, YODA applies more refinement steps to the location $(i, j)$ than to $(i', j')$. Since $0 \leq \mathbf{A}_{i,j} \leq 1$, the number of diffusion steps employed to a specific position $(i, j)$ is determined as a proportion of the maximum time steps, $T$. For instance, $\mathbf{A}_{i,j} = 0.7$ means $(i, j)$ is refined during 70% of all diffusion steps. In addition, a lower bound hyperparameter $0 < l < 1$ ensures that every region undergoes a minimum level of refinement. In other words, the hyperparameter $l$ guarantees that every spatial position is refined at least $l \cdot T$ times. Note that the backward diffusion process starts from time step $T$ and ends in time step 0. We define the time-dependent masking process at time step $0 \leq t \leq T$ with $t \to 0$ as follows:

$$\mathbf{M}(t)_{i,j} = \begin{cases} 1, \text{ if } T \cdot (\mathbf{A}_{i,j} + l) \geq t \\ 0, \text{ otherwise} \end{cases} \tag{8}$$

This equation ensures that the diffusion process gets applied a variable number of times for different regions, allowing the salient areas to diffuse over a longer time span. It is important to highlight that once a spatial position is marked for refinement, it continues to undergo refinement across all subsequent steps as $t$ approaches 0: $\mathbf{M}(t)_{i,j} \geq \mathbf{M}(t - k)_{i,j} \, \forall k > 0$. Figure 1 shows an example of our time-dependent masking. For each time step $t$, we can determine with $\mathbf{M}(t)_{i,j} = 1$ whether a given spatial position $(i, j)$ should be refined or not.

## 3.3 Guided Backward Diffusion

YODA's guided diffusion process iteratively refines the image from a noisy state $\mathbf{z}_T$ to a high-resolution state $\mathbf{z}_0$. This phase involves selectively refining areas based on the current time step's mask, $\mathbf{M}(t)$, and blending these refined areas with the unrefined, remaining LR regions. YODA ensures a seamless transition between refined and unrefined areas, improving image quality with a focus on key regions.
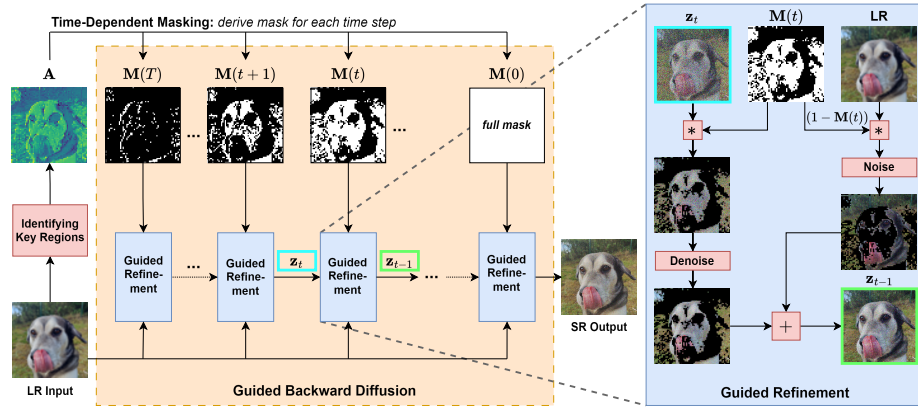
At each time step $t$, the areas that will be refined for transitioning from $t$ to $(t - 1)$ are determined based on the current iteration $\mathbf{z}_t$ and the current mask $\mathbf{M}(t)$:

$$\widetilde{\mathbf{z}}_t \leftarrow \mathbf{M}(t) \odot \mathbf{z}_t \tag{9}$$

Next, we divide the image into two components that will later be combined as the output for the next time step: $\mathbf{z}_{t-1}^{SR}$, which is the refined image prediction, and $\mathbf{z}_{t-1}^{LR}$, the complementary LR image. The state $\mathbf{z}_{t-1}^{LR}$ represents unchanged LR areas but is sampled using $\mathbf{x}$ as the mean. Both components acquire the same noise level $\Sigma_\theta(\widetilde{\mathbf{z}}_t, \mathbf{x}, \gamma_t)$, and can be described by:

$$\mathbf{z}_{t-1}^{SR} \sim \mathcal{N}\left(\mu_\theta(\widetilde{\mathbf{z}}_t, \mathbf{x}, \gamma_t), \Sigma_\theta(\widetilde{\mathbf{z}}_t, \mathbf{x}, \gamma_t)\right) \tag{10}$$

$$\mathbf{z}_{t-1}^{LR} \sim \mathcal{N}\left(\mathbf{x}, \Sigma_\theta(\widetilde{\mathbf{z}}_t, \mathbf{x}, \gamma_t)\right) \tag{11}$$

**Fig. 1:** Overview of YODA. First, extract an attention map $\mathbf{A}$ from the LR input. Next, use the values of $\mathbf{A}$ to produce a time-dependent masking $\mathbf{M}(t)$. For $t : T \to 0$, the area of selected pixels expands from detail-rich regions to the whole image. Our diffusion process uses these masks for dynamic and attention-guided refinement, emphasizing important regions. More specifically, it starts with masked areas that need refinement (derived from $\mathbf{z}_t$ and $\mathbf{M}(t)$) and LR regions, which retain the noise level needed for the next time step. Finally, the SR and LR areas are combined to form a whole image with no masked-out regions for the next iteration.

Finally, YODA combines the complementing, non-overlapping image regions[4]:

$$\mathbf{z}_{t-1} \leftarrow \mathbf{M}(t) \odot \mathbf{z}_{t-1}^{SR} + (1 - \mathbf{M}(t)) \odot \mathbf{z}_{t-1}^{LR} \tag{12}$$

Consequently, the areas refined by $\mathbf{z}_t^{SR}$ expand as $t \to 0$, whereas the areas described by $\mathbf{z}_t^{LR}$ shrink in size. The new state, $\mathbf{z}_{t-1}$, now contains both SR and LR areas and, importantly, does not have any masked-out regions. As a result, $\mathbf{z}_{t-1}$ can be used in the next iteration step. This guided refinement is depicted on the right part of Figure 1.
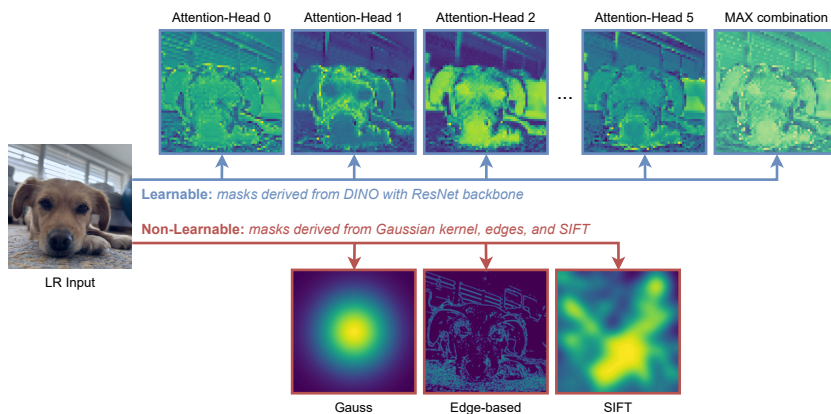
### 3.4 Optimization

To confine the backward diffusion process to specific image regions as determined by the current time step $0 \leq t \leq T$ and the corresponding mask $\mathbf{M}(t)$, YODA improves the training objective Equation 7 as follows:

$$\mathcal{L}\left(\theta\right) = \mathop{\mathbb{E}}_{(\mathbf{x},\mathbf{y})} \mathop{\mathbb{E}}_{t} \left\| \mathbf{M}(t) \odot \left[\varepsilon_t - f_\theta\left(\mathbf{x}, \mathbf{z}_t, \gamma_t\right)\right] \right\|_1 \tag{13}$$

The loss function focuses on the highlighted spatial regions within the mask $\mathbf{M}(t)$. Consequently, YODA optimizes only areas described by $\mathbf{M}(t)$.

---

[4] This formulation is similar to RePaint [31], a method for diffusion-based inpainting. While RePaint uses a constant mask for all time steps, YODA uses a time-dependent procedure that enables dynamic control over the refined image regions.

**Fig. 2:** Comparison of attention maps (blue = low attention; yellow = high attention). Top row denotes maps derived from ResNet-50 using DINO. It shows various attention head outputs and the max aggregation of all attention maps (MAX). Bottom row denotes non-learnable methods, namely Gaussian, Edge-based, and using SIFT's points of interest. The maps are subsequently used to produce time-dependent binary masks.

## 4 Experiments

We start by analyzing different methods for obtaining attention maps for YODA. Then, we evaluate YODA and compare its performance in tandem with SR3 [38] for face-only, and SRDiff [25] for general SR. We chose SR3 and SRDiff because they are the most prominent representative diffusion models for image SR in the respective tasks, where YODA can be integrated straightforwardly. However, YODA can be theoretically applied to any existing method. We present quantitative and qualitative results for both tasks. Our method achieves high-quality results for both tasks and outperforms the baselines using standard metrics such as PSNR, SSIM, and LPIPS [34]. All experiments were run on a single NVIDIA A100-80GB GPU. In the supplementary materials, we discuss the complexity of YODA and explore its potential synergies with other diffusion model research.

### 4.1 Analysis of Attention Maps

Since the dynamic masking in YODA relies on attention maps, we thoroughly evaluate different choices. We considered the pre-trained attention heads from the last layer of DINO with the respective neural network, i.e., ResNet and ViT [14,20]. For ResNet-50, we used a dedicated method to extract the attention maps from its weights [19]. A qualitative visual comparison of attention maps generated with DINO and ResNet-50 is shown in Figure 2. It shows that YODA captures and highlights perceptually essential areas. More visual results can be found in the supplementary materials. In addition to different neural networks, we test non-learnable methods to extract attention maps, also shown in Figure 2:
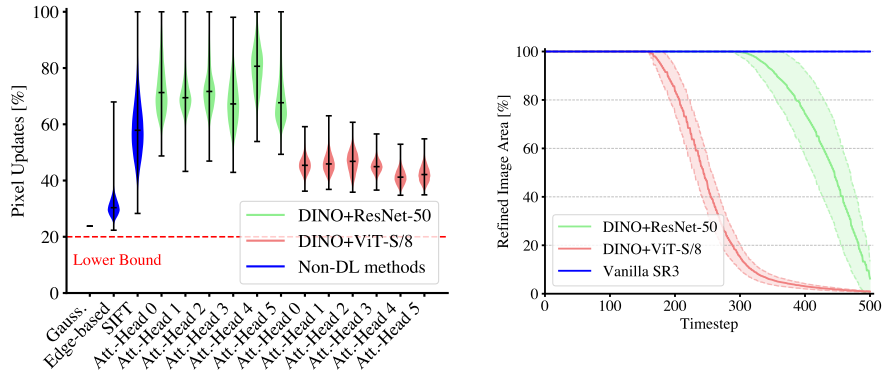
**Table 1:** Results of using different attention maps with SR3+YODA for $16 \times 16 \rightarrow$ $128 \times 128$ on CelebA-HQ. Aggregating the attention maps extracted with DINO and ResNet-50 backbone under the MAX strategy performs best.

| | Attention Maps for YODA | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|
| | PULSE (No Mask) | 16.88 | 0.440 | n.a. |
| | FSRGAN (No Mask) | 23.01 | 0.620 | n.a. |
| | SR3 Reported (No Mask) | 23.04 | 0.650 | n.a. |
| | SR3 Reproduced (No Mask) | 22.35 | 0.646 | 0.082 |
| | Non-learnable, fixed: Gaussian | 22.13 | 0.602 | 0.260 |
| | Non-learnable, adaptive: Edge-based | 22.93 | 0.648 | 0.151 |
| | Non-learnable, adaptive: SIFT | 22.84 | 0.678 | 0.095 |
| ViT-S/8 | Attention-Head 0 | 22.91 | 0.650 | 0.105 |
| | Attention-Head 1 | 22.43 | 0.616 | 0.130 |
| | Attention-Head 2 | 22.55 | 0.633 | 0.111 |
| | Attention-Head 3 | 22.73 | 0.641 | 0.110 |
| | Attention-Head 4 | 22.85 | 0.645 | 0.097 |
| | Attention-Head 5 | 22.86 | 0.648 | 0.101 |
| | Attention-AVG | 23.25 | 0.663 | 0.122 |
| | Attention-MAX | 23.46 | <u>0.683</u> | 0.103 |
| ResNet-50 | Attention-Head 0 | 22.82 | 0.649 | 0.115 |
| | Attention-Head 1 | 22.54 | 0.627 | 0.117 |
| | Attention-Head 2 | 22.84 | 0.650 | 0.107 |
| | Attention-Head 3 | 22.78 | 0.645 | 0.105 |
| | Attention-Head 4 | 22.38 | 0.620 | 0.127 |
| | Attention-Head 5 | 22.50 | 0.630 | 0.119 |
| | Attention-AVG | <u>23.55</u> | 0.682 | <u>0.093</u> |
| | **Attention-MAX** | **23.84** | **0.695** | **0.072** |

- **Gaussian:** Placing a simple 2D Gaussian pattern at the center of the image provides a straightforward approach, which relies on the assumption that the essential parts of an image are centered.
- **Edge-based:** Using the Canny edge detector, the attention maps are defined by the edges of the image, where adjacent and near edges are connected and blurred.
- **Scale-Invariant Feature Transform (SIFT):** Through Gaussian differences, SIFT [30] provides an attention map characterized by scale invariance. It produces an attention map by applying 2D Gaussian patterns around the points of interest.

Table 1 presents the results of our study with several baselines and masking variants for $16 \times 16 \rightarrow 128 \times 128$ scaling on the CelebA-HQ dataset. The straightforward Gaussian approach performs worst as it does not adapt to image features. The edge-based segmentation and SIFT methods showed improved performance over the reproduced baseline using a small batch size. However, they underperformed relative to the reported SR3 [38] results, which used a larger batch size.

**Fig. 3: (Left)** Ratio comparison between diffused pixels using our time-dependent masking approach and the total number of pixel updates in standard diffusion. On average, DINO with a ResNet-50 backbone leads to more pixel updates than the VIT-S/8 backbone. The lower bound, defined by $l$, is a threshold to eliminate areas that would never undergo diffusion. **(Right)** Refined image area in percentage across time steps for the MAX combination. Note that the sampling process goes from $T = 500$ to $T = 0$. ResNet-50 initiates the refinement process much earlier, advances more rapidly toward refining the entire image, and has a higher standard deviation.

In contrast, using DINO [6] to extract attention maps improves performance, regardless of the choice of backbone (ResNet-50 and ViT-S/8). We tested individual attention heads (0 to 5) independently, along with combination strategies that include averaging (AVG) and selecting the maximum value (MAX). The MAX combination achieved the best results compared to individual heads or the AVG combination. DINO [6] with a ResNet-50 backbone performs best, and we used the MAX aggregation strategy for all main experiments.

The left part of Figure 3 provides more information on the ratio of diffused pixels using our time-dependent masking. The upper bound is 100%, which is the case if diffusion is applied across all pixel locations throughout every time step (as in standard diffusion). Thus, any result under 100% shows that not all pixels are diffused during all time steps. As can be seen, DINO [6] coupled with ResNet-50 produces higher attention values, thereby leading to more total pixel updates. Also, the high variance of the ResNet-50 backbone indicates a high adaption capability. It can employ 100% of the updates for some samples, a characteristic not observed with the ViT-S/8 backbone. Nevertheless, the improved performance with the ViT-S/8 backbone compared to non-learnable methods and its low ratio of diffused pixels makes it an attractive candidate for future work regarding inference speed optimization based on sparser diffusion.

The right part of Figure 3 shows the ratio of diffused pixels depending on time steps with the MAX aggregation for the ResNet-50 and the ViT-S/8 backbone. As the backward diffusion goes from $T$ to 0, ResNet-50 backbone initiates and incorporates the refinement of the whole image areas more quickly than ViT-S/8.

**Table 2:** Results on face SR with 4× scaling ($64 \times 64 \rightarrow 256 \times 256$) and 8× scaling ($64 \times 64 \rightarrow 512 \times 512$) on CelebA-HQ. The models were trained for 1M steps on FFHQ and a reduced batch size of 4 and 8 instead of 256 to fit on a single A100-80GB GPU.
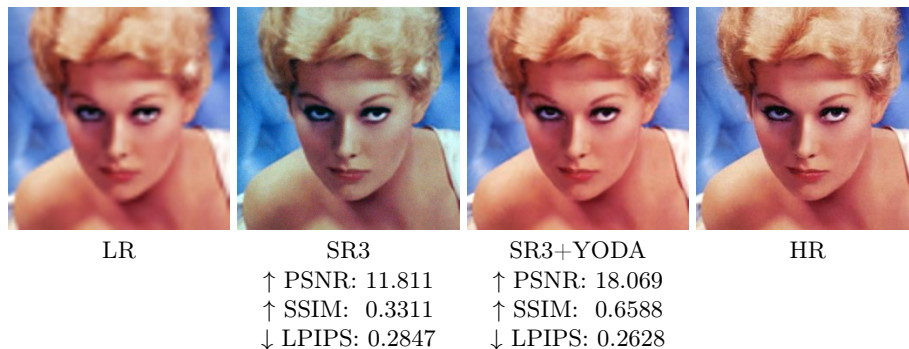
| Scaling | Model | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---------|-------|--------|--------|---------|
| 4× | SR3 [38] | 17.98 | 0.607 | 0.138 |
|  | **SR3 + YODA** | **26.33** | **0.838** | **0.090** |
| 8× | PartDiff ($K$=25) [48] | - | - | 0.222 |
|  | PartDiff ($K$=50) [48] | - | - | 0.217 |
|  | SR3 [38] | 17.44 | 0.631 | 0.147 |
|  | **SR3 + YODA** | **25.04** | **0.800** | **0.126** |

For the first 200 time steps, less than 20% of the image area is addressed by the attention map derived by ViT-S/8, whereas ResNet-50 has already developed to 100%. Therefore, this observation leads to our assumption that the superior performance of the ResNet-50 backbone may be attributed to its faster progression in refining the image. In other words, it performs better by progressing from 0% to 100% of the image area at an earlier stage. Intermediate diffusion results can be found in the supplementary materials.

### 4.2  Face Super-Resolution

**Details**  We use the FFHQ dataset [23] for training, which contains 50,000 high-quality facial images. We adopted the AdamW [29] optimizer, using a weight decay of 0.0001 and a learning rate of 5e-5. The number of sampling steps is set to $T_{\mathrm{train}} = 500$. We use the CelebA-HQ dataset [22] for evaluation, which contains 30,000 facial images. The number of sampling steps is set to $T_{\mathrm{eval}} = 200$. We trained all models for 1M iterations as in SR3 [38]. We evaluated three scenarios: $16 \times 16 \rightarrow 128 \times 128$, $64 \times 64 \rightarrow 256 \times 256$, and $64 \times 64 \rightarrow 512 \times 512$. Due to hardware limitations and missing quantitative results in the original publication of SR3, our experiments required a reduction from the originally used batch size of 256: we used a batch size of 4 for the $64 \times 64 \rightarrow 512 \times 512$, and 8 for the $64 \times 64 \rightarrow 256 \times 256$ scenario to fit on a single A100-80GB GPU.

**Results**  The results in Table 2 show a significant improvement when SR3 is coupled with YODA across all examined metrics. We explain the poor performance of SR3 with a phenomenon that is also observed by other authors [8, 42]. SR3 suffers from a color shift issue, which we attribute to the reduced batch size necessitated by our hardware limitations. Examples are shown in Figure 4 and Figure 5. This color shift manifests in a pronounced deviation in pixel-based metrics such as PSNR and SSIM but only slightly decreased perceptual quality in terms of LPIPS. Our results suggest that YODA's role extends beyond mere performance enhancement. It actively mitigates the color shift issue and stabilizes training, especially when faced with hardware constraints. With YODA,

| LR | SR3 | SR3+YODA | HR |
|----|-----|----------|----|
|    | ↑ PSNR: 11.811 | ↑ PSNR: 18.069 | |
|    | ↑ SSIM:  0.3311 | ↑ SSIM:  0.6588 | |
|    | ↓ LPIPS: 0.2847 | ↓ LPIPS: 0.2628 | |

**Fig. 4:** A comparison of LR, HR, SR3, and SR3+YODA images in the $64 \times 64 \rightarrow 256 \times 256$ setting (4x). SR3 suffers from color shift issues, as also observed by [8, 42]. Our method solves this issue and produces higher quality reconstructions.



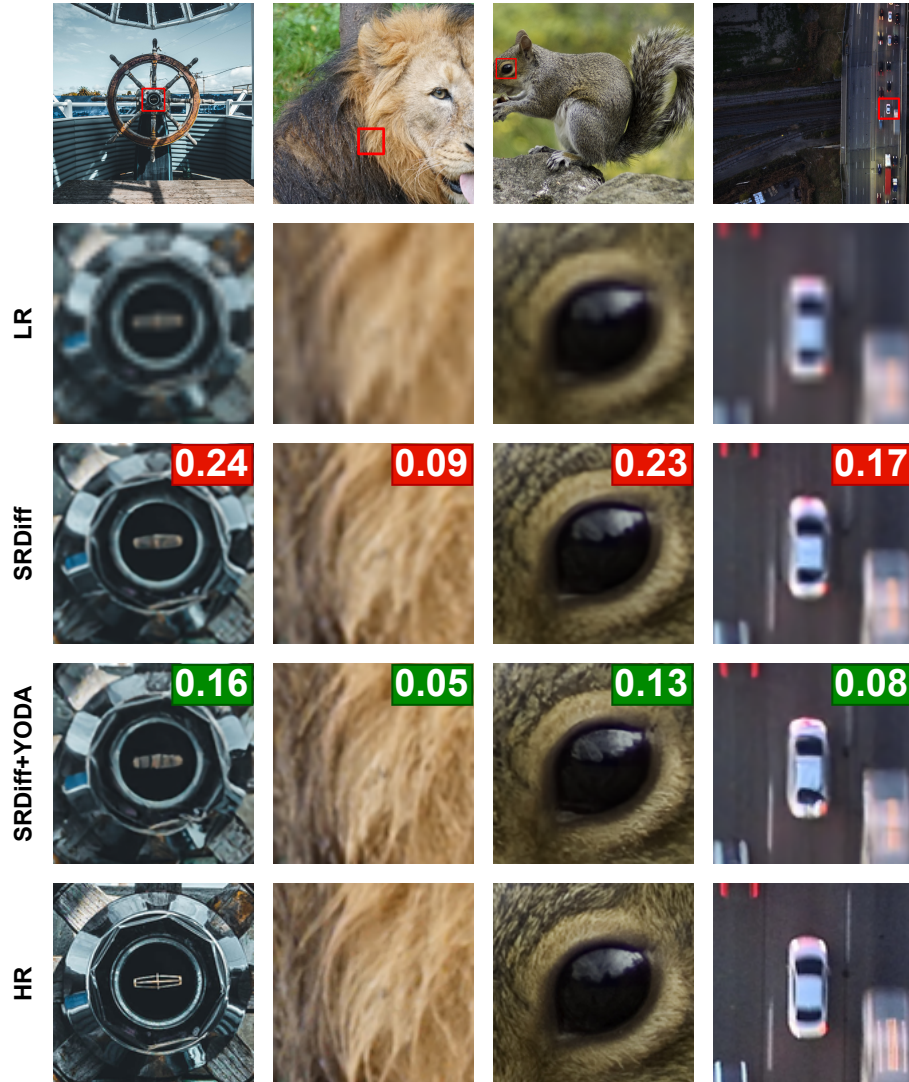| LR | SR3 | SR3+YODA | HR |
|----|-----|----------|----|
|    | ↑ PSNR: 23.061 | ↑ PSNR: 23.289 | |
|    | ↑ SSIM:  0.6208 | ↑ SSIM:  0.6334 | |
|    | ↓ LPIPS: 0.0504 | ↓ LPIPS: 0.0502 | |

**Fig. 5:** A comparison of LR, HR, SR3, and SR3+YODA images in the $16 \times 16 \rightarrow 128 \times 128$ setting (8x). The color shift in SR3 can still be observed (e.g., see upper right corners). YODA produces higher-quality images without color shift issues.

SR3 can be trained with a much smaller batch size and still achieves strong performance. Figure 5 offers another qualitative comparison between SR3 and SR3+YODA, highlighting subtle yet important differences, e.g., around the eyes, mouth, and hair.

### 4.3   General Super-Resolution

**Details** The experimental design follows SRDiff [25] and its hyperparameters, which are originally based on the experimental design of SRFlow [32]. For training, we employed 800 2K resolution HQ images from DIV2K [1] and 2,650 images from Flickr2K [41]. For testing, we used the DIV2K validation set of 100 images. Furthermore, we evaluated SR3, which was not originally tested on DIV2K. As in SRDiff, we extracted $40 \times 40$ sub-images with a batch size of 16, AdamW [29], a channel size of 64 with channel multipliers [1, 2, 2, 4] and $T = 100$.

**Fig. 6:** Example zoom-in regions of images (first row) from the DIV2K dataset. LPIPS is reported in the boxes (the lower, the better). YODA consistently produces better texture and more high-frequency details.

**Results** Table 3 shows the $4\times$ scaling general image SR results on the DIV2K validation set. The reported values include regression-based methods, which typically yield higher pixel-based scores than generative approaches [38]. This disparity is due to PSNR/SSIM penalizing misaligned high-frequency details, a known and significant challenge in the wider SR field [34].

**Table 3:** Quantitative results of 4× general SR on the DIV2K validation set.

| Type | Methods | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|
| Interpolation | Bicubic | 26.70 | 0.77 | 0.409 |
| Regression | EDSR [27] | 28.98 | 0.83 | 0.270 |
| | LIIF [7] | 29.24 | 0.84 | 0.239 |
| | RRDB [43] | 29.44 | 0.84 | 0.253 |
| GAN | RankSRGAN [47] | 26.55 | 0.75 | 0.128 |
| | ESRGAN [43] | 26.22 | 0.75 | 0.124 |
| Flow | SRFlow [32] | 27.09 | 0.76 | 0.120 |
| | HCFlow [26] | 27.02 | 0.76 | 0.124 |
| Flow + GAN | HCFlow++ [26] | 26.61 | 0.74 | 0.110 |
| VAE + AR | LAR-SR [18] | 27.03 | 0.77 | 0.114 |
| Diffusion | SR3 [38] | 14.14 | 0.15 | 0.753 |
| | **SR3 + YODA** | **27.24** | **0.77** | **0.127** |
| | SRDiff [25] | 27.41 | 0.79 | **0.136** |
| | **SRDiff + YODA** | **27.62** | **0.80** | 0.146 |

We observe that results from SR3 underperform compared to YODA. Equipping SR3 with YODA improves image quality even without an extensive hyperparameter search. We hypothesize that a vanilla SR3 strongly depends on larger batch sizes and longer diffusion times.

Combining SRDiff with YODA demonstrates improved performance in PSNR (+0.21db) and SSIM (+0.01), with a minor deterioration in LPIPS (+0.01). Thus, our approach excels in pixel-based metrics but sees a marginal decline in the perceptual metric. Nonetheless, when looking qualitatively at the predictions, one can observe that SRDiff benefits from YODA, as shown in Figure 6. YODA produces much better LPIPS values for perceptually essential areas, i.e., hair and cars, but falls short in background areas, i.e., blurry grass or dark ground. More results can be found in the supplementary materials.

Overall, YODA's strengths are more significant when combined with SR3 than with SRDiff. A critical distinction between SR3 and SRDiff is the handling of conditional information, i.e., the LR image, which we identify as a potential contributor to the reduced perceptual score LPIPS. SRDiff employs an LR encoder that generates an embedding during the denoising phase. Meanwhile, SR3 directly uses the LR image during the backward diffusion.

Furthermore, SR3 operates on the full LR images, whereas SRDiff diffuses residuals. Hence, the attention maps generated beforehand might not accurately capture the essential regions of the sparse residual inputs. Another possible reason for the lower perceptual score could be the image size during DINO training, i.e., $224 \times 224$ for the teacher network, which is much smaller than the 2K resolu-

tion images in DIV2K. As such, fine-tuning DINO on larger-scale images might be essential to capture more meaningful semantic features.

## 5    Conclusion

In this work, we presented "You Only Diffuse Areas" (YODA), an attention-guided diffusion-based image SR approach that emphasizes specific areas through time-dependent masking. YODA first extracts attention maps that reflect the pixel-wise importance of each scene using a self-supervised, general-purpose vision encoder. The attention maps are then used to guide the diffusion process by focusing on key regions in each time step while providing a fusion technique to ensure that masked and non-masked image regions are correctly connected between two successive time steps. This targeting allows for a more efficient transition to high-resolution outputs, prioritizing areas that gain the most from iterative refinements, such as detail-intensive regions. Beyond performance enhancement, YODA stabilizes training and mitigates the color shift issue that emerges when a reduced batch size constrains vanilla SR3. As a result, YODA consistently outperforms strong SR baselines such as SR3 and SRDiff, while requiring less computational resources by using smaller batch sizes.

## 6    Limitations & Future Work

A notable constraint of this study is its dependence on a good saliency estimation. Even though DINO is trained to be a generic vision encoder, it has known limitations that will reflect on the quality of YODA [11]. Meanwhile, the modularity of YODA allows for the saliency model to be switched once a better one becomes available. Additionally, DINO is explicitly trained on input image resolutions such as $224 \times 224$, which may not suffice for image SR applications with much larger spatial sizes of the LR image. An ideal solution would be a scale-invariant extraction of attention maps, e.g., a more extended version of our SIFT-adapted approach. Also, we use a linear correspondence between saliency and the number of steps. We posit that non-linear mappings may affect the quality of results. Lastly, YODA has to be extended for other recent state-of-the-art diffusion-based methods, e.g., latent-based methods like LDM [37] or unsupervised methods based on singular value decomposition like DDRM [24] or DDNM [45], which is orthogonal to our work (see supplementary materials for more details).

## Acknowledgment

# References

1. Agustsson, E., Timofte, R.: Ntire 2017 challenge on single image super-resolution: Dataset and study. In: CVPR Workshop (2017)
2. Anwar, S., Barnes, N.: Densely residual laplacian super-resolution. In: IEEE TPAMI (2020)
3. Baldassarre, F., El-Nouby, A., Jégou, H.: Variable rate allocation for vector-quantized autoencoders. In: ICASSP (2023)
4. Bansal, A., Borgnia, E., Chu, H.M., Li, J., Kazemi, H., Huang, F., Goldblum, M., Geiping, J., Goldstein, T.: Cold diffusion: Inverting arbitrary image transforms without noise. Adv. Neural Inform. Process. Syst. **36** (2024)
5. Bar-Tal, O., Yariv, L., Lipman, Y., Dekel, T.: Multidiffusion: Fusing diffusion paths for controlled image generation (2023)
6. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: ICCV (2021)
7. Chen, Y., Liu, S., Wang, X.: Learning continuous image representation with local implicit image function. In: CVPR (2021)
8. Choi, J., Lee, J., Shin, C., Kim, S., Kim, H., Yoon, S.: Perception prioritized training of diffusion models. 2022 ieee. In: CVPR (2022)
9. Chung, H., Lee, E.S., Ye, J.C.: Mr image denoising and super-resolution using regularized reverse diffusion. In: IEEE Transactions on Medical Imaging (2022)
10. Chung, H., Sim, B., Ye, J.C.: Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction. In: CVPR (2022)
11. Darcet, T., Oquab, M., Mairal, J., Bojanowski, P.: Vision transformers need registers. arXiv preprint arXiv:2309.16588 (2023)
12. Delbracio, M., Milanfar, P.: Inversion by direct iteration: An alternative to denoising diffusion for image restoration. arXiv preprint arXiv:2303.11435 (2023)
13. Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. In: IEEE TPAMI (2015)
14. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021)
15. Du, R., Chang, D., Hospedales, T., Song, Y.Z., Ma, Z.: Demofusion: Democratising high-resolution image generation with no $$$. arXiv preprint arXiv:2311.16973 (2023)
16. El-Shafai, W., Ali, A.M., El-Nabi, S.A., El-Rabaie, E.S.M., Abd El-Samie, F.E.: Single image super-resolution approaches in medical images based-deep learning: a survey. Multimedia Tools and Applications pp. 1–37 (2023)
17. Frolov, S., Hinz, T., Raue, F., Hees, J., Dengel, A.: Adversarial text-to-image synthesis: A review. Neural Networks **144**, 187–209 (2021)
18. Guo, B., Zhang, X., Wu, H., Wang, Y., Zhang, Y., Wang, Y.F.: Lar-sr: A local autoregressive model for image super-resolution. In: CVPR (2022)
19. Gur, S., Ali, A., Wolf, L.: Visualization of supervised and self-supervised neural networks via attribution guided factorization. In: AAAI. vol. 35, pp. 11545–11554 (2021)
20. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
21. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: Adv. Neural Inform. Process. Syst. (2020)

22. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. In: ICLR (2018)
23. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR (2019)
24. Kawar, B., Elad, M., Ermon, S., Song, J.: Denoising diffusion restoration models. Adv. Neural Inform. Process. Syst. **35**, 23593–23606 (2022)
25. Li, H., Yang, Y., Chang, M., Chen, S., Feng, H., Xu, Z., Li, Q., Chen, Y.: Srdiff: Single image super-resolution with diffusion probabilistic models. In: Neurocomputing (2022)
26. Liang, J., Lugmayr, A., Zhang, K., Danelljan, M., Van Gool, L., Timofte, R.: Hierarchical conditional flow: A unified framework for image super-resolution and image rescaling. In: ICCV (2021)
27. Lim, B., Son, S., Kim, H., Nah, S., Mu Lee, K.: Enhanced deep residual networks for single image super-resolution. In: CVPR Workshop (2017)
28. Liu, G.H., Vahdat, A., Huang, D.A., Theodorou, E.A., Nie, W., Anandkumar, A.: I $^2$ sb: Image-to-image schrödinger bridge. arXiv preprint arXiv:2302.05872 (2023)
29. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (2019)
30. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. In: IJCV (2004)
31. Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., Van Gool, L.: Repaint: Inpainting using denoising diffusion probabilistic models. In: CVPR (2022)
32. Lugmayr, A., Danelljan, M., Van Gool, L., Timofte, R.: Srflow: Learning the super-resolution space with normalizing flow. In: ECCV (2020)
33. Moser, B.B., Frolov, S., Raue, F., Palacio, S., Dengel, A.: Dwa: Differential wavelet amplifier for image super-resolution. In: Iliadis, L., Papaleonidas, A., Angelov, P., Jayne, C. (eds.) ICANN (2023)
34. Moser, B.B., Raue, F., Frolov, S., Palacio, S., Hees, J., Dengel, A.: Hitchhiker's guide to super-resolution: Introduction and recent advances. In: IEEE TPAMI (2023)
35. Moser, B.B., Shanbhag, A.S., Raue, F., Frolov, S., Palacio, S., Dengel, A.: Diffusion models, image super-resolution and everything: A survey. arXiv preprint arXiv:2401.00736 (2024)
36. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023)
37. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR (2022)
38. Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D.J., Norouzi, M.: Image super-resolution via iterative refinement. In: IEEE TPAMI (2022)
39. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: ICML (2015)
40. Sun, W., Chen, Z.: Learned image downscaling for upscaling using content adaptive resampler. In: IEEE Transactions on Image Processing (2020)
41. Timofte, R., Gu, S., Wu, J., Van Gool, L.: Ntire 2018 challenge on single image super-resolution: Methods and results. In: CVPR Workshop (2018)
42. Wang, J., Yue, Z., Zhou, S., Chan, K.C., Loy, C.C.: Exploiting diffusion prior for real-world image super-resolution. arXiv:2305.07015 (2023)
43. Wang, X., Yu, K., Wu, S., Gu, J., Liu, Y., Dong, C., Qiao, Y., Change Loy, C.: Esrgan: Enhanced super-resolution generative adversarial networks. In: ECCV Workshop (2018)

44. Wang, X., Yi, J., Guo, J., Song, Y., Lyu, J., Xu, J., Yan, W., Zhao, J., Cai, Q., Min, H.: A review of image super-resolution approaches based on deep learning and applications in remote sensing. Remote Sensing **14**(21), 5423 (2022)
45. Wang, Y., Yu, J., Zhang, J.: Zero-shot image restoration using denoising diffusion null-space model. arXiv preprint arXiv:2212.00490 (2022)
46. Whang, J., Delbracio, M., Talebi, H., Saharia, C., Dimakis, A.G., Milanfar, P.: Deblurring via stochastic refinement. In: CVPR (2022)
47. Zhang, W., Liu, Y., Dong, C., Qiao, Y.: Ranksrgan: Generative adversarial networks with ranker for image super-resolution. In: ICCV (2019)
48. Zhao, K., Hung, A.L.Y., Pang, K., Zheng, H., Sung, K.: Partdiff: Image super-resolution with partial diffusion models. arXiv:2307.11926 (2023)

## Appendix

**Complexity of YODA**

We discuss the resource implications of the core components of YODA: Identifying key regions, time-dependent masking, and the guided diffusion process. Additionally, we explore potential avenues for future enhancements aimed at optimizing computational efficiency.

**Identifying Key Regions.** To avoid the computational burden of on-the-fly generation, we pre-compute the attention maps prior to training. Table 4 shows the parameter count and throughput of different DINO backbones.

**Time-Dependent Masking and Guided Diffusion Process.** The integration of attention masks within the diffusion framework introduces minimal computational overhead, thanks to the inherently parallelizable nature of element-wise multiplication and addition, as demonstrated in Equation 9 and Equation 12. Consequently, the predominant factor influencing the overall computational complexity remains the choice of diffusion model, whether it be SR3, SRDiff, or another variant.
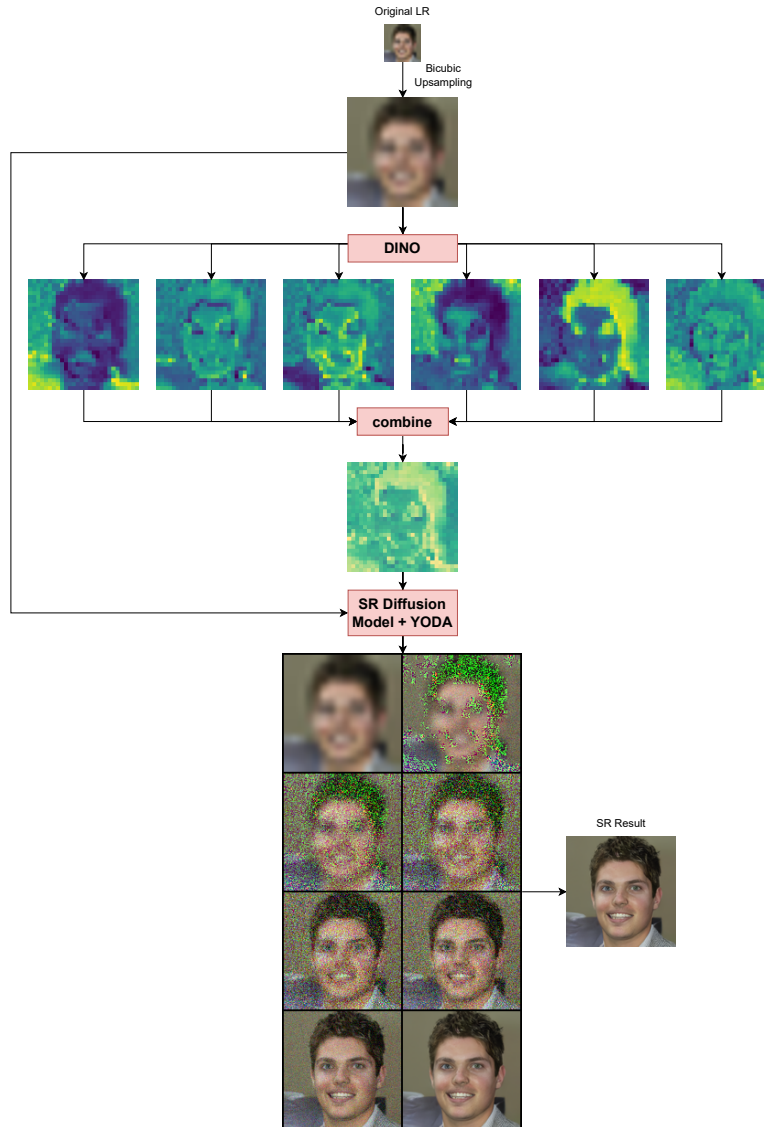
**Potential Future Improvements.** YODA notably decreases computational requirements by enabling the use of smaller batch sizes during training, which in turn reduces VRAM usage without compromising performance. Looking ahead, YODA paves the way for leveraging sparse diffusion techniques. Such approaches promise further computational savings by focusing computation efforts on selectively identified regions (through YODA), thereby streamlining the diffusion process. Currently, in PyTorch, applying masks to regions within a matrix does not result in computational savings.

**Table 4:** Details of different DINO backbones, values directly extracted from the original work [6]. Throughput was measured with a NVIDIA V100 GPU.

| Model | Parameters [M] | Throughput [img/s] |
|---|---|---|
| ResNet-50 | 23M | 1237 |
| ViT-S/8 | 21M | 180 |

This examination of YODA's complexity highlights its efficiency and the strategic decisions made to balance performance with computational demands. The potential for future improvements underscores YODA's adaptability and its capacity for integration with emerging sparse diffusion methodologies.

**Pipeline of YODA with DINO**



**Fig. 7:** An Overview of integrating YODA and DINO within SR diffusion models. Our process begins with using DINO to extract multiple attention maps. These maps are then combined to form a singular comprehensive attention map, denoted as **A**. Subsequently, leveraging **A**, YODA defines a unique diffusion schedule through time-dependent masks $\mathbf{M}(t)$ for every spatial location, as detailed within our method section.

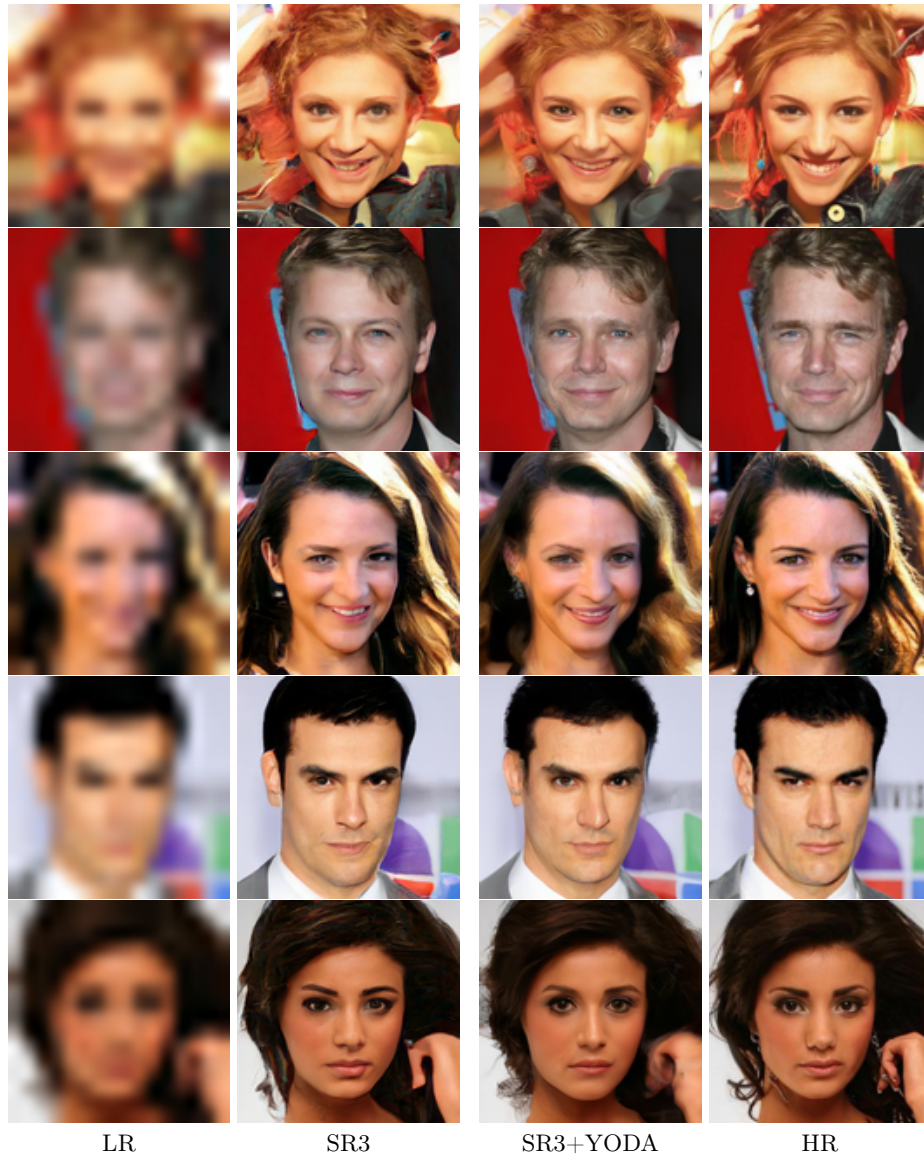**More Visual Results on Face SR ($16 \times 16 \rightarrow 128 \times 128$)**



| LR | SR3 | SR3+YODA | HR |

**Fig. 8:** A comparison of LR, HR, SR3, and SR3+YODA images for $16 \rightarrow 128$.

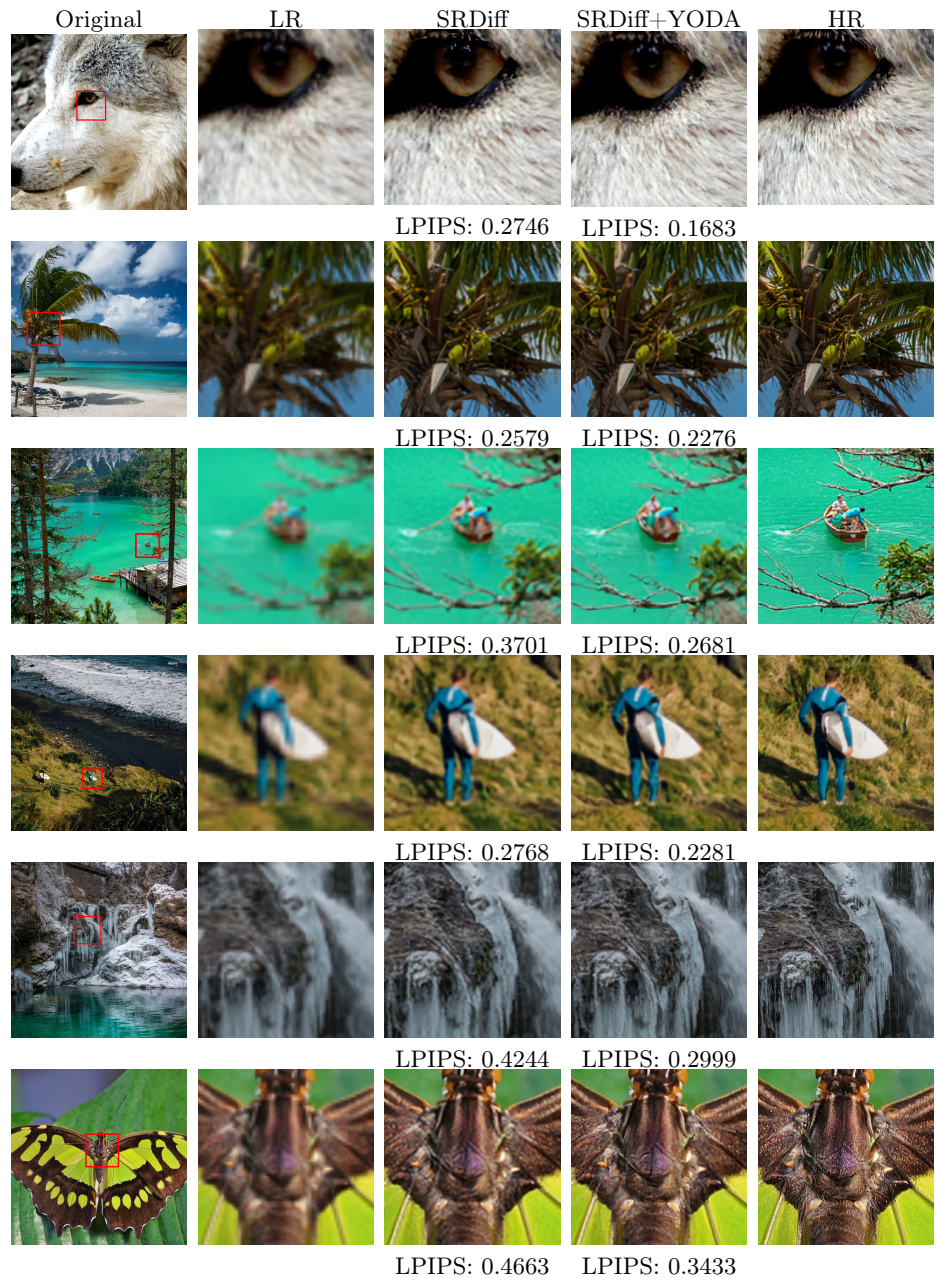## More Visual Results on DIV2K Validation



**Fig. 9:** Zoomed-in comparison of LR, HR, SRDiff, and SRDiff+YODA on DIV2K.

**Intermediate Results with YODA**



**Fig. 10:** Intermediate results of the YODA's guided diffusion process on CelebA-HQ.

**Other State-Of-The-Art Diffusion Models**

**Table 5:** ×4 upscaling results on ImageNet-Val. (256 × 256). Values directly derived from the original work of LDM [37].

| Method | IS ↑ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|
| Image Regression [38] | 121.1 | **27.9** | **0.801** |
| SR3 [38] | **180.1** | <u>26.4</u> | <u>0.762</u> |
| LDM-4 (100 steps) [37] | 166.3 | $24.4_{\pm 3.8}$ | $0.69_{\pm 0.14}$ |
| LDM-4 (big, 100 steps) [37] | <u>174.9</u> | $24.7_{\pm 4.1}$ | $0.71_{\pm 0.15}$ |
| LDM-4 (50 steps, guiding) [37] | 153.7 | $25.8_{\pm 3.7}$ | $0.74_{\pm 0.12}$ |

As shown in the study of Moser et al. [35], many approaches apply to image SR. In this section, we want to discuss their potential in combination with YODA and possible limitations for future work.

**Latent Diffusion Models.** Despite the significant advancements brought by Latent Diffusion Models (LDMs) [37], their efficacy in the realm of image SR competes closely with that of SR3 [38], as shown in Table 5. Unfortunately, recent research in this direction focused primarily on text-to-image tasks [17], which makes further comparisons with image SR methods challenging, e.g., SDXL [36], MultiDiffusion [5], or DemoFusion [15]. Nevertheless, their potential for image SR is undeniable. Concerning YODA, we also see great research avenues in combination with LDMs. A critical prerequisite for this synergy is the conversion of attention maps from pixel to latent representations. This aspect has to be investigated in more detail in future work.
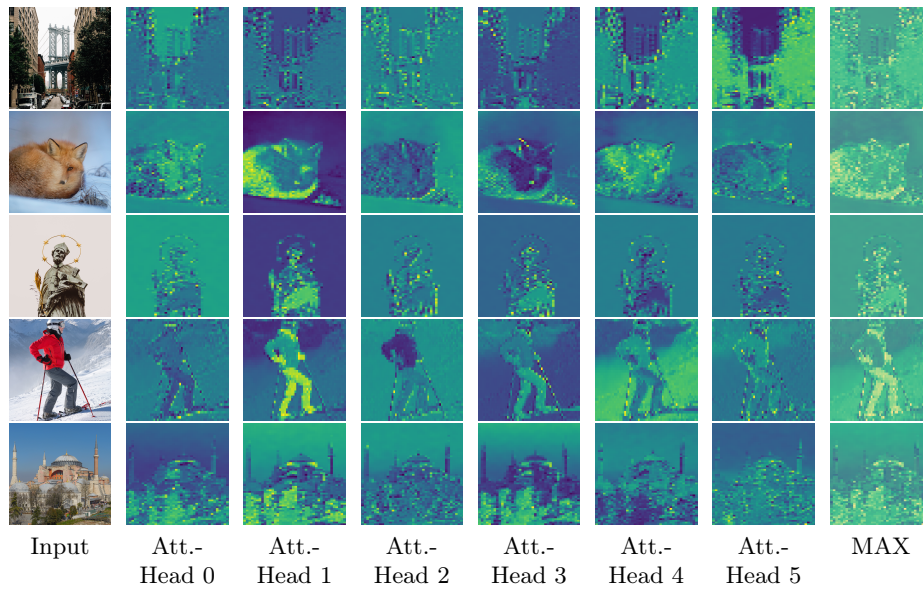
**Unsupervised Diffusion Models.** Another interesting research avenue is unsupervised diffusion models for image SR, exemplified by DDRM [24] or DDNM [45]. Interestingly, they use a pre-trained diffusion model to solve any linear inverse problem, including image SR, but they rely on singular value decomposition (SVD). Similar to the challenges faced with LDMs, integrating YODA into unsupervised diffusion models presents another interesting research avenue. The core of this challenge lies in devising a method for effectively translating attention maps into a format compatible with the SVD process used by these models. This transformation is crucial for harnessing the power of attention-based enhancements in unsupervised diffusion frameworks for image SR. Future work needs to conceptualize and implement a seamless integration strategy that combines the dynamic attention modulation offered by YODA with SVD.

**Alternative Corruption Spaces.** Applying YODA with alternative corruption spaces (not pure Gaussian noise), such as used in InDI [12], I$^2$SB [28], CCDF [10], or ColdDiffusion [4] is also an interesting future research direction. Although our primary focus has been refining and enhancing specific models' diffusion process through attention-guided masks, we acknowledge the orthogonal potential these approaches represent within the broader context of SR.

**More Details on DINO**

DINO [6] is a self-supervised learning approach. involving a teacher and student network. While both networks share the same architecture, their parameters differ. The student network is optimized to match the teacher's output via cross-entropy loss. During training, both receive two random views of the same input image: the teacher is trained on global views, i.e., $224 \times 224$ crops, while the student receives local views, i.e., $96 \times 96$ crops. This setup encourages the student to learn "local-to-global" correspondences. In other words, by predicting the teacher's output, the student learns to infer global information from local views. To prevent mode collapse due to identical architectures, the teacher's parameters are updated as a moving average of the student's parameters.

**More DINO Attention Maps for YODA**



**Fig. 11:** Comparison of attention maps derived from different DINO layers.