

# DemoFusion: Democratising High-Resolution Image Generation With No \$\$\$

Ruoyi Du<sup>1,4†</sup>, Dongliang Chang<sup>2\*</sup>, Timothy Hospedales<sup>3</sup>, Yi-Zhe Song<sup>4</sup>, Zhanyu Ma<sup>1</sup>

<sup>1</sup>PRIS, Beijing University of Posts and Telecommunications, China

<sup>2</sup>Tsinghua University, China <sup>3</sup>University of Edinburgh, UK <sup>4</sup>SketchX, University of Surrey, UK

{duruoyi, mazhanyu}@bupt.edu.cn, changdongliang@pris-cv.cn,

t.hospedales@ed.ac.uk, y.song@surrey.ac.uk

<https://ruoyidu.github.io/demofusion/demofusion.html>



Figure 1. Selected landscape samples of DemoFusion versus SDXL [24] (all images in the figure are presented at their actual sizes). SDXL can synthesize images up to a resolution of  $1024^2$ , while DemoFusion extends SDXL to generate images at  $4\times$ ,  $16\times$ , and even higher resolutions without any fine-tuning or prohibitive memory demands. All generated images are produced using a single RTX 3090 GPU. Best viewed ZOOMED-IN.

## Abstract

High-resolution image generation with Generative Artificial Intelligence (GenAI) has immense potential but, due to the enormous capital investment required for training, it is increasingly centralised to a few large corporations, and hidden behind paywalls. This paper aims to democratise high-resolution GenAI by advancing the frontier of high-resolution generation while remaining accessible to a broad audience. We demonstrate that existing Latent Diffusion

Models (LDMs) possess untapped potential for higher-resolution image generation. Our novel DemoFusion framework seamlessly extends open-source GenAI models, employing Progressive Upscaling, Skip Residual, and Dilated Sampling mechanisms to achieve higher-resolution image generation. The progressive nature of DemoFusion requires more passes, but the intermediate results can serve as “previews”, facilitating rapid prompt iteration.

## 1. Introduction

Generating high-resolution images with Generative Artificial Intelligence (GenAI) models has demonstrated remark-

<sup>†</sup>The work is done while Ruoyi Du visiting the People-Centred AI Institute at the University of Surrey

\*Corresponding Author

able potential [1, 19, 22]. However, these capabilities are increasingly centralised. Training high-resolution image generation models requires substantial capital investments in hardware, data, and energy that are beyond the reach of individual enthusiasts and academic institutions. For example, training Stable Diffusion 1.5, at a resolution of  $512^2$ , entails over 20 days of training on 256 A100 GPUs. Companies that make these investments understandably want to recoup their costs and increasingly hide the resulting models behind paywalls. This trend toward centralisation and pay-per-use access is accelerating as GenAI image synthesis advances in quality since the investment required to train image generators increases rapidly with image resolution.

In this paper we reverse this trend and re-democratise GenAI image synthesis by introducing *DemoFusion*, which pushes the frontier of high-resolution image synthesis from  $1024^2$  in SDXL [24], Midjourney [19], DALL-E [22], etc to  $4096^2$  or more. DemoFusion requires no additional training and runs on a single consumer-grade RTX 3090 GPU (hardware for the “working class” in the GenAI era), as shown in Fig. 1. The only trade-off? A little more patience.

Specifically, we start with the open source SDXL [24] model, capable of generating images of  $1024^2$ . DemoFusion is a plug-and-play extension to SDXL that enables  $4\times$ ,  $16\times$ , or more increase in generation resolution (Fig 1) – all with zero additional training, and only a few simple lines of code. Off-the-shelf SDXL fails if directly prompted to produce higher-resolution images (Fig. 2 (a)). However, we observe that text-to-image LDMs encounter many cropped photos during their training process. These cropped photos either exist inherently in the training set or are intentionally cropped for data augmentation. Consequently, models like SDXL occasionally produce outputs that focus on localised portions of objects [24], as illustrated in Fig. 2 (b). In other words, existing open-source LDMs already contain sufficient prior knowledge to generate high-resolution images, if only we can unlock them by fusing multiple such high-resolution patches into a complete scene.

However, achieving coherent patch-wise high-resolution generation is non-trivial. A recent study, MultiDiffusion [2] showcased the potential of fusing multiple overlapped denoising paths to generate panoramic images. Yet, when directly applying this approach to generate specific high-resolution object-centric images, results are repetitive and distorted without global semantic coherence [41], as illustrated in Fig. 2 (c). We conjecture the underlying reason is that overlapped patch denoising merely reduces the seam issue without a broad perception of the global context required for semantic coherence. DemoFusion builds upon the same idea of fusing multiple denoising paths from a pre-trained SDXL model to achieve high-resolution gener-

---

Information from: <https://huggingface.co/runwayml/stable-diffusion-v1-5>



**Figure 2. Examples of  $4\times (2048^2)$  generation based on SDXL [24].** (a) Directly prompting SDXL to generate a  $4\times$  image. (b) SDXL [24] inferences on non-overlapping patches at the original resolution. It fails, but reveals that the SDXL possesses prior knowledge of localized patches at higher resolutions. (c) MultiDiffusion [2] fuses multiple overlapping denoising paths to generate higher-resolution images without edge effects, but lacks the global context for semantic coherence. (d) Our proposed DemoFusion achieves global semantic coherence in high-resolution generation.

ation. It introduces three key mechanisms to achieve global semantic coherence together with rich local detail (Fig. 2 (d) vs (a, c)): (i) *Progressive Upscaling*: Starting with the low-resolution input, DemoFusion iteratively enhances images through an “upsample-diffuse-denoise” loop, using the noise-inversed lower-resolution image as a better initialisation for generating the higher-resolution image. (ii) *Skip Residual*: Within the same iteration, we additionally utilise the intermediate noise-inversed representations as skip residuals, maintaining global consistency between high and low-resolution images. (iii) *Dilated Sampling*: We extend MultiDiffusion to increase global semantic coherence by using dilated sampling of denoising paths. These three techniques to modify inference are simple to implement on a pre-trained SDXL and provide a dramatic boost in high-resolution image generation quality and coherence. Fig. 3 illustrates the framework.

The caveat is that generating high-resolution images does require more runtime (users need to exercise more patience). This is partially due to the progressive upscaling requiring more passes; however, primarily because the time required grows exponentially with resolution (as per

any patch-wise LDM [2]), and thus, the highest resolution pass dominates the cost. Nevertheless, the memory cost is low enough for consumer-grade GPUs, and progressive generation allows the users to preview low-resolution results rapidly, facilitating rapid iteration on the prompt until satisfaction with the general layout and style, prior to waiting for a full high-resolution generation.

## 2. Related Work

With the progress of several years, diffusion model (DM) [32] has recently reached its own “tipping point” – with the emergence of works like DDPM [33], DDIM [33], ADM [5], DM has shown great potential in image generation due to its outstanding generation quality and diversity. Subsequently, using a pre-trained autoencoder, the latent diffusion model (LDM) [28] applies a diffusion model in the latent space, achieving efficient training and inference. This enabled the emergence of high-performance generative models trained on billions of data, such as the Stable Diffusion series. LDM’s excellent generalisation capability has led to subsequent research on controllable generation [21, 29, 40] and editable generation [3, 9, 20]; it has also been widely applied in numerous downstream generative tasks, such as text-to-video [8, 11, 36], text-to-3D [18, 25, 37], text-to-avatar [6, 17, 35], and text-to-human sketch [14, 26], *etc.*

Despite achieving numerous successes, current LDMs like Stable Diffusion 1.5 and Stable Diffusion XL are still confined to generating images at resolutions of  $512^2$  and  $1024^2$ , respectively [24]. Escalating resolution significantly increases training expenses and computational load, making such models impractical for most researchers and users. An intuitive solution to generate high-resolution images involves using LDMs for initial image generation, followed by enhancement through a super-resolution (SR) model. Cascaded Diffusion Models [12] cascades several diffusion-based SR models behind a diffusion model, but its application remains capped at  $256^2$  resolution images. We attempted to enhance state-of-the-art LDMs with SR models [34, 39], but found that images generated at lower resolutions were deficient in detail. Upscaling these images with SR failed to yield the high-resolution detail desired. Another attempt is to retrain/fine-tune open-source DMs to achieve satisfactory results [13, 41], but fine-tuning still brings a non-negligible cost.

Recently, MultiDiffusion [2] fuses multiple overlapped denoising paths of LDMs, achieving seamless panorama generation in a training-free manner. Subsequently, SyncDiffusion [16] further constrains the consistency between denoising paths using a gradient descent approach. However, these methods are limited to generating scene images through repetition; when applied to generating specific objects, they lead to local repetition and structural distortion.

Valuing the training-free characteristic of such methods, we proposed DemoFusion based on MultiDiffusion in this paper towards democratising high-resolution generation.

Note that a recent concurrent work, SCALECRAFTER [7], with the same motivation, proposed a tuning-free framework for high-resolution image generation. It ingeniously adapts the diffusion model for higher resolutions by dilating its convolution kernels at specific layers. Despite a smart move, our experiments indicate that SCALECRAFTER somewhat degrades the model’s performance and does not bring about the local details expected at higher resolutions. In contrast, DemoFusion has demonstrated better results.

## 3. Methodology

### 3.1. Preliminaries

**Latent Diffusion Model:** Given an image  $\mathbf{x}$ , an LDM first encodes it to the latent space via the encoder of the pre-trained autoencoder, *i.e.*,  $\mathbf{z} = \mathcal{E}(\mathbf{x})$ ,  $\mathbf{z} \in \mathbb{R}^{c \times h \times w}$ .

Following this, the two core components of the diffusion model, the diffusion and the denoising process, take place in the latent space. The diffusion process comprises a sequence of  $T$  steps with Gaussian noise incrementally introducing into the latent distribution at each step  $t \in [0, T]$ . With a prescribed variance schedule  $\beta_1, \dots, \beta_T$ , the diffusion process can be formulated as

$$q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I}). \quad (1)$$

In contrast, the denoising process aims to recover the cleaner version  $\mathbf{z}_{t-1}$  from  $\mathbf{z}_t$  by estimating the noise, which can be expressed as

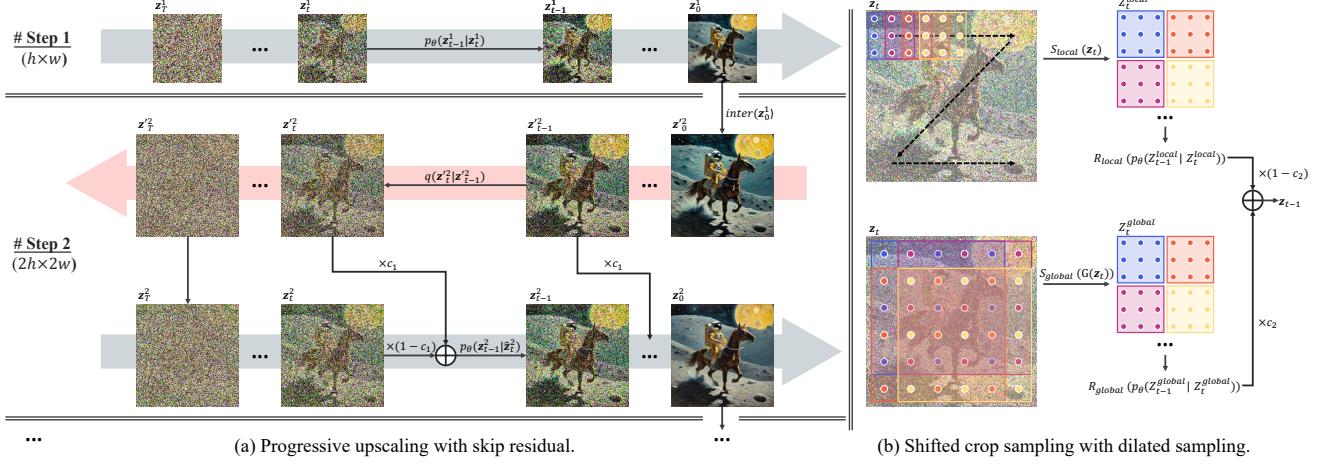
$$p_\theta(\mathbf{z}_{t-i} | \mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t-i}; \mu_\theta(\mathbf{z}_t, t), \Sigma_\theta(\mathbf{z}_t, t)), \quad (2)$$

where  $\mu_\theta$  and  $\Sigma_\theta$  are determined through estimation procedures and  $\theta$  denotes the parameters of the denoise model.

**MultiDiffusion:** MultiDiffusion [2] extends LDMs such as SDXL to produce high-resolution panoramas by overlapped patch-based denoising.

In simple terms, MultiDiffusion defines a latent space  $\mathbb{R}^{c \times H \times W}$  with  $H > h$  and  $W > w$ . For arbitrary denoising step  $t$  with  $\mathbf{z}_t \in \mathbb{R}^{c \times H \times W}$ , MultiDiffusion first applies a shifted crop sampling  $\mathcal{S}_{local}(\cdot)$  to obtain a series of local latent representations, *i.e.*,  $Z_t^{local} = [\mathbf{z}_{0,t}, \dots, \mathbf{z}_{n,t}, \dots, \mathbf{z}_{N,t}] = \mathcal{S}_{local}(\mathbf{z}_t)$ ,  $\mathbf{z}_{n,t} \in \mathbb{R}^{c \times h \times w}$ , where  $N = (\frac{(H-h)}{d_h} + 1) \times (\frac{(W-w)}{d_w} + 1)$ ,  $d_h$  and  $d_w$  is the vertical and horizontal stride, respectively.

After that, the conventional denoising process is independently applied to these local latent representations via  $p_\theta(\mathbf{z}_{n,t-1} | \mathbf{z}_{n,t})$ . And then  $Z_{t-1}^{local}$  is reconstructed to the original size with the overlapped parts averaged as  $\mathbf{z}_{t-1} = \mathcal{R}_{local}(Z_{t-1}^{local})$ , where  $\mathcal{R}_{local}$  denotes the reconstruction



**Figure 3. The proposed DemoFusion framework.** (a) Starting with conventional resolution generation, DemoFusion engages an “upsample-diffuse-denoise” loop, taking the low-resolution generated results as the initialization for the higher resolution through noise inversion. Within the “upsample-diffuse-denoise” loop, a noise-inverted representation from the corresponding time-step in the preceding diffusion process serves as skip-residual as global guidance. (b) To improve the local denoising paths of MultiDiffusion, we introduce dilated sampling to establish global denoising paths, promoting more globally coherent content generation.

process. Eventually, a higher-resolution panoramic image can be obtained by directly decoding  $\mathbf{z}_0$  into image  $\hat{\mathbf{x}}$ .

MultiDiffusion provides effective panorama generation, thanks to smoothing the edge effects between generated patches. However, as discussed by [41], and illustrated in Fig. 2, it struggles with generating coherent semantic content for specific objects. The fundamental reason for this is that each patch/diffusion path is constrained only by the text condition and lacks awareness of the global context of the other patches.

We introduce three modifications to the inference procedure of SDXL that enable a patch-wise high-resolution image generation strategy to achieve both global semantic coherence and rich local details. These are: *Progressive Upscaling* (see Sec. 3.2), *Skip Residual* (see Sec. 3.3) and *Dilated Sampling* (see Sec. 3.4). The overall flow of DemoFusion is summarised in Appendix A.

### 3.2. Progressive Upscaling

Progressively generating images from low to high resolution is a well-established concept [15]. By initially synthesizing a semantically coherent overall structure at low resolution, and subsequently increasing resolution to add detailed local features, models can produce coherent yet rich images. In this paper, we present a novel *progressive upscaling* generation process tailored for LDMs (Fig 3 (a)).

Consider a pre-trained latent diffusion model with parameters  $\theta$ , operating on the latent space  $\mathbb{R}^{c \times h \times w}$  to produce images with a resolution magnified by a factor of  $K$ . The scaling factor for the side length should be  $S = \sqrt{K}$ . And the target latent space is  $\mathbb{R}^{c \times H \times W}$  where  $H = Sh$  and  $W = Sw$ . Instead of directly synthesizing  $\mathbf{z}_t \in \mathbb{R}^{c \times H \times W}$ ,

we break the generation process into  $S$  distinct phases, each consisting of an “upsample-diffuse-denoise” loop, except for the first phase which follows an “initialise-denoise” scheme. Specifically, given diffusion and denoising process as  $q(\mathbf{z}_T | \mathbf{z}_0) = \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})$  and  $p_\theta(\mathbf{z}_0 | \mathbf{z}_T) = \prod_{t=T}^1 p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t)$ . Then, we can formulate the proposed progressive upscaling generation process as

$$p_\theta(\mathbf{z}_0^S | \mathbf{z}_T^1) = p_\theta(\mathbf{z}_0^1 | \mathbf{z}_T^1) \prod_{s=2}^S (q(\mathbf{z}'_T^s | \mathbf{z}'_0^s) p_\theta(\mathbf{z}_0^s | \mathbf{z}'_T^s)), \quad (3)$$

where  $\mathbf{z}'_0^s$  is obtained through explicit upsampling as  $\mathbf{z}'_0^s = inter(\mathbf{z}_0^{s-1})$  and  $inter(\cdot)$  is an arbitrary interpolation algorithm (e.g., bicubic). In essence, we first run a regular LDM such as SDXL as  $p_\theta(\mathbf{z}_0^1 | \mathbf{z}_T^1)$ . We then iteratively for each scale  $s$ : (i) upscale the low-resolution image  $\mathbf{z}_0^{s-1}$  to  $\mathbf{z}'_0^s$ , (ii) reintroduce noise via the diffusion process to obtain  $\mathbf{z}'_T^s$ , and (iii) denoise to obtain  $\mathbf{z}_0^s$ . By repeating this process, we can compensate for the artificial interpolation-based upsampling and gradually fill in more and more local details.

### 3.3. Skip Residual

The “diffuse-denoise” process has parallels in some image editing works – people attempt to find the initial noise of an image using specialized noise inversion techniques, ensuring that the unedited parts remain consistent with the original image during the denoising editing process [9, 20]. However, these inversion techniques are less practical to DemoFusion’s denoising process. Therefore, we instead simply use a conventional diffusion process by adding random Gaussian noise.

However, directly diffusing  $\mathbf{z}_0^s$  to  $\mathbf{z}'_T^s$  as initialization would result in most information loss. In contrast, diffusing to an intermediate  $t$  and then starting denoise from  $\mathbf{z}'_t^s$  might be better. However, it is challenging to determine the optimal intersection time-step  $t$  of the “upsample-diffuse-denoise” loop – the larger the  $t$ , the more information is lost, which weakens the global perception; the smaller the  $t$ , the stronger the noise introduced by upsampling (refer to Appendix C). It is a difficult trade-off and could be example-specific. Therefore, we introduce the skip residual as a general solution, which can be informally considered as a weighted fusion of multiple “upsample-diffuse-denoise” loops with a series of different intersection time-steps  $t$  (Fig. 3 (a)).

For each generation phase  $s$ , we have already obtained a series of noise-inversed versions of  $\mathbf{z}_0^s$  as  $\mathbf{z}'_t^s$  with  $t \in [1, T]$ . During the denoising process, we introduce the corresponding noise-inversed versions as *skip residuals*. In other words, we modify  $p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)$  to  $p_\theta(\mathbf{z}_{t-1}|\hat{\mathbf{z}}_t)$  with

$$\hat{\mathbf{z}}_t^s = c_1 \times \mathbf{z}'_t^s + (1 - c_1) \times \mathbf{z}_t^s, \quad (4)$$

where  $c_1 = ((1 + \cos(\frac{T-t}{T} \times \pi))/2)^{\alpha_1}$  is a scaled cosine decay factor with a scaling factor  $\alpha_1$ . This essentially utilizes the results from the previous phase to guide the generated image’s global structure during the initial steps of the denoising process. Meanwhile, we gradually reduce the impact of the noise residual, allowing the local denoising paths to optimize the finer details more effectively in the later steps.

### 3.4. Dilated Sampling

Beyond the explicit integration of global information as a residual, we introduce *dilated sampling* to give each denoising path more global context. The technique of dilating convolutional kernels to expand their receptive field is conventional in various dense prediction tasks [38]. The concurrent tuning-free method, SCALECRAFTER [7], similarly uses dilated convolutional kernels for adapting trained latent diffusion models to higher-resolution image generation. However, our approach diverges here: rather than dilating the convolutional kernel, we directly dilate the sampling within the latent representation. After that, the global denoising paths, derived through dilated sampling, are processed analogously to local denoising paths in MultiDiffusion.

As depicted in Fig. 3 (b), we applied shifted dilated sampling to obtain a series of global latent representation, *i.e.*,  $Z_t^{global} = [\mathbf{z}_{0,t}, \dots, \mathbf{z}_{m,t}, \dots, \mathbf{z}_{M,t}] = \mathcal{S}_{global}(\mathbf{z}_t)$ ,  $\mathbf{z}_{m,t} \in \mathbb{R}^{c \times h \times w}$ . To sample from the whole latent representation, the dilation factor is set to be  $s$  and  $M = s^2$ . Similarly, we apply the general denosing process on these global latent representations as  $p_\theta(\mathbf{z}_{m,t-1}|\mathbf{z}_{m,t})$ . Then, the reconstructed global representations are fused with the reconstructed local representations to form the final latent rep-

resentation:

$$\mathbf{z}_{t-1} = c_2 \times \mathcal{R}_{global}(Z_{t-1}^{global}) + (1 - c_2) \times \mathcal{R}_{local}(Z_{t-1}^{local}), \quad (5)$$

where  $c_2 = ((1 + \cos(\frac{T-t}{T} \times \pi))/2)^{\alpha_2}$  is a scaled cosine decay factor with a scaling factor  $\alpha_2$ , also chosen based on the characteristic of the diffusion model where earlier steps mainly reconstruct the overall structure, while later steps focus on refining the details.

It is noteworthy that directly using dilated sampling can lead to grainy images. This is because, unlike the local denoising paths, which have overlaps, the global denoising paths operate independently of each other. To address this issue, we employ a straightforward yet intuitive approach – applying a Gaussian filter  $\mathcal{G}(\cdot)$  to the latent representation before performing dilated sampling as  $Z_t^{global} = \mathcal{S}_{global}(\mathcal{G}(\mathbf{z}_t))$ . The kernel size of the Gaussian filter is set to be  $4s - 3$ , making it sufficient at every phase. Moreover, the standard deviation of the Gaussian filter will decrease from  $\sigma_1$  to  $\sigma_2$  as  $c_3 \times (\sigma_1 - \sigma_2) + \sigma_2$ , where  $c_3 = ((1 + \cos(\frac{T-t}{T} \times \pi))/2)^{\alpha_3}$  is also a scaled cosine decay factor with a scaling factor  $\alpha_3$ , ensuring that the effect of the filter gradually diminishes as the directions of global denoising paths become consistent, preventing the final image from becoming blurry.

## 4. Experiments

Here, we report qualitative and quantitative experiments and ablation studies. For more details and results, please refer to Appendix: implementation details in Appendix B, more discussions in Appendix C, more visualisations in Appendix D, more applications in Appendix E, and all prompts we use in Appendix F.

### 4.1. Comparison

We compared DemoFusion with the following methods (i) **SDXL** [24], which is designed to generate images of  $1024^2$ . In the quantitative experiments, we also report the results of inferencing it at higher resolutions. (ii) **MultiDiffusion** [2], our baseline method based on overlapped local patch denoising. (iii) **SDXL+BSRGAN**. Using a super-resolution model is an intuitive solution to directly upscale SDXL results. Here, we choose BSRGAN [39], a representative SR method, for comparison. (iv) **SCALECRAFTER** [7], a concurrent training-free high-resolution generation method built on SDXL, which upscales by dilating convolutional kernels at specific layers.

**Qualitative Results:** As shown in Fig. 4, each model is asked to generate images at  $4\times$  and  $16\times$  resolutions (compared to SDXL). We chose three prompts about realistic content rather than showcasing DemoFusion’s prowess in artistic creation, as such content is more objective and facilitates a fair comparison.

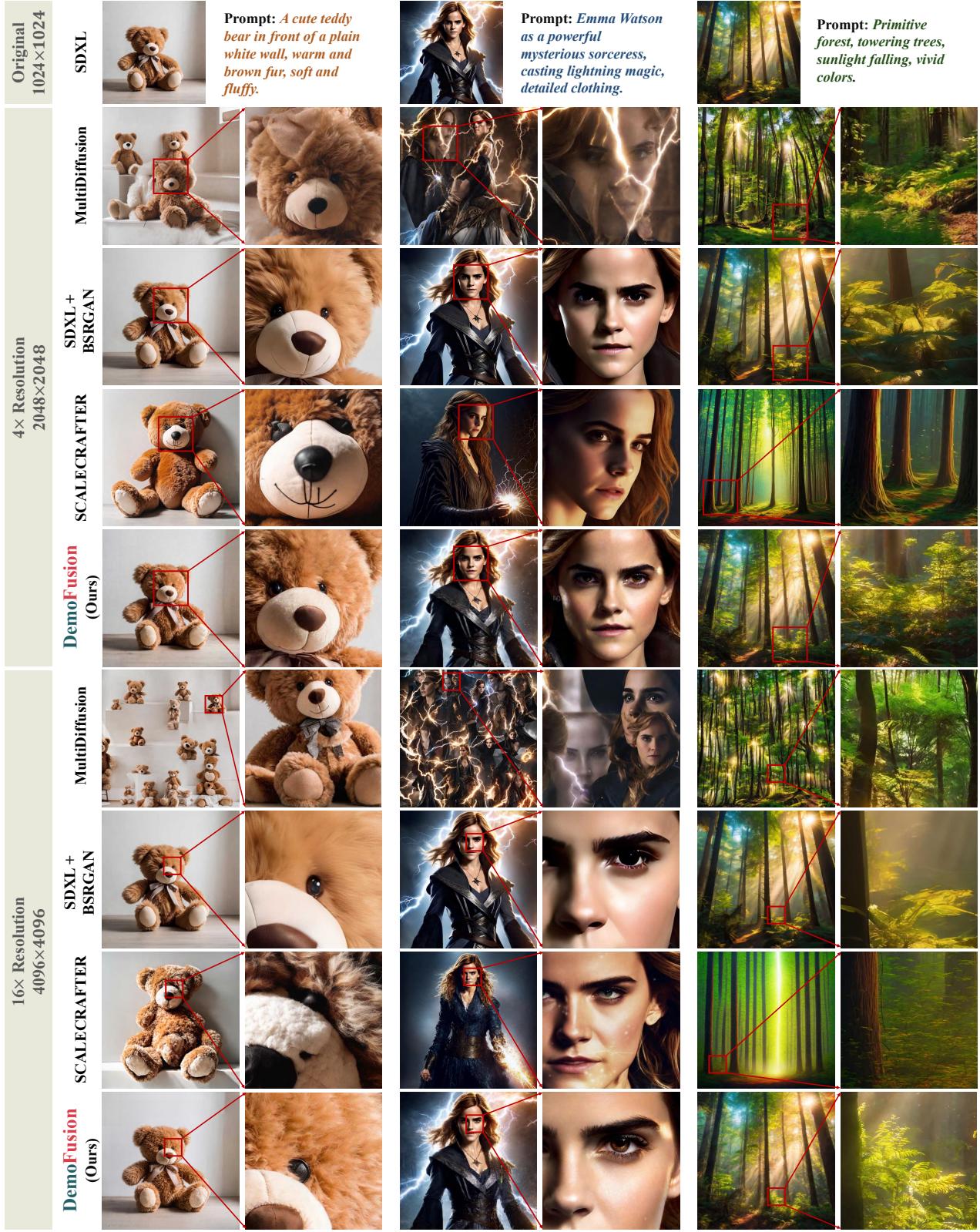


Figure 4. **Qualitative comparison with other baselines.** Local details have already been zoomed in, but it's still recommended to **ZOOM IN** for a closer look.

Method	2048 × 2048						2048 × 4096						4096 × 4096					
	FID ↓	IS ↑	FID <sub>crop</sub> ↓	IS <sub>crop</sub> ↑	CLIP ↑	Time	FID ↓	IS ↑	FID <sub>crop</sub> ↓	IS <sub>crop</sub> ↑	CLIP ↑	Time	FID ↓	IS ↑	FID <sub>crop</sub> ↓	IS <sub>crop</sub> ↑	CLIP ↑	Time
SDXL Direct Inference [24]	79.66	13.47	73.91	17.38	28.12	1 min	97.08	14.12	96.41	18.01	27.29	3 min	105.65	14.01	98.59	19.47	25.64	8 min
MultiDiffusion [2]	75.93	14.56	70.93	17.85	28.97	3 min	89.38	14.17	82.78	18.87	28.66	6 min	97.98	13.84	79.45	19.73	28.62	15 min
SDXL + BSRGAN [39]	<u>66.41</u>	<u>16.22</u>	<u>67.42</u>	<u>21.11</u>	<u>29.61</u>	1 min	<b>68.70</b>	<u>16.29</u>	<u>75.03</u>	<u>21.76</u>	<u>29.01</u>	1 min	<b>66.44</b>	<b>16.21</b>	<u>77.20</u>	<u>22.42</u>	<b>29.63</b>	1 min
SCALECRAFTER [7]	69.91	15.72	68.36	19.44	29.51	1 min	80.16	15.29	83.08	19.56	28.87	6 min	87.50	15.20	84.36	20.32	29.04	19 min
DemoFusion (Ours)	<b>65.73</b>	16.41	<b>64.81</b>	<b>21.40</b>	<b>29.68</b>	3 min	<u>73.15</u>	<b>16.37</b>	<b>71.35</b>	<b>23.55</b>	<b>29.05</b>	11 min	<b>74.11</b>	<u>16.11</u>	<b>70.34</b>	<b>24.28</b>	<u>29.57</u>	25 min

Table 1. Quantitative comparison results. The best results are marked in **bold**, and the second best results are marked by underline.



Figure 5. Ablation studies on the three components of DemoFusion: Progressive Upscaling (PU), Skip Residual (SR), and Dilated Upsampling (DS). All images are generated at  $3072^2$  ( $9 \times$  resolutions). Best viewed ZOOMED-IN.

Firstly, as previously mentioned, MultiDiffusion tends to generate repetitive content lacking semantic coherence. For SDXL+BSRGAN, we observe that the SR model effectively eliminates the blurriness and jagged edges of up-sampling, resulting in sharp and pleasing outcomes. However, the goal of the SR model is to produce images consistent with the input, which limits its performance in high-resolution generation – needing more detail for true high-resolution visuals beyond simple smoothing. Checking the zoomed-in results of  $4096^2$  – compared to SDXL+BSRGAN, DemoFusion generates much richer details in the fur of the teddy bear, gives much richer details to Hermoine’s eyes, and adds much more detail to the forest vegetation. This comparison confirms that high-resolution generation cannot be substituted by simple image super-resolution. As for SCALECRAFTER, while it partially addresses the issue of MultiDiffusion’s repetitive content, it still needs improvement in semantic coherence. *E.g.*, the teddy bear has multiple arms, eyes, or mouths. Additionally, directly dilating the convolutional kernels has somewhat affected the performance of the LDM, resulting in an overall image quality degradation, and local details exhibit many repetitive patterns (*e.g.*, the trunks of the trees). In summary, the proposed DemoFusion achieves both rich local detail and strong global semantic coherence by modifying MultiDiffusion style patch-wise denoising paths to maximise the global context available for each path.

**Quantitative Results:** For quantitative comparison, we adopt 3 widely-used metrics: FID (Fréchet Inception Distance) [10], IS (Inception Score) [30], and CLIP Score [27].

Considering that FID and IS require resizing images to  $299^2$ , which is not very suitable for high-resolution image assessment, inspired by [4], we additionally crop local patches of  $1 \times$  resolution and then resize them to calculate these metrics, termed  $\text{FID}_{\text{crop}}$  and  $\text{IS}_{\text{crop}}$ . The CLIP Score assesses the entire image’s semantics; thus, we do not consider evaluating local patches here. We evaluate on the LAION-5B dataset [31] with  $1K$  randomly sampled captions. Note that the results of FID and IS are related to the number of samples; therefore, the scores of  $\text{FID}_{\text{crop}}$  and  $\text{IS}_{\text{crop}}$  might be better than FID and IS due to more samples. The inference time is evaluated on an RTX 3090 GPU.

As shown in Tab. 1, DemoFusion achieved the best overall performance – securing first or second place across all metrics. As the resolution increases, DemoFusion may score slightly lower than SDXL+BSRGAN on FID and IS because BSRGAN is designed to adhere strictly to low-resolution inputs, and these metrics also downsample images to low resolution for evaluation. However, DemoFusion significantly outperforms SDXL+BSRGAN on  $\text{FID}_{\text{crop}}$  and  $\text{IS}_{\text{crop}}$ , indicating that DemoFusion can provide high-resolution local details. Besides, we observed that MultiDiffusion surpassed SCALECRAFTER on crop-based metrics due to these metrics’ lack of an assessment of the overall structure of the image. Therefore, we keep the general FID and IS metrics. Regarding efficiency, since DemoFusion is based on MultiDiffusion and operates progressively, it requires a longer inference time. We discuss this point further in Sec. 5.

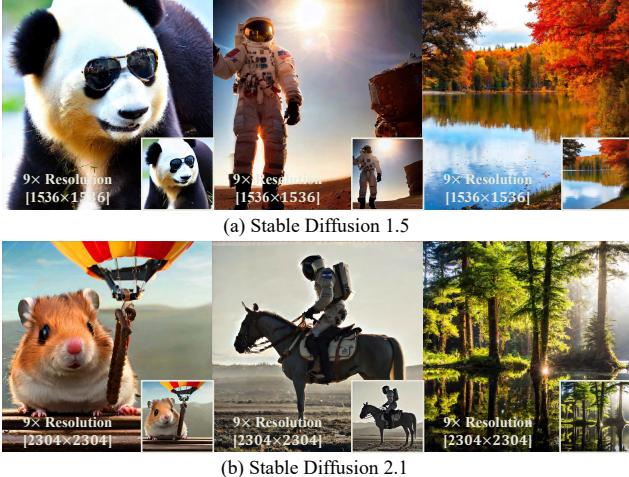


Figure 6. **Results of DemoFusion on other LDMs**, *i.e.*, Stable Diffusion 1.5 (default resolution of  $512^2$ ) and Stable Diffusion 2.1 (default resolution of  $768^2$ ). All images are generated at  $9\times$  resolutions. Best viewed **ZOOMED-IN**.

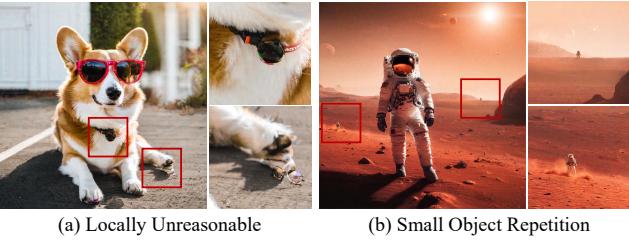


Figure 7. **Failure cases of DemoFusion**. (a) Irrational content appears locally in images with a sharp focus. (b) Small objects are repetitively present against a sparse background. All images are generated at  $9\times$  resolutions. Best viewed **ZOOMED-IN**.

## 4.2. Ablation Study

The proposed DemoFusion consists of three components: (i) progressive upscaling, (ii) skip residual, and (iii) dilated sampling. To visually demonstrate the effectiveness of these three components, we conducted experiments on all possible combinations, as shown in Fig. 5. All images are generated at  $3072^2$  ( $9\times$  resolutions). When all three components are removed, we generate at the original resolution first and then achieve higher resolutions via an “upsample-diffuse-denoise” loop. The results obtained under this setting are similar to naively generating via MultiDiffusion, with much repetitive content. However, this issue is gradually mitigated by incorporating the three proposed techniques, resulting in high-resolution images consistent with their original resolution counterparts.

Specifically, we found that continuously introducing information from the low resolution via skip residual dramatically helps maintain the overall structure to obtain acceptable results. On this basis, dilated sampling can further introduce denoising paths with global perception during the denoising process, guiding local denoising paths towards

the global optimal direction. However, these mutually independent global denoising paths introduce two drawbacks (even though we have introduced Gaussian filtering to alleviate this): (i) bringing grainy textures when generating from Gaussian noises and (ii) amplifying the artificial noises introduced during the upscaling process. The former can be alleviated by introducing skip residuals, while the latter can be addressed by progressive upscaling, which prevents the strong artificial noises brought by direct large-scale upscaling. Overall, the three proposed techniques are complementary and indispensable. It is fascinating to see how well they work together.

## 5. Limitations and Opportunities

DemoFusion exhibits limitations in the following aspects: (i) The nature of MultiDiffusion-style inference requires high computational load due to the overlapped denoising paths, and the progressive upscaling also prolongs inference times. (ii) As a tuning-free framework, DemoFusion’s performance is directly correlated with the underlying LDM. In Fig. 6, we show the results based on other LDMs (Stable Diffusion 1.5 and Stable Diffusion 2.1), where DemoFusion is still effective, but the results are less astonishing than those on SDXL. (iii) DemoFusion entirely depends on the LDMs’ prior knowledge of cropped images, and therefore, local irrational content may appear when generating sharp close-up images, as depicted in Fig. 7 (a). (iv) Although we have significantly mitigated the issue of repetitive content, the possibility of small repetitive content in background regions remains (see Fig. 7 (b)).

Behind these limitations, opportunities exist: (i) DemoFusion functions by fusing multiple denoising paths of the original size. This allows it to implement each denoising step in mini-batches, preventing the expected exponential increase in memory requirements. (ii) Although progressive upscaling requires more passes, users can acquire low-resolution intermediate results as “previews” within several seconds, facilitating rapid prompt iteration. (iii) The priors of current LDMs regarding image crops are solely derived from the general training scheme, which has already resulted in impressive performance. Training a bespoke LDM for a DemoFusion-like framework may be a promising direction to explore.

## 6. Conclusion

In this paper, we introduce DemoFusion, a tuning-free framework that integrates plug-and-play with open-source GenAI models to achieve higher-resolution image generation. DemoFusion is built upon MultiDiffusion and introduces *Progressive Upscaling*, *Skip Residual*, and *Dilated Sampling* techniques to enable generation with both global semantic coherence and rich local details. DemoFusion per-

suasively demonstrates the possibility of LDMs generating images at higher resolutions than those used for training and the untapped potential of existing open-source GenAI models. By advancing the frontier of high-resolution image generation without additional training or prohibitive memory requirements for inference, we hope that DemoFusion can help democratize high-resolution image generation.

## References

- [1] Stability AI. Stable diffusion: A latent text-to-image diffusion model. <https://stability.ai/blog/stable-diffusion-public-release>, 2022. 2
- [2] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. In *ICML*, 2023. 2, 3, 5, 7
- [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *CVPR*, 2023. 3
- [4] Lucy Chai, Michael Gharbi, Eli Shechtman, Phillip Isola, and Richard Zhang. Any-resolution training for high-resolution image synthesis. In *ECCV*, 2022. 7
- [5] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021. 3
- [6] Xiao Han, Yukang Cao, Kai Han, Xiatian Zhu, Jiankang Deng, Yi-Zhe Song, Tao Xiang, and Kwan-Yee K Wong. Headsculpt: Crafting 3d head avatars with text. In *NeurIPS*, 2023. 3
- [7] Yingqing He, Shaoshu Yang, Haoxin Chen, Xiaodong Cun, Menghan Xia, Yong Zhang, Xintao Wang, Ran He, Qifeng Chen, and Ying Shan. Scalecrafter: Tuning-free higher-resolution visual generation with diffusion models. *arXiv preprint arXiv:2310.07702*, 2023. 3, 5, 7, 11, 16
- [8] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv preprint arXiv:2211.13221*, 2023. 3
- [9] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *ICLR*, 2022. 3, 4
- [10] 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. *NeurIPS*, 2017. 7
- [11] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3
- [12] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *The Journal of Machine Learning Research*, 2022. 3
- [13] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for high resolution images. *arXiv preprint arXiv:2301.11093*, 2023. 3
- [14] Ajay Jain, Amber Xie, and Pieter Abbeel. Vectorfusion: Text-to-svg by abstracting pixel-based diffusion models. In *CVPR*, 2023. 3
- [15] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018. 4
- [16] Yuseung Lee, Kunho Kim, Hyunjin Kim, and Minhyuk Sung. Syncdiffusion: Coherent montage via synchronized joint diffusions. *arXiv preprint arXiv:2306.05178*, 2023. 3
- [17] Tingting Liao, Hongwei Yi, Yuliang Xiu, Jiaxaling Tang, Yangyi Huang, Justus Thies, and Michael J Black. Tada! text to animatable digital avatars. *arXiv preprint arXiv:2308.10899*, 2023. 3
- [18] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *CVPR*, 2023. 3
- [19] MidJourney. Midjourney: An independent research lab. <https://www.midjourney.com/>, 2022. 2
- [20] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *CVPR*, 2023. 3, 4
- [21] Chong Mou, Xiantao Wang, Liangbin Xie, Jian Zhang, Zhonggang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 3
- [22] OpenAI. Dall-e: Creating images from text. <https://openai.com/blog/dall-e/>, 2021. 2
- [23] OpenAI. Chatgpt: Large-scale language models. <https://www.openai.com/blog/chatgpt>, 2022. 16
- [24] 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. *arXiv preprint arXiv:2307.01952*, 2023. 1, 2, 3, 5, 7, 11, 13, 14
- [25] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2022. 3
- [26] Zhiyu Qu, Tao Xiang, and Yi-Zhe Song. Sketchdreamer: Interactive text-augmented creative sketch ideation. In *BMVC*, 2023. 3
- [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 7
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3
- [29] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023. 3
- [30] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *NeurIPS*, 2016. 7

- [31] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 2022. 7
- [32] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 3
- [33] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 3
- [34] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. *arXiv preprint arXiv:2305.07015*, 2023. 3
- [35] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing Shen, Dong Chen, Fang Wen, Qifeng Chen, et al. Rodin: A generative model for sculpting 3d digital avatars using diffusion. In *CVPR*, 2023. 3
- [36] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *CVPR*, 2023. 3
- [37] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaohu Qie, and Shenghua Gao. Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models. In *CVPR*, 2023. 3
- [38] Fisher Yu and Koltun Vladlen. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016. 5
- [39] Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *CVPR*, 2021. 3, 5, 7, 16
- [40] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 3, 16
- [41] Qingping Zheng, Yuanfan Guo, Jiankang Deng, Jianhua Han, Ying Li, Songcen Xu, and Hang Xu. Any-size-diffusion: Toward efficient text-driven synthesis for any-size hd images. *arXiv preprint arXiv:2308.16582*, 2023. 2, 3, 4

## Appendix



Figure 8. During the progressive upscaling process, we diffuse  $\mathbf{z}_0^s$  to different time-steps  $t$ , and then denoise it back to obtain the results. The number of training time-step is 1000.

PU	SR	DS	FID	IS	$\text{FID}_{crop}$	$\text{IS}_{crop}$	CLIP
✗	✗	✗	95.28	13.92	80.11	19.61	28.13
✓	✗	✗	90.77	13.95	79.23	20.08	28.52
✗	✓	✗	79.93	15.19	74.17	22.48	29.40
✗	✗	✓	94.35	14.89	82.32	19.64	28.85
✓	✓	✗	75.92	15.66	72.98	23.20	29.50
✓	✗	✓	89.26	15.02	80.04	21.86	28.87
✗	✓	✓	76.53	15.71	73.22	23.09	29.48
✓	✓	✓	<b>74.11</b>	<b>16.11</b>	<b>70.34</b>	<b>24.28</b>	<b>29.57</b>

Table 2. **Quantitative results of the ablation study.** The best results are marked in **bold**. Impact of components: Progressive Upscaling (PU), Skip Residual (SR), and Dilated Upsampling (DS).

## A. Pseudo Code

We further illustrate the image synthesis process of DemoFusion in Algorithm 1.

## B. Implementation Details

In cases where it is not explicitly stated, all the results in this paper are obtained based on SDXL with a DDIM scheduler of 50 steps. The guidance scale for all denoising paths is set to 7.5. The crop size of MultiDiffusion is set to be aligned with the maximum training size of pre-trained LDMs, *e.g.*,  $h = w = 128$  for SDXL, and the stride is set to be  $d_h = \frac{h}{2}$  and  $d_w = \frac{w}{2}$ . Each crop’s position is subjected to a slight random perturbation, with maximum offsets of

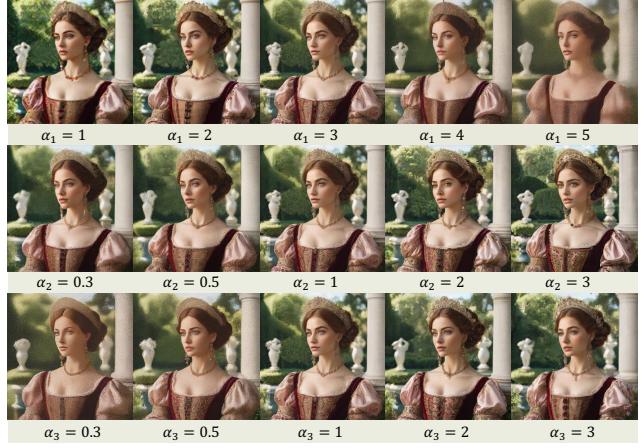


Figure 9. **Results with different  $\alpha_1$ ,  $\alpha_2$ , and  $\alpha_3$ .** All images are generated at  $3072^2$  ( $9 \times$  resolutions). Best viewed **ZOOMED-IN**.

$\frac{h}{16}$  and  $\frac{w}{16}$  in vertical and horizontal directions, respectively, further preventing the occurrence of seam issues.

When generating images with varying aspect ratios, we ensure that the longer side aligns with the maximum training size. Three scale factors  $\alpha_1$ ,  $\alpha_2$ ,  $\alpha_3$  were set to 3, 1, and 1 respectively. The Gaussian filter’s standard deviation decreases from  $\sigma_1 = 1$  to  $\sigma_2 = 0.01$ . To decode high-resolution images, we also employed a tiled decoder strategy as [7] and some open-source projects. To eliminate seams between tiles, we sample a larger range of features around each tile during the decoding process.

Note that SDXL permits the input of a coarse cropping condition [24], *i.e.*, the coordinates of the top-left corner of the cropping area. Therefore, when utilizing SDXL, we additionally input the coordinates of the top-left corner of the corresponding patch as a condition for local denoising paths. However, we observed that the presence or absence of this condition does not significantly impact the results. Besides, DemoFusion initiates the generation process from the highest resolution of the LDM during the first phase. When generating images with varying aspect ratios, we ensure that the longer side aligns with the highest resolution.

## C. More Experimental Results

### C.1. Diffusing to Different Time-step $t$

In Fig. 8, we illustrate that, when skip residual is removed, the effects of different time-steps we denoise to within the “upsample-diffuse-denoise” loop. This provides evidence for our discussion in the main text – the larger the  $t$ , the more information is lost, which weakens the global perception; the smaller the  $t$ , the stronger the noise introduced by upsampling.

<https://github.com/pkuliyi2015/multidiffusion-upscaler-for-automatic1111>

---

**Algorithm 1** Image Synthesis Process of DemoFusion

---

```

1: ##### Phase 1 #####
2:  $\mathbf{z}_T^0 \sim \mathcal{N}(0, I)$  ▷ Random Initialization
3: for  $t = T$  to 1 do
4:    $p_\theta(\mathbf{z}_{t-i}^1 | \mathbf{z}_t^1)$  ▷ Denoising Step
5: end for
6: ##### Phase 2 to S #####
7: for  $s = 2$  to  $S$  do
8:    $\text{inter}(\mathbf{z}'_0^s | \mathbf{z}_0^{s-1})$  ▷ Upsampling
9:   for  $t = 1$  to  $T$  do
10:     $q(\mathbf{z}'_t^s | \mathbf{z}'_{t-1}^s)$  ▷ Diffusion Step
11:   end for
12:   for  $t = T$  to 1 do
13:      $\hat{\mathbf{z}}_t^s = c_1 \times \mathbf{z}'_t^s + (1 - c_1) \times \mathbf{z}_t^s$  ▷ Skip Residual
14:      $\mathcal{S}_{\text{local}}(\hat{\mathbf{z}}_t^s) \rightarrow Z_t^{\text{local}}$  ▷ Crop Sampling (MultiDiffusion)
15:      $\mathcal{S}_{\text{global}}(\hat{\mathbf{z}}_t^s) \rightarrow Z_t^{\text{global}}$  ▷ Dilated Sampling
16:     for  $\hat{\mathbf{z}}_{n,t}^s$  in  $Z_t^{\text{local}}$  do
17:        $p_\theta(\mathbf{z}_{n,t-i}^s | \hat{\mathbf{z}}_{n,t}^s)$  ▷ Local Path Denoising Step (MultiDiffusion)
18:     end for
19:     for  $\hat{\mathbf{z}}_{m,t}^s$  in  $Z_t^{\text{global}}$  do
20:        $p_\theta(\mathbf{z}_{m,t-i}^s | \hat{\mathbf{z}}_{m,t}^s)$  ▷ Global Path Denoising Step
21:     end for
22:      $\mathcal{R}_{\text{local}}(Z_{t-1}^{\text{local}}) \times (1 - c_2) + \mathcal{R}_{\text{global}}(Z_{t-1}^{\text{global}}) \times c_2 \rightarrow \mathbf{z}_t^s$  ▷ Fusing Local and Global Paths
23:   end for
24: end for
25: return  $\mathbf{x}_0^S = \mathcal{D}(\mathbf{z}_0^S)$  ▷ Decoding to Image

```

---

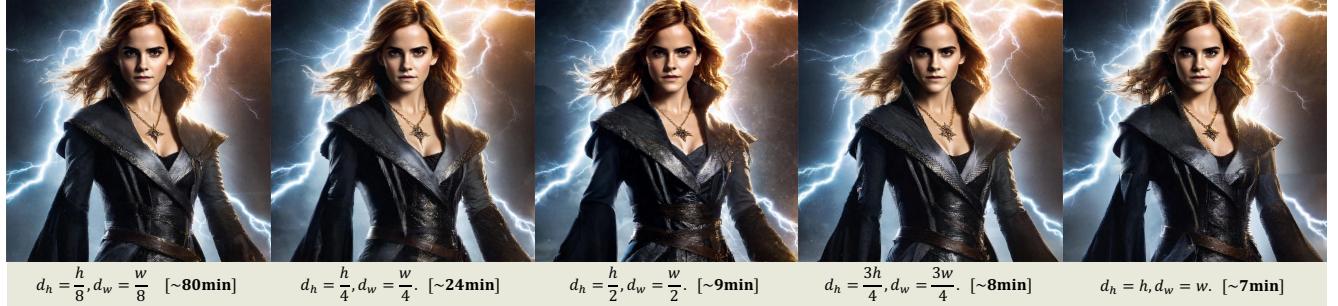


Figure 10. Results with different strides  $d_h$  and  $d_w$ . All images are generated at  $3072^2$  ( $9 \times$  resolutions). Best viewed **ZOOMED-IN**.

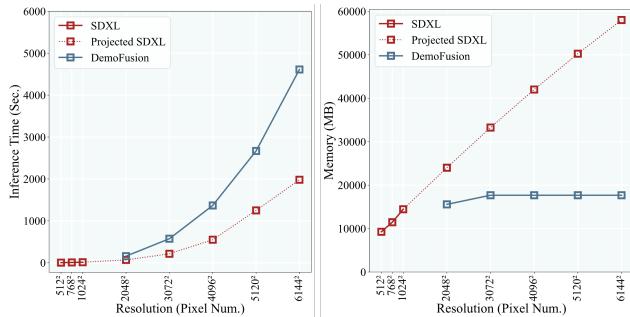


Figure 11. Inference time of SDXL versus DemoFusion (Left). Memory demands of SDXL versus DemoFusion (Right).

## C.2. Quantitative Results of Ablation Study

The quantitative results of the ablation study are shown in Tab. 2. Here, we only experiment with the resolution of  $4096^2$ .

## C.3. Effects of Scale Factors $\alpha_1$ , $\alpha_2$ , and $\alpha_3$

A shared understanding of the DM’s denoising process is that the DM first determines the coarse details and then gradually refines the local details. In line with this understanding, we adopt a unified strategy: utilizing cosine descending weights, we assign greater weights to skip residuals, dilated sampling, and accompanying Gaussian filtering in the early stages of the denoising process, gradually de-

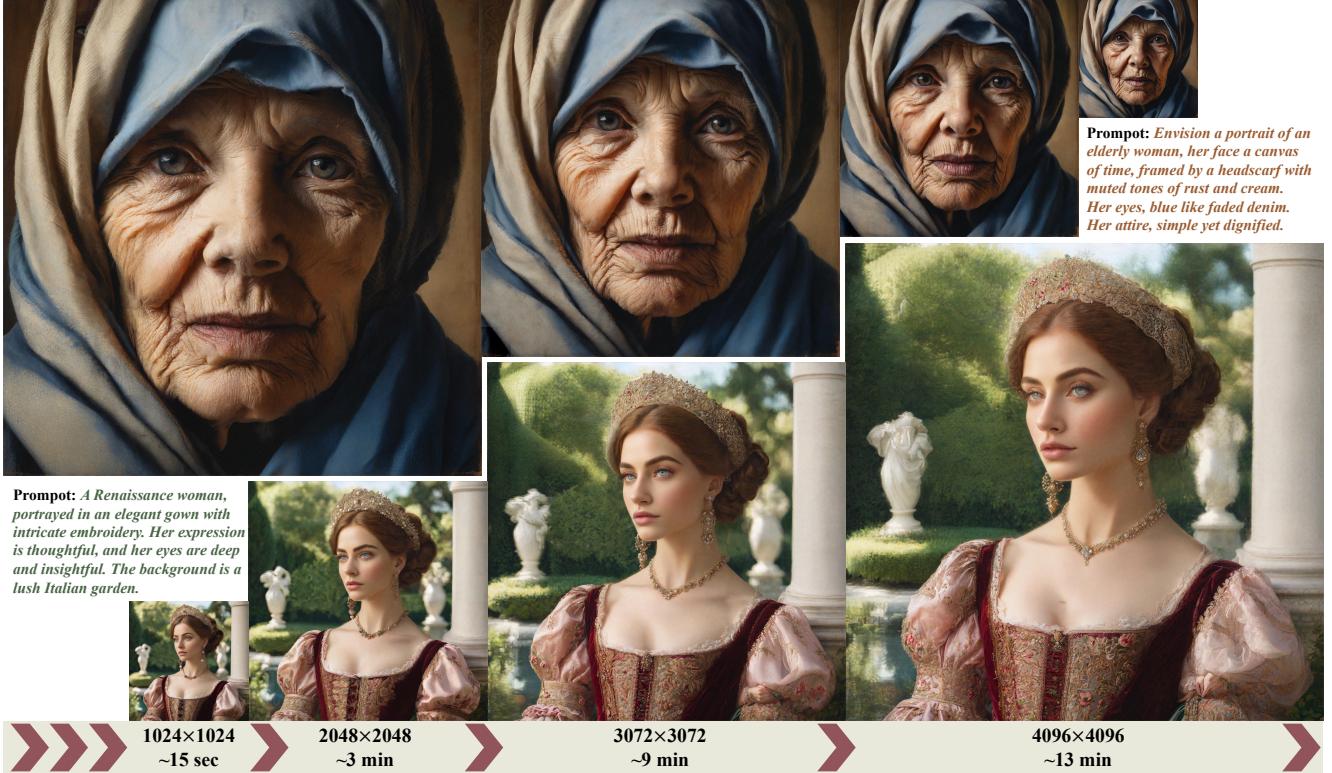


Figure 12. **Illustration of the progressive upscaling process.** The time required for each phase is indicated. Best viewed ZOOMED-IN.

creasing the weights as denoising progresses. Despite this unified approach, the three components still need distinct scale factors to control the descent rate.

Through grid search, we obtained the globally optimal parameter combination. In Fig. 9, we varied only one parameter at a time while keeping the others at their optimal values to demonstrate the impact of each parameter on the results. Note that a larger scale factor means a faster decline, which weakens the effect of this item, and vice versa.

According to the experimental results, when the skip residual effect is too strong (*i.e.*,  $\alpha_1 = 1$ ), we observe significant artificial noise caused by upsampling. Because these factors interact with each other, when  $\alpha_1 = 5$ , we observe that the results are close to the one when  $\alpha_3 = 1$  – Gaussian filtering leads to excessive smoothing of latent representation. The trade-off of dilated sampling is – too large a weight (*i.e.*,  $\alpha_2 = 1$ ) can result in grainy images, while too small a weight (*i.e.*,  $\alpha_2 = 5$ ) fails to provide sufficient global perception, leading to noticeable issues of content repetition. Regarding Gaussian filtering, excessive strength can lead to over-smoothing of the latent representation, while too small a strength can weaken the global denoising paths due to lack of interaction, resulting in content repetition and grainy appearance.

#### C.4. Effect of Stride Sizes $d_h$ and $d_w$

In general, the stride size  $d_h$  and  $d_w$  in MultiDiffusion determines the extent of the seam issue of images – a smaller stride means more seamless images, but at the same time, brings more overlapping computation. For DemoFusion, due to the proposed *Progressive Upscaling*, *Skip Residual*, and *Dilated Sampling techniques*, there is a better consistency between patches, relaxing the stride size requirement. In Fig. 10, we showcase the performance of DemoFusion and the corresponding inference time for different stride sizes. It can be observed that even in the case of  $d_h = h$  and  $d_w = w$ , *i.e.*, no overlap between patches, we still achieve good global semantic coherence, while when  $d_h < h$  and  $d_w < w$ , all the generated images have no noticeable seams; ultimately, in order to balance the performance and efficiency, we chose  $d_h = \frac{h}{2}$  and  $d_w = \frac{w}{2}$ .

#### C.5. Resource Demands of DemoFusion

In Fig. 11, we illustrate the resource demand comparison of DemoFusion and the original SDXL [24]. Note that SDXL cannot generate valid content at resolutions higher than  $1024^2$ . We just calculate the expected resource demands by assuming we have a high-resolution SDXL under the current framework. It can be seen that DemoFusion achieves high-resolution image generation on limited com-



Figure 13. **More selected landscape samples of DemoFusion versus SDXL [24]** (all images in the figure are presented at their actual sizes). All generated images are produced using a single RTX 3090 GPU. Best viewed **ZOOMED-IN**.

putational resources while paying a little bit more time cost.

## D. More Visualizations

### D.1. The Progressive Upscaling Process

To better demonstrate how the model progressively generates images with different resolutions, in Fig. 12, we show the model outputs at each phase of the generation process.

We can observe that DemoFusion does an excellent job of achieving global consistency under different resolutions, indicating the reason for its success in resolving content repetition.

### D.2. More Landscape Samples

In Fig. 13, we have supplemented more samples to show the performance of DemoFusion, and in particular, we further

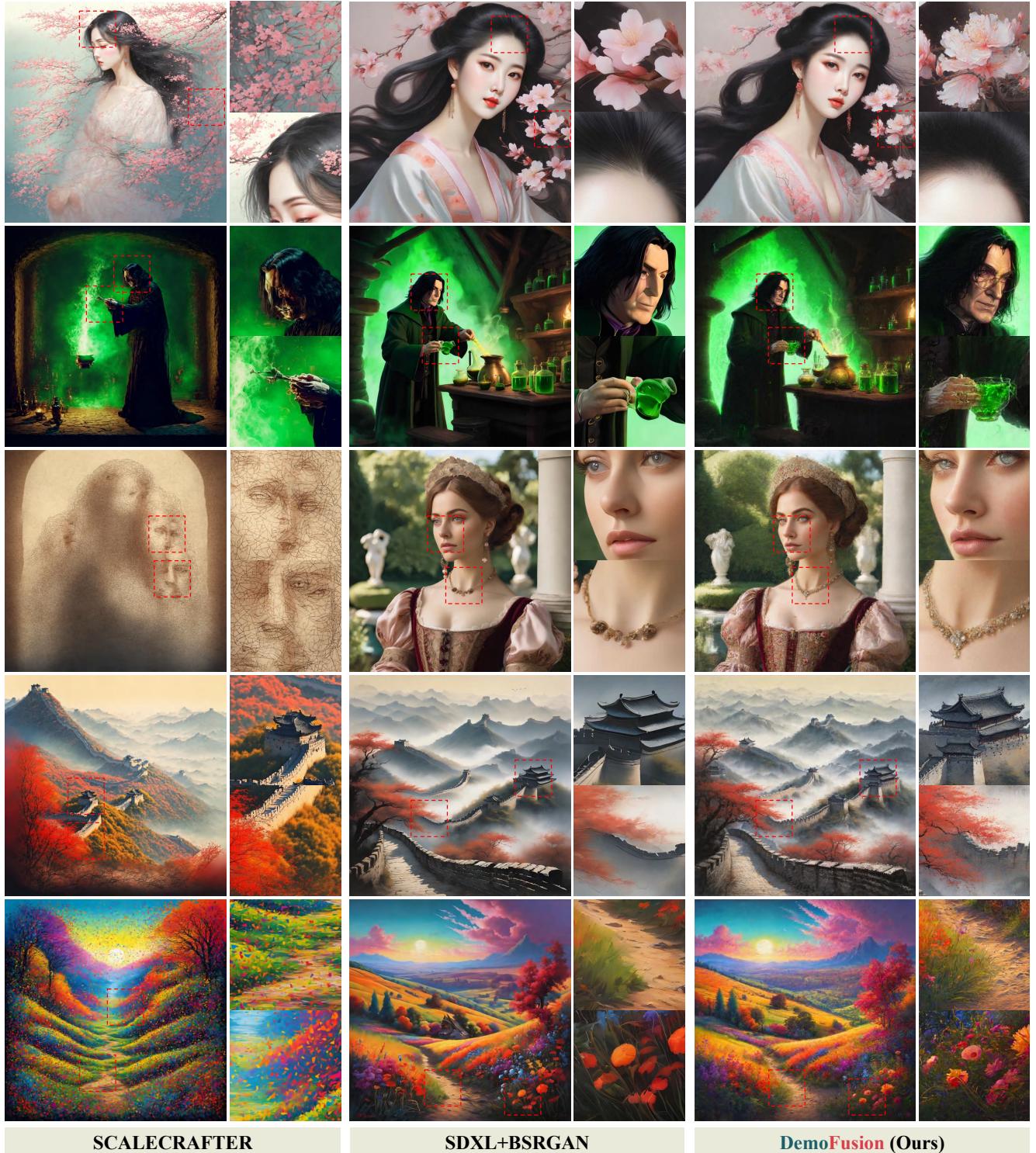


Figure 14. **More Qualitative comparison results.** All images are generated at  $4096^2$  ( $16\times$  resolutions). Local details have already been zoomed in, but it's still recommended to **ZOOM IN** for a closer look.



Figure 15. Results of DemoFusion combining with ControlNet [40]. All images are generated at  $3072^2$  ( $9\times$  resolutions). Best viewed ZOOMED-IN.



Figure 16. Results of upscaling real images. All images are upscaled to  $3072^2$ . Best viewed ZOOMED-IN.

show results at the resolution of  $8192 \times 4096$  ( $64\times$  upscaling compared to the initial resolution of  $512 \times 1024$ ).

### D.3. More Comparison Results

In Fig. 14, we have supplemented more comparison results with **SDXL+BSRGAN** [39] and **SCALECRAFTER** [7] to demonstrate the effectiveness of DemoFusion. The results are consistent with those in the main text. Compared to SDXL+BSRGAN, DemoFusion provides better local details; while compared to SCALECRAFTER, DemoFusion better preserves the performance of SDXL during upscale.

## E. More Applications

### E.1. Combining with ControlNet

The tuning-free characteristic of DemoFusion enables seamless integration with many LDM-based applications. E.g., DemoFusion combined with ControlNet [40] can

achieve controllable high-resolution generation. In Fig. 15, we showcase examples using Canny edge and human pose as conditions.

### E.2. Upscaling Real Images

Since DemoFusion works in a progressive manner, we can replace the output of phase 1 with representations obtained by encoding real images, thereby achieving upscaling of real images. However, we carefully avoid using the term “super resolution”, as the outputs tend to lean towards the latent data distribution of the base LDM, making this process more akin to image generation based on a real image. The results are shown in the Fig. 16.

## F. Prompts Used in This Paper

All prompts used in this paper are taken from the internet or generated by ChatGPT [23]. They are summarised here.

### Fig. 1 in the main text:

- *Steampunk makeup, in the style of vray tracing, colorful impasto, uhd image, indonesian art, fine feather details with bright red and yellow and green and pink and orange colours, intricate patterns and details, dark cyan and amber makeup. Rich colourful plumes. Victorian style.*
- *Stunning feminine body, commercial image, beautiful girl from Spain, holographic photography shoots, large body of water sprayed, liquid splashing all over the places, street pop, luminous palette, close up, realistic impressionism, shiny/glossy, extreme colorsplash, behind that a universe of vortex of fire waves and ice waves, around fire splashes and ice splashes and floral, bonsais, roots, smoke swirls, dust swirls, tentacles of fire and ice, s-curve composition, leading lines, cinematic, style of hokusai, unreal engine, octane render, asymmetric, golden ratio, style of hokusai, liquid splashes, merging, melting, splashing, droplets, mixing, fading away, exploding,*

*swirling, intricate detail, modelshoot style, dreamlikeart, dramatic lighting. 8k, highly detailed, trending artstation.*

- *The beautiful scenery of Seattle, painting by Al Capp.*
- *By Tang Yau Hoong, ultra hd, realistic, vivid colors, highly detailed, UHD drawing, pen and ink, perfect composition, beautiful detailed intricate insanely detailed octane render trending on artstation, 8k artistic photography, photorealistic concept art, soft natural volumetric cinematic perfect light, ultra hd, realistic, vivid colors, highly detailed, UHD drawing, pen and ink, perfect composition, beautiful detailed intricate insanely detailed octane render trending on artstation, 8k artistic photography, photorealistic concept art, soft natural volumetric cinematic perfect light.*
- *A cute and adorable fluffy puppy wearing a witch hat in a halloween autumn evening forest, falling autumn leaves, brown acorns on the ground, halloween pumpkins spider-webs, bats, a witch's broom.*
- *A robot standing in the rain reading newspaper, rusty and worn down, in a dystopian cyberpunk street, photorealistic, urbanpunk.*
- *Einstein, a bronze statue, with a fresh red apple on his head, by Bruno Catalano.*
- *A woman in a pink dress walking down a street, cyberpunk art, inspired by Victor Mosquera, conceptual art, style of raymond swanland, yume nikki, restrained, robot girl, ghost in the shell.*
- *Photo of a rhino dressed suit and tie sitting at a table in a bar with a bar stools, award winning photography, Elke vogelsang.*
- *An astronaut riding a horse on the moon, oil painting by Van Gogh.*
- *Classic traditional cornucopia at the fall harvest festival, farm in the background, high quality masterful still-life painting, American pastoral, oil painting, festive spirit, vibrant cultural tradition, Autumnal atmosphere, vibrant rich colors.*

#### **Fig. 2 in the main text:**

- *An astronaut riding a horse on the moon, oil painting by Van Gogh.*

#### **Fig. 3 in the main text:**

- *An astronaut riding a horse on the moon, oil painting by Van Gogh.*

#### **Fig. 4 in the main text:**

- *A cute teddy bear in front of a plain white wall, warm and brown fur, soft and fluffy.*
- *Emma Watson as a powerful mysterious sorceress, casting lightning magic, detailed clothing.*
- *Primitive forest, towering trees, sunlight falling, vivid colors.*

#### **Fig. 5 in the main text:**

- *Astronaut in a jungle, cold color palette, muted colors,*

*detailed, 8k.*

- *Emma Watson as a powerful mysterious sorceress, casting lightning magic, detailed clothing.*

#### **Fig. 6 in the main text:**

- *A panda wearing sunglasses.*
- *Astronaut on Mars During sunset.*
- *A serene lakeside during autumn, with trees displaying a palette of fiery colors.*
- *A hamster piloting a tiny hot air balloon.*
- *An astronaut riding a horse.*
- *A deep forest clearing with a mirrored pond reflecting a galaxy-filled night sky.*

#### **Fig. 7 in the main text:**

- *A corgi wearing cool sunglasses.*
- *Astronaut on Mars During sunset.*

#### **Fig. 8 in Appendix:**

- *Astronaut in a jungle, cold color palette, muted colors, detailed, 8k.*

#### **Fig. 9 in Appendix:**

- *A Renaissance noblewoman, portrayed in an elegant gown with intricate embroidery. Her expression is thoughtful, and her eyes are deep and insightful. The background is a lush Italian garden, reflecting the artistic style of the High Renaissance.*

#### **Fig. 10 in Appendix:**

- *Emma Watson as a powerful mysterious sorceress, casting lightning magic, detailed clothing.*

#### **Fig. 12 in Appendix:**

- *Envision a portrait of an elderly woman, her face a canvas of time, framed by a headscarf with muted tones of rust and cream. Her eyes, blue like faded denim. Her attire, simple yet dignified.*
- *A Renaissance noblewoman, portrayed in an elegant gown with intricate embroidery. Her expression is thoughtful, and her eyes are deep and insightful. The background is a lush Italian garden, reflecting the artistic style of the High Renaissance.*

#### **Fig. 13 in Appendix:**

- *Realistic oil painting of a stunning model merged in multicolor splash made of finely torn paper, eye contact, walking with class in a street.*
- *A painting of a beautiful graceful woman with long hair, a fine art painting, by Qiu Ying, no gradients, flowing sakura silk, beautiful oil painting.*
- *Katsushika Hokusai's Japanese depiction of a very turbulent sea with massive waves. The background shows a beautiful dark night over a illuminated village. The colors are red and yellow, mood lighting Imagine a dreamlike scene blending the swirling cosmic colors of Vincent van Gogh's Starry Night with the surreal celestial precision of Salvador Dalí.*
- *Character of lion in style of saiyan, mafia, gangsta, city-lights background, Hyper detailed, hyper realistic, unreal*

*engine ue5, cgi 3d, cinematic shot, 8k.*

- *Santa Claus riding on top of a turkey, with very large bag of gifts, snow, ice, very cold place, realistic digital art, blurred background, expansive lighting, 4k, light gray and blue color palette, sharp and fine intricate details defined.*
- *Best Quality, Masterpiece, steampunk theme, centered, front cover of fashion magazine, concept art, design, magazine design, 1girl, cute, blonde ponytail hair, gothic steampunk dress, model pose, (epic composition, epic proportion), vibrant color, text, diagrams, advertisements, magazine title, typography.*
- *Burning pile of money, epic composition, digital painting, emotionally profound, thought-provoking, intense and brooding tones, high quality, masterpiece.*
- *A pastoral scene with shepherds, flocks, and rolling hills, in the tradition of a Jean-François Millet landscape.*
- *A painting of brooklyn new york 1940 storefronts, by John Kay, highly textured, rich colour and detail, ballard, deep colours, style of raymond swanland, trio, oil painting, h 768, well worn, displayed, detailed 4 k oil painting, glenn barr, textured oil on canvas, looking cute.*
- *A swirling night sky filled with bright stars and a small village below, inspired by Vincent van Gogh's Starry Night.*
- *Portrait of a bear as a roman general, with a helmet, decorative, fantasy environment, oil painting, masterpiece, detailed, sharp, clear, cinematic lights.*
- *The Great Wall of China winding through mist-covered mountains, captured in the delicate brushwork and harmonious colors of a traditional Chinese landscape painting.*
- *RAW photo of a mountain lake landscape, clean water, 8k, UHD.*

#### **Fig. 14 in Appendix:**

- *A painting of a beautiful graceful woman with long hair, a fine art painting, by Qiu Ying, no gradients, flowing sakura silk, beautiful oil painting.*
- *Professor Snape brewing a potion in the dungeon, the room illuminated by the green glow of the cauldron.*
- *A Renaissance noblewoman, portrayed in an elegant gown with intricate embroidery. Her expression is thoughtful, and her eyes are deep and insightful. The background is a lush Italian garden, reflecting the artistic style of the High Renaissance.*
- *The Great Wall of China winding through mist-covered mountains, captured in the delicate brushwork and harmonious colors of a traditional Chinese landscape painting.*
- *Summer landscape, vivid colors, a work of art, grotesque, Mysterious.*

#### **Fig. 15 in Appendix:**

- *A Samoyed wearing a sunglasses, sticking out its tongue,*

*dslr image, 8k.*

- *A Corgi wearing a sunglasses, sticking out its tongue, dslr image, 8k.*
- *A German Shepherd wearing a sunglasses, sticking out its tongue, on the grass, dslr image, 8k.*
- *A Husky wearing a sunglasses, sticking out its tongue, dslr image, 8k.*
- *An Australian Cattle Dog wearing a sunglasses, sticking out its tongue, style of Gian Lorenzo Bernini.*
- *A medieval knight standing in a lush forest, oil painting by Van Gogh.*
- *A gardener tending to a colorful, blooming garden, oil painting by Van Gogh.*
- *A robot exploring the ruins of an ancient civilization, watercolor by MORILAND.*
- *A ghost haunting an abandoned Victorian mansion, watercolor by MORILAND.*
- *An astronaut floating in space, watercolor by MORILAND.*

#### **Fig. 16 in Appendix:**

- *A cute corgi on the lawn.*
- *A portrait of Mr. Bean (Rowan Atkinson).*
- *Japanese Ukiyo-e, Kanagawa Surfing Sato.*
- *A cute panda on a tree trunk.*
- *A Portrait of Albus Dumbledore.*
- *A Chinese Painting of the Great Wall.*