
High-Resolution Image Editing via Multi-Stage Blended Diffusion

Johannes Ackermann*

The University of Tokyo

ackermann@ms.k.u-tokyo.ac.jp

Minjun Li

Preferred Networks Inc.

minjunli@preferred.jp

Abstract

Diffusion models have shown great results in image generation and in image editing. However, current approaches are limited to low resolutions due to the computational cost of training diffusion models for high-resolution generation. We propose an approach that uses a pre-trained low-resolution diffusion model to edit images in the megapixel range. We first use Blended Diffusion to edit the image at a low resolution, and then upscale it in multiple stages, using a super-resolution model and Blended Diffusion. Using our approach, we achieve higher visual fidelity than by only applying off the shelf super-resolution methods to the output of the diffusion model. We also obtain better global consistency than directly using the diffusion model at a higher resolution.

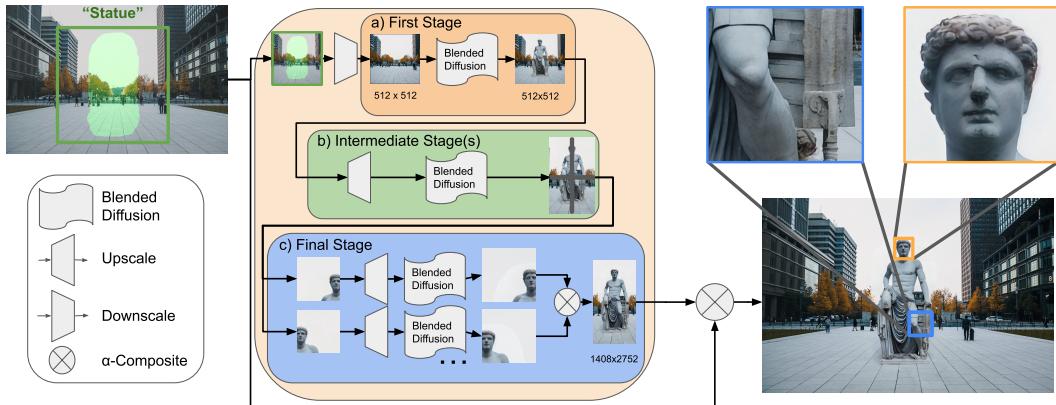


Figure 1: Our approach performs high-resolution text-guided image editing in multiple stages. In the first stage a), we apply Blended Diffusion [2], given a masked region and a text prompt. In each following stage b), we first upscale the image using an off the shelf super-resolution model and then use Blended Diffusion, starting at an intermediate diffusion step, to improve the image quality and ensure consistency with the input prompt. c) When the output resolution of a stage is too large to fit into the GPU memory, we split the image into multiple segments, apply upscaling and Blended Diffusion to them separately and alpha-composite the results.

Over the last year, along with great advances in text-guided image generation [13, 17], text-guided image editing approaches have been proposed with impressive results [2, 10]. Avrahami et al. [2] published a method called *Blended Diffusion*, which allows us to reuse pre-trained text-guided diffusion models to edit images given a masked region and a text prompt. However, as training large

*Work done during an internship at Preferred Networks Inc.

15 diffusion models is computationally expensive, publicly available models are limited to resolutions
 16 of at most 512x512 pixels. Directly applying these models to higher resolutions leads to incoherent
 17 results, with repeated patterns and elements. To overcome this issue, we propose a multi-stage
 18 approach that, by first editing the image at a low resolution and then upscaling the image in multiple
 19 stages, is able to edit large areas coherently at a high resolution.

20 **Approach** As illustrated in Fig. 1, we begin by cropping the image to a square area around the
 21 masked region and downscale this region to match the training resolution of the diffusion model.
 22 In the first stage (Fig. 1a), we use Blended Latent Diffusion [2] to obtain a batch of five edited
 23 outputs. We found it crucial to utilize Repaint [8] together with Blended Diffusion here, leading
 24 to a significantly higher consistency between the masked and unmasked region. From the batch of
 25 results we select the best sample by similarity to the prompt using CLIP [12], following related work
 26 [1, 14]. We then perform one or more intermediate stages (Fig. 1b) of upscaling using Real-ESRGAN
 27 [19] followed by Blended Diffusion. We first take a forward step in the diffusion process to an
 28 intermediate timestep [9] and then denoise with Blended Diffusion. As background input for Blended
 29 Diffusion, we found it helpful to not just add the same amount of noise as to the masked region, but
 30 also to first low-pass filter the masked region to match the spatial resolution of the edited region.
 31 Without this preprocessing, the diffusion process tends to blur the region we are editing, as we show
 32 in the Appendix A.4. When the output resolution of a stage is too large (Fig. 1c), we divide the image
 33 into a grid of overlapping regions, processing each region separately. We then use alpha compositing
 34 [11] to recombine the processed segments. We also found performing decoder optimization [1] after
 35 each stage to significantly improve the blending between the edited region and the background.

36 **Evaluation** We implement our approach based on the latent diffusion implementation provided
 37 by Rombach et al. [16] and use Stable Diffusion v1.4 [16]. In Fig. 2, we compare our approach to
 38 two baselines: First, we directly apply Blended Diffusion [2] to the highest resolution we can fit
 39 in 32GB of VRAM, 960x960 pixels, and bilinearly upscale to the original resolution. Due to the
 40 resolution being higher than in training, the image is not globally consistent with multiple heads
 41 being generated. As a second baseline, we use the editing function of the Dall-E 2 [13] Web UI to
 42 paint in the masked region. As the Dall-E 2 model is limited to 1024x1024 pixels, we need to edit the
 43 image in six separate segments, leading to poor consistency between parts of the image. Two statues
 44 have been generated, one of which is floating above the other. Finally, using our approach, we are
 45 able to produce a high-resolution result that is globally consistent.

46 We provide a discussion of the background and related work in the Appendix, as well as additional
 47 sample images, implementation details, ablations, and attribution of the used images. Our source
 48 code will be released after acceptance and we are planning to hold a live demo at the workshop.



Figure 2: Comparison of our approach to two baselines for the prompt “Statue of Roman Emperor, Canon 5D Mark 3, 35mm, flickr”. From left to right: Blended Diffusion applied to 960x960 pixels followed by bilinear upscaling, Dall-E 2 editing in multiple segments, our proposed approach. The size of the mask is 1166x2297 pixels. Applying Blended Diffusion directly to the higher resolution leads to incoherent generation with repeated elements. Similarly, Dall-E 2 generates two statues, with one floating above the other. Our approach is able to generate a detailed, coherent image.

49 **Ethical Impact** Image editing, as well as image generation, can have significant ethical implications.
50 While they can be used for many positive uses, for example for concept art generation, general image
51 retouching or entertainment, they can also be used by malicious actors to fabricate images to pass off
52 as real. However, at the current state it is possible with closer inspection to tell which images are
53 edited. Another issue is that the pretrained model we use in our approach, Stable Diffusion v1.4 [16],
54 was trained on the LAION dataset, which is known to have significant biases with regard to ethnicity
55 and gender. See Birhane et al. [3] for an investigation of this issue.

56 References

- 57 [1] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *arXiv preprint arXiv:2206.02779*, 2022.
- 58 [2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of
59 natural images. In *CVPR*, 2022.
- 60 [3] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny,
61 pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021.
- 62 [4] Lei Cai, Hongyang Gao, and Shuiwang Ji. Multi-stage variational auto-encoders for coarse-to-
63 fine image generation. In *SIAM*, 2019.
- 64 [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*,
65 2020.
- 66 [6] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim
67 Salimans. Cascaded diffusion models for high fidelity image generation. *JMLR*, 2022.
- 68 [7] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- 69 [8] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc
70 Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*, 2022.
- 71 [9] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano
72 Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In
73 *ICLR*, 2021.
- 74 [10] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
75 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing
76 with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 77 [11] Thomas Porter and Tom Duff. Compositing digital images. In *Proceedings of the 11th annual
78 conference on Computer graphics and interactive techniques*, 1984.
- 79 [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
80 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
81 models from natural language supervision. In *ICML*, 2021.
- 82 [13] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
83 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 84 [14] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
85 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021.
- 86 [15] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images
87 with vq-vae-2. *NeurIPS*, 2019.
- 88 [16] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
89 High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- 90 [17] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed
91 Kamyar Seyed Ghaseipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al.
92 Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint
93 arXiv:2205.11487*, 2022.

- 95 [18] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In
 96 *ICLR*, 2021.
- 97 [19] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world
 98 blind super-resolution with pure synthetic data. In *ICCV Workshops*, 2021.
- 99 [20] Chenfei Wu, Jian Liang, Xiaowei Hu, Zhe Gan, Jianfeng Wang, Lijuan Wang, Zicheng Liu,
 100 Yuejian Fang, and Nan Duan. Nuwa-infinity: Autoregressive over autoregressive generation for
 101 infinite visual synthesis. *arXiv preprint arXiv:2207.09814*, 2022.

102 **A Appendix**

103 **A.1 Background**

104 Diffusion models [5] are a class of generative models, which define a forward process $q(x_t|x_{t-1})$
 105 that adds noise to a given sample, such that the marginal distribution at timestep $t = T$ is a
 106 symmetric Gaussian distribution $p(x_T) = \mathcal{N}(x_T; \mathbf{0}; \mathbf{I})$ and the marginal distribution at timestep
 107 $t = 0$ is the data distribution that we aim to learn $p(x_0) = p_D(x_0)$. To generate samples from the
 108 data distribution, Denoising Diffusion Probabilistic Model (DDPM) [5] learn the reverse process
 109 $p_\theta(x_{t-1}|x_t)$, which allows a sample to be generated by sampling $x_T \sim \mathcal{N}(\mathbf{0}; \mathbf{I})$ and iteratively
 110 denoising it with $p_\theta(x_{t-1}|x_t)$. As the number of timesteps is usually large ($T \approx 1000$), each
 111 step requiring one forward pass of the network parametrizing p_θ , both training and inference are
 112 computationally expensive. Latent diffusion models (LDMs) [16] speed up training and sampling
 113 by splitting their approach into two stages: First, they train a variational autoencoder (VAE) with
 114 decoder D_θ and encoder E_θ to obtain a lower dimensional representation of the input images. In the
 115 second stage, they define the forward process in the latent space and train a DDPM model for the
 116 reverse process.

117 Blended Diffusion [2] provides a way to edit images with text-guided diffusion models. Given a
 118 mask m , an input image x_{in} and text conditioning c , they alter the sampling process by combining
 119 the intermediate result x_t with a noised version \hat{x}_t of the input after each denoising step: $x_t \leftarrow$
 120 $m \odot x_t + (1 - m) \odot \hat{x}_t$. \hat{x}_t here is a sample of the forward diffusion process at time step t , starting
 121 from x_{in} ; $\hat{x}_t \sim q(x_t|x_0 = x_{in})$. In a subsequent work, Avrahami et al. [1] apply Blended Diffusion
 122 to LDMs. They show that directly applying Blended Diffusion leads to poor blending with the
 123 original image due to the lossy encoding of the VAE. To address this issue, they propose to optimize
 124 the weights of the VAE decoder for each edited image to minimize the following loss

$$L_{DO} = \|m \odot (x_{out} - D_\theta(z_{out}))\|_2^2 + \lambda \|(1 - m) \odot (x_{in} - D_\theta(z_{out}))\|_2^2, \quad (1)$$

125 where z_{out} is the result of the editing in latent space, $x_{out} = D_\theta(z_{out})$ is the result in pixel space
 126 when using the original decoder weights, x_{in} is the input image, and λ weights the two loss terms. In
 127 other words, this loss optimizes the unmasked area to be the similar to the unedited image and the
 128 masked area to be similar to the edited image.

129 SDEdit [9] is a method that allows unconditional or text-conditional diffusion models to be utilized as
 130 image-to-image models, projecting an input image x_{in} onto the data distribution p_D that a diffusion
 131 model p_θ was trained on. To do so, they sample from the forward diffusion process at time $T' < T$,
 132 using x_{in} as sample at time $t = 0$, obtaining $x_{T'} \sim q(x_{T'}|x_0 = x_{in})$. They then use the trained
 133 diffusion model $p_\theta(x_{t-1}|x_t)$ to sample from the reverse process starting with $x_{T'}$, obtaining the
 134 result at $t = 0$.

135 Repaint [8] is a method that utilizes unconditional diffusion models for inpainting. Starting from
 136 Gaussian noise at $t = T$, similar to Blended Diffusion [2], after every denoising step they replace
 137 the unmasked area with a noised version of the input image. However, instead of sampling from the
 138 reverse process by going from step $t = T$ to $t = 0$ linearly, they repeat each denoising step R times,
 139 by taking forward and backward steps in the diffusion process iteratively. While this increases the
 140 runtime, they show that it also significantly increases the quality of the inpainting result.

141 **A.2 Implementation Details**

142 See Figure 3 for a detailed diagram of our approach, with outputs of all intermediate processing steps
 143 being shown, and Algorithm 1 for an algorithm showing our approach.

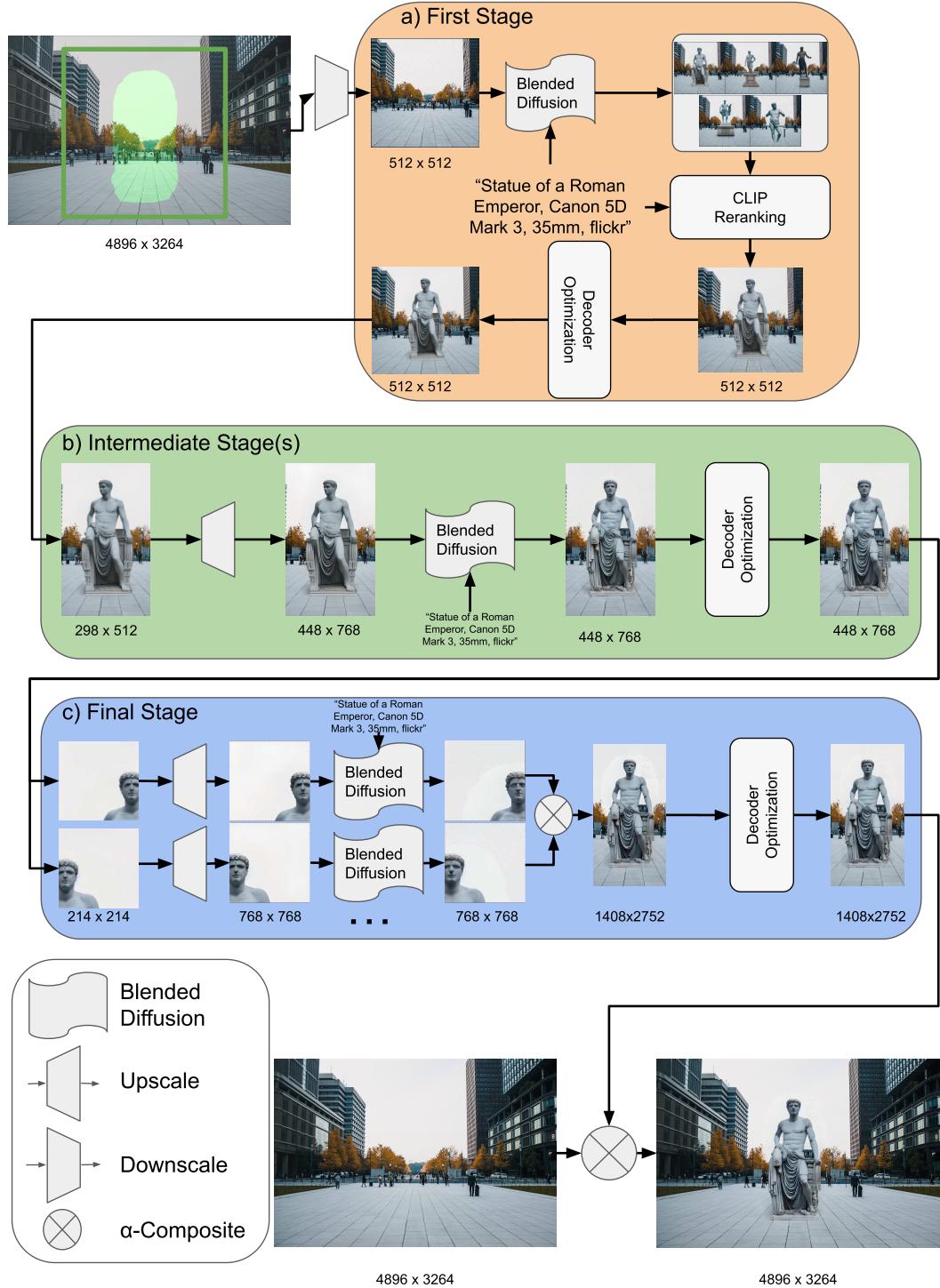


Figure 3: Our approach, with all intermediate results being shown. Best viewed zoomed in.

Algorithm 1: Multi-Stage Blended Diffusion

Input : input image x_{in} , mask m , text prompt c , DDPM p_θ , forward process q , super-resolution model SR, low-pass filter LP, stage count n_s , stage-wise start timestep T'_s , repaint step R , first stage batch size B

```

1  $\tilde{x} \leftarrow \text{Downsample}(\text{Crop}(x_{in}, m))$ 
2  $\mathcal{X}^1 \leftarrow \emptyset$ 
3 for  $i \in [1, \dots, B]$  do
4    $x_T \sim q(x_T | x_0 = \tilde{x})$ 
5   for  $t \in [T, \dots, 0]$  do           // First stage, Blended Diffusion with Repaint
6      $x_{t-1} \sim p_\theta(x_{t-1} | x_t, c)$ 
7      $\hat{x} \sim q(x_{t-1} | x_0 = \tilde{x})$ 
8      $x_{t-1} \leftarrow m \odot x_{t-1} + (1 - m) \odot \hat{x}$ 
9     for  $r \in [1, \dots, R]$  do
10        $x_t \sim q(x_t | x_{t-1} = x_{t-1})$ 
11        $x_{t-1} \sim p_\theta(x_{t-1} | x_t, c)$ 
12        $x_{t-1} \leftarrow m \odot x_{t-1} + (1 - m) \odot \hat{x}$ 
13     end
14   end
15    $\mathcal{X}^1 \leftarrow \mathcal{X}^1 \cup x_0$ 
16 end
17  $x^1 \leftarrow \text{ClipReranking}(\mathcal{X}^1, c)$ 
18  $x^1 \leftarrow \text{DecoderOptimization}(x^1, x_{in}, p_\theta)$ 
19 for  $s \in [2, \dots, n_s - 1]$  do // Intermediate stages, Blended Diffusion without Repaint
20    $\tilde{x} \leftarrow m \odot \text{SR}(x^{s-1}) + (1 - m) \odot \text{LP}(x_{in})$ 
21    $x_{T'_s} \sim q(x_{T'_s} | x_0 = \tilde{x})$ 
22   for  $t \in [T'_s, \dots, 0]$  do
23      $x_{t-1} \sim p_\theta(x_{t-1} | x_t, c)$ 
24      $\hat{x} \sim q(x_{t-1} | x_0 = \tilde{x})$ 
25      $x_{t-1} \leftarrow m \odot x_{t-1} + (1 - m) \odot \hat{x}$ 
26   end
27    $x^s \leftarrow \text{DecoderOptimization}(x_0, x_{in}, p_\theta)$ 
28 end
// Final stage processed in segments, omitted for brevity, see Fig. 3
29  $x_{out} \leftarrow m \odot x^{n_s} + (1 - m) \odot x_{in}$ 

```

- 144 We build our implementation based on the latent diffusion model (LDM) implementation by Rombach
145 et al. [16] and use Stable Diffusion v1.4 [16] for our experiments. However, our approach is also
146 applicable to pixel-space diffusion models. As outlined in the main text, given an input image x_{in} , a
147 mask m , and a text conditioning c , we begin by cropping the image to a square region around the
148 mask. Because the diffusion model only receives the selected area as input, it is important that this
149 area includes the context necessary for consistent inpainting. We, therefore, set a margin around
150 the masked region for each example which ensures the necessary context is included. This area is
151 then downsampled to the resolution the LDM was originally trained on. In the first stage, we use a
152 combination of Blended Diffusion [2] and Repaint [8], as shown in lines 4-14 of Algorithm 1. We
153 generate a batch of B images and select the best one by similarity to the prompt using CLIP ViT/L14
154 [12], following related work [1, 14].
- 155 We then perform multiple stages $s \in [2, \dots, n_s]$ of upscaling, consisting of applying a pre-trained
156 super-resolution model SR and then Blended Diffusion with SDEdit [9]. As SR, we use "realesrgan-
157 x4plus" [19] followed by bilinear downscaling. As the RealESRGAN model was trained to upscale
158 natural images with natural distortions, artifacts caused by the diffusion model in the previous stage
159 usually are present in the output of SR. Therefore, we then use Blended Diffusion to improve the
160 visual fidelity and ensure consistency of the generation with the prompt. As input \tilde{x} for Blended
161 Diffusion, we use the output of SR, but replace the unmasked region with a low-pass filtered version
162 of the input image: $\tilde{x} = m \odot \text{SR}(x^{s-1}) + (1 - m) \odot \text{LP}(x_{in})$, where LP is a low-pass filter and
163 x^{s-1} is the output of the previous stage. We discuss the reason for the filtering in more detail in

164 Section A.4. We implement the low-pass filter LP by bilinearly downsampling the input image x_{in}
165 to the input resolution of the current stage and then upsampling it to the target resolution. After each
166 stage, we perform decoder optimization to improve the blending with the original image.

167 While we could, barring memory constraints, apply the diffusion model to arbitrarily large resolutions,
168 in practice even a generous graphics memory of 32GB limits us to at most 960x960 pixel when
169 using Stable Diffusion. Therefore, for higher resolutions we split the input into multiple overlapping
170 segments and apply the upscaling and Blended Diffusion to them separately. We then combine the
171 upscaled segments by alpha-compositing [11] the overlapping regions. Finally, as we also can not
172 perform decoder optimization at high resolutions directly, we again split the image into multiple
173 segments and optimize the loss L_{DO} across all segments. Using the optimized weights, we then
174 decode each input segment and again use alpha compositing to combine them.

175 **Hyperparameters** In our experiments, we use one intermediate stage with the longer edge of the
176 output being 768 pixels long, followed by the final full resolution which is processed in segments.
177 The segments are each 768x768 pixels large overlapping with each neighbor by 128 pixels. To speed
178 up inference, we use DDIM [18] as sampler with $T = 50$ steps. In the first stage we generate a batch
179 of size $B = 5$ as input for the clip reranking. In the intermediate stages we start the denoising from
180 $T' = 0.4T$ and in the final, segmented stage we start from $T' = 0.25T$. We run decoder optimization
181 for 100 steps, using Adam [7] with learning rate 1×10^{-4} , as suggested in [1]. We use $R = 5$ repaint
182 steps in the first stage and don't use Repaint in the following stages.

183 A.3 Baseline Comparisons

184 To validate our approach, we compare it to two baselines: 1) Using Blended Diffusion with Stable
185 Diffusion v1.4 naively by running the first stage at the highest resolution we can fit in VRAM,
186 960x960 pixels in our case, and then upscaling the edited images bilinearly. 2) Using Dall-E 2 [13]
187 via their Web-UI, we can inpaint the masked region at high resolutions. However, as Dall-E 2 is
188 limited to 1024x1024 pixels, we have to inpaint in multiple segments separately. The results are
189 shown in Figure 4. We find that directly applying Blended Diffusion to a resolution significantly
190 higher than the model was trained on leads to repeated elements, such as the two mountains or two
191 paintings in our examples. Similarly, while producing a high-fidelity image, applying Dall-E 2 in
192 segments also produced repeated elements (statues, paintings). Our multi-stage approach allows us to
193 obtain a globally consistent result.

194 A.4 Ablations

195 **Upscaling** To upscale the image or image segments after the first stage, we upscale the low-
196 resolution input with Real-ESRGAN and then use text-guided Blended Diffusion with a low-pass
197 filtered background image. To validate our design decisions, we compare a) only bilinear upscaling,
198 b) only ESRGAN upscaling, c) ESRGAN + unconditional diffusion, d) ESRGAN + text-conditional
199 diffusion, e) ESRGAN + text-conditional Blended Diffusion with an unfiltered background image
200 f) ESRGAN + text-conditional Blended Diffusion with a low-pass filtered background image. We
201 show the results for the different ablations in Figure 5, showing that using our proposed method
202 significantly improves blending with the background and the visual fidelity of the results.

203 We also note that if we do not low-pass filter the background, Blended Diffusion tends to produce
204 more blurred outputs. We hypothesize that this is due to the difference in "sharpness" between the
205 high-resolution background image and the upscaled output of the previous stage. This difference then
206 seems to be amplified in the denoising process, leading to a more blurry edited region. We stress that
207 this is an empirical observation and more work is needed to properly understand this phenomenon.

208 **Repaint** As we use Repaint [8] in the first stage of our approach, we provide an ablation of the
209 effect of different repaint steps R on the final result. We do not use clip-reranking in this experiment
210 in order to isolate the effect of Repaint. The results are shown in Figure 6. While not using repaint
211 often leads to poor blending of the edited area, such as in the corgi, painting, and hair examples,
212 using a too high number of repainting steps can also lead to poor results (hair, corgi). Unlike the
213 other samples, the river scene was generated by Stable Diffusion and editing it with Stable Diffusion
214 seems to work well even without repaint. In practice we find that setting the number of repaint steps

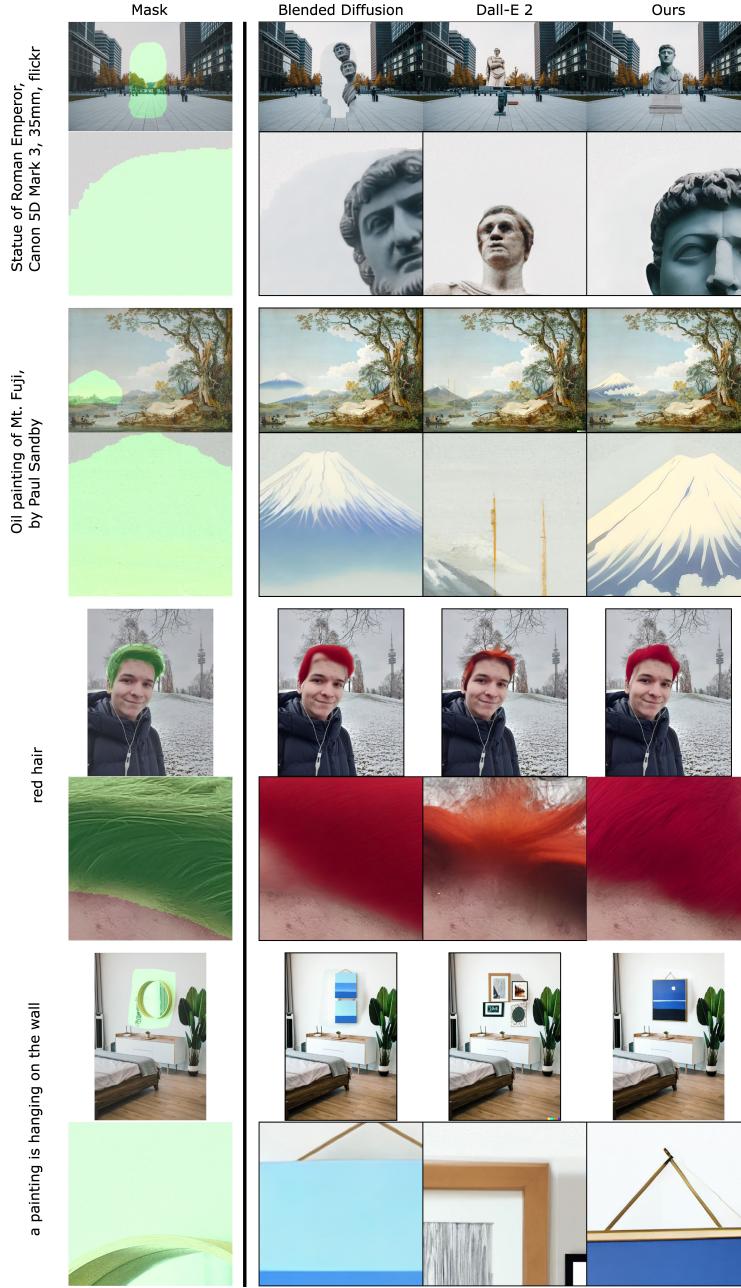


Figure 4: Comparison of our approach to two baselines. Left: we directly apply Blended Diffusion to the highest resolution we can fit into the VRAM (960x960 pixels) and then bilinearly upscale the output. Middle: We use the Dall-E 2 web UI to edit the image at its full resolution. Due to the edited region being larger than the 1024x1024 generation window, we have to apply Dall-E to multiple independent segments. Right: Our proposed approach. We find that directly applying Blended Diffusion leads to repeated elements (two heads, two mountains, two pictures) and fails to produce fine details (hair). Dall-E 2 produces visually high-fidelity images, but fails to produce globally coherent images (floating statues, four paintings). Our method produces globally consistent images while providing a similar visual fidelity. Note that the full images are downsampled. The zoomed-in regions measure 512x512 pixels and are shown at full resolution.

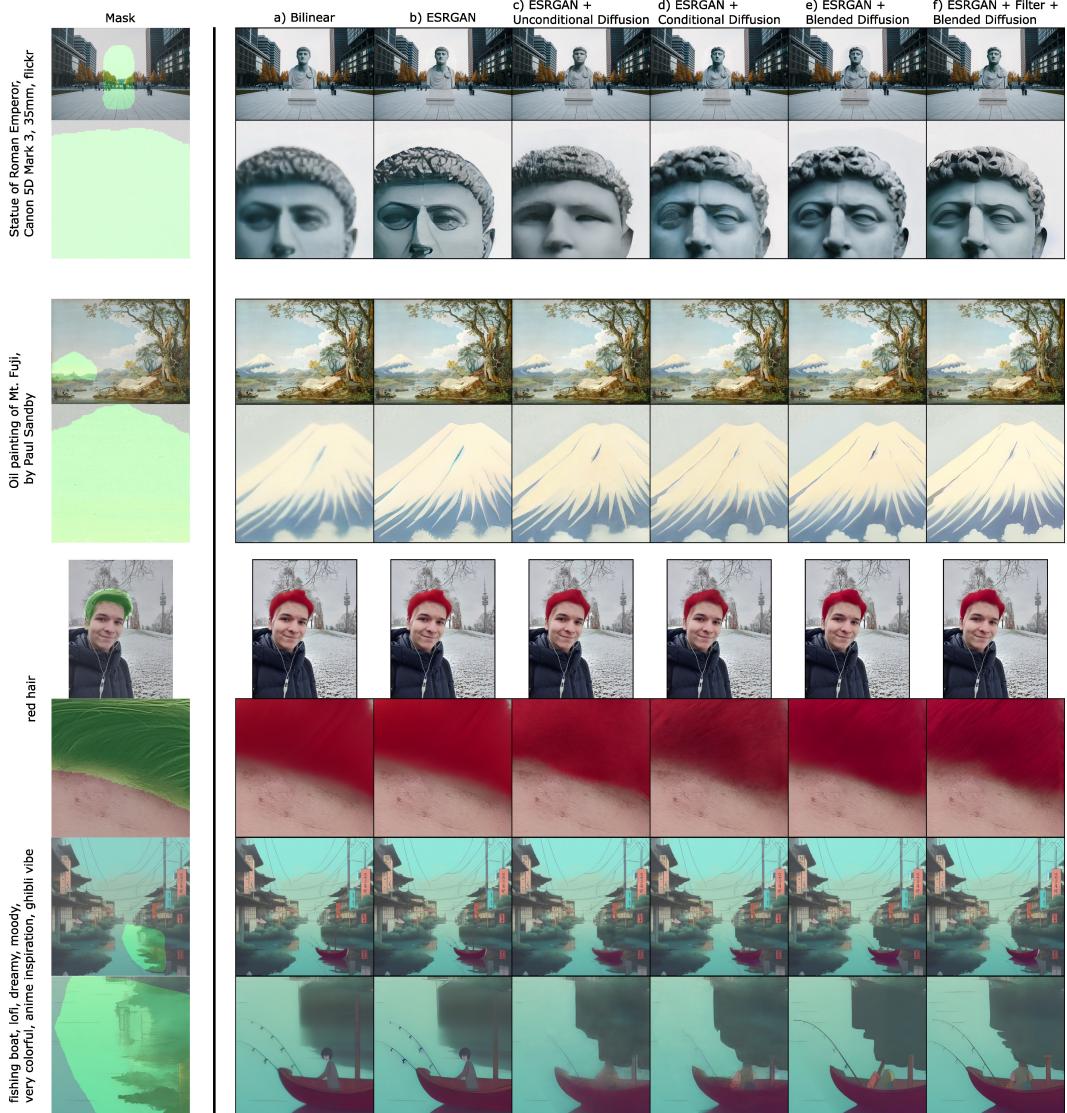


Figure 5: Ablation of different upscaling methods, applied after the Blended Diffusion in the first stage with fixed seeds. a) Bilinear upscaling, b) ESRGAN, c) ESRGAN + unconditional diffusion, d) ESRGAN + text-conditional diffusion, e) ESRGAN + text-conditional Blended Diffusion, f) ESRGAN + text-conditional Blended Diffusion with a low-pass filtered background. Note that the images are downsampled from the full resolution.

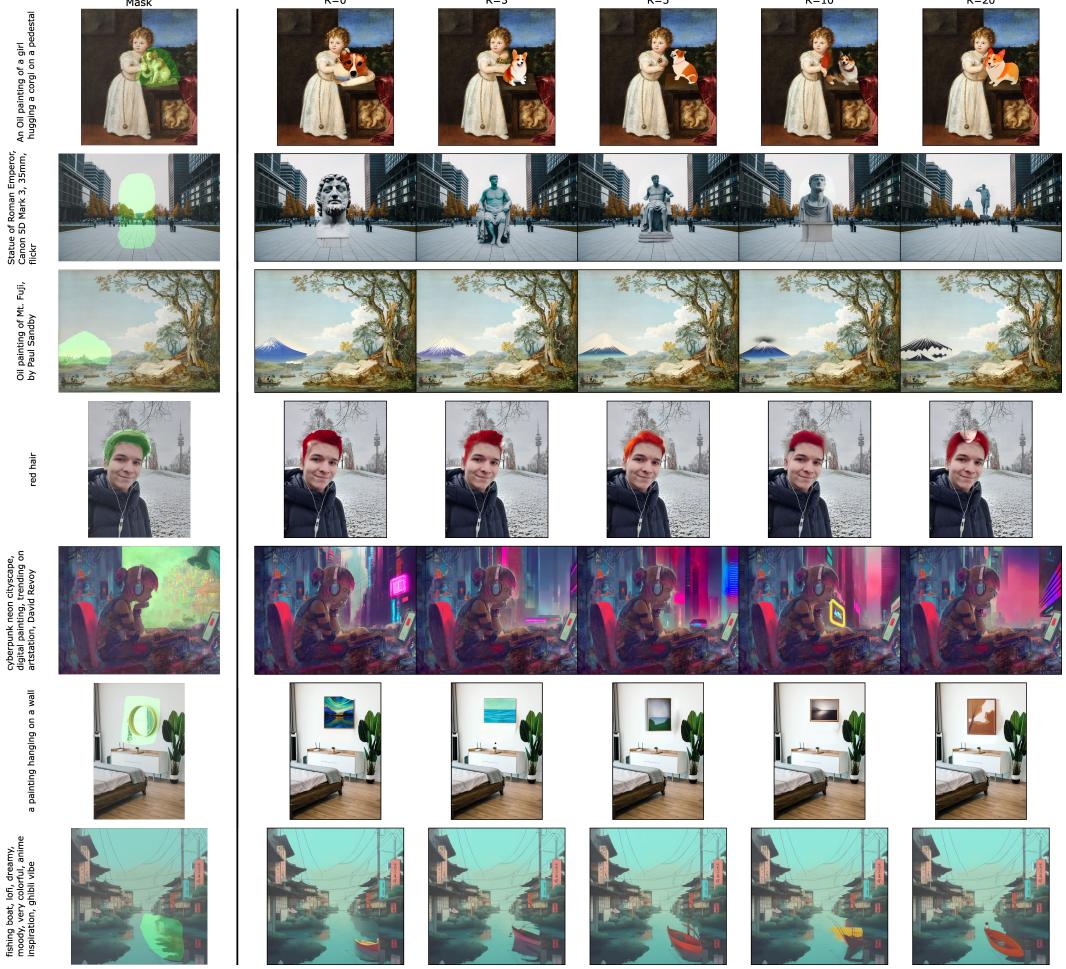


Figure 6: Comparison of different repaint steps R without clip reranking. Without repaint, the edited region often does not blend well with the background, while too many repaint steps can also cause unrealistic images (corgi, hair). We generally find $R = 5$ to work best. Note that the images are downsampled from the full resolution.

215 to $R = 5$ works well in most cases, however, a practitioner may want to adjust this parameter on a
216 case-by-case basis to obtain the best possible results.

217 **Random Seeds** To show the variation of our approach to random seeds we show results for different
218 inputs for 8 different seeds each in Figure 7. The quality of the results does vary between seeds, but
219 for a practitioner it should be easy to run the model multiple times and pick the most appropriate
220 output.

221 A.5 Related Work

222 Multi-stage approaches have been previously proposed for image generation [4, 15]. Cascaded
223 Diffusion [6] trains one model to generate images at a low resolution and then two super-resolution
224 networks which are conditioned on the low-resolution input through channel-wise concatenation. A
225 similar architecture is also utilized in Imagen [17] and DALL-E 2 [13]. Unlike these methods we use
226 the same model in all stages and do not explicitly condition on the low-resolution image, allowing us
227 to use any pre-trained text-conditional model. Additionally, these methods are in practice still limited
228 to 1024x1024 pixels by computational constraints. Wu et al. [20] present an approach that is able to
229 generate arbitrarily high-resolution images and perform high-resolution outpainting by learning a
230 representation of the nearby context, allowing global consistency. However, their work is focused on

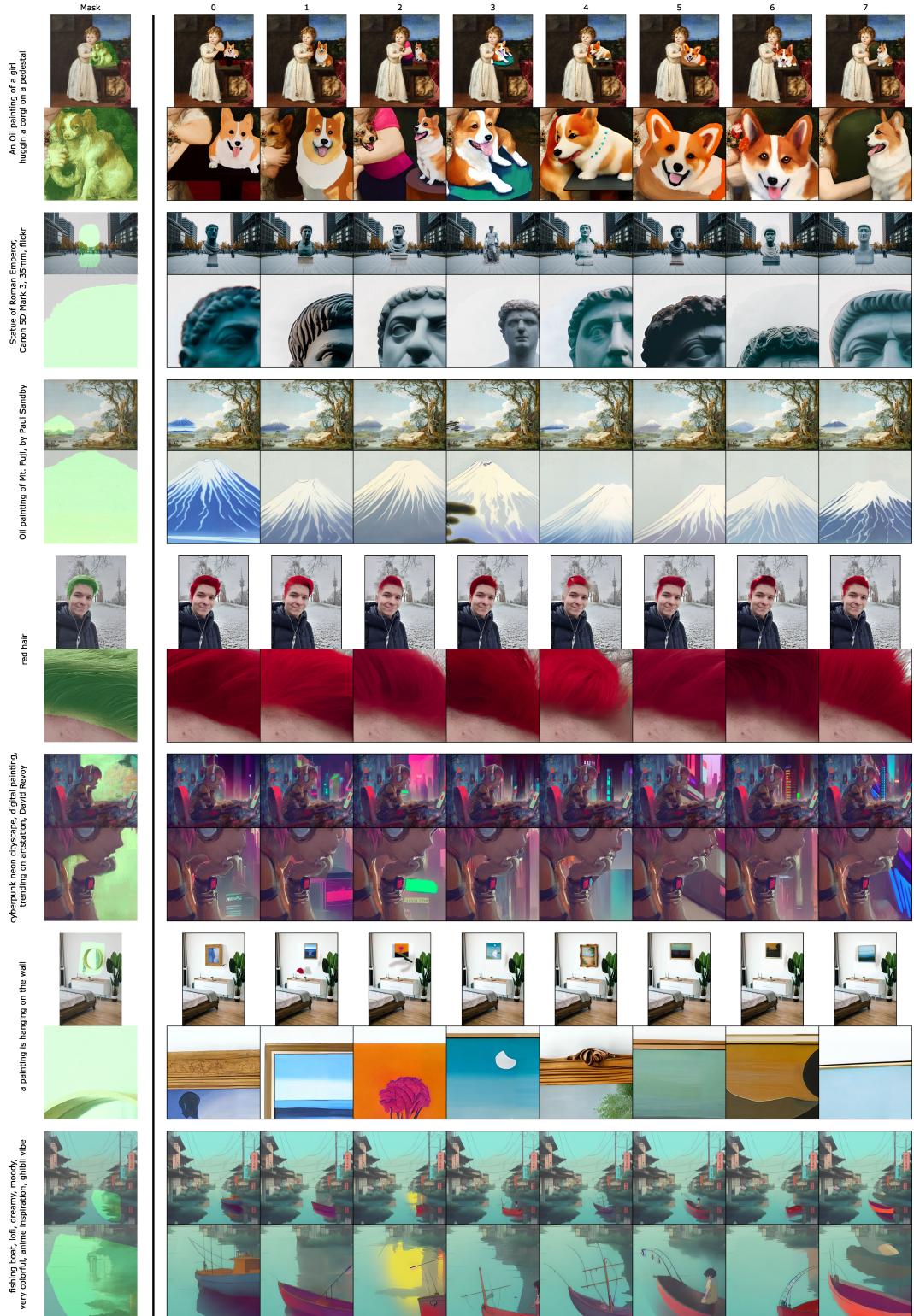


Figure 7: Outputs of our model for different random seeds. Note that the images are downsampled from the full resolution.

231 out-painting or image generation and they do not investigate text-guided editing. Glide [10] presents
232 an approach to generate and edit images using a text-conditional diffusion model. Similar to [2], their
233 approach is limited to the resolution that their diffusion model was trained on, and can not be easily
234 extended to higher resolutions.

235 **B Limitations**

236 Our current approach still has limitations, in three main ways: Firstly, we reuse pretrained diffusion
237 models and are therefore limited to editing images that they perform well at generating. For example,
238 we found stable diffusion to perform poorly on night-scenes and in general it seems to not perform as
239 well on photographs as it does on artwork. Secondly, the blending between the unedited region and
240 edited region is often still visible upon close inspection. Decoder optimization works well to improve
241 the blending, but once we reach resolutions around 1000x2000 or above the optimization starts to
242 take a significant amount of time (10 minutes on a V100 for 3000x2000 pixels). Using a pixel-space
243 diffusion model instead would remove the need for decoder optimization. Thirdly, our approach uses
244 Real-ESRGAN to upscale the images. While this works well in most cases, in some cases fine details
245 are removed, leading to poor blending with the original image.

246 **C Images used**

247 All images used in our paper are used in accordance with their licenses and attributed below, in the
248 order of Figure 7:

- 249 1. "Portrait of Clarissa Strozzi", by Titian Vecelli, via
250 https://commons.wikimedia.org/wiki/File:Clarissa_Strozzi,_por_Tiziano.jpg (1803x2117)
- 251 2. "people walking on sidewalk near high rise buildings during daytime", by Nat Weearwong,
252 via <https://unsplash.com/photos/0cZgvYHirBg> (4896x3264)
- 253 3. "The River Severn at Shrewsbury, Shropshire", by Paul Sandby, via
254 <https://unsplash.com/photos/HEEvYhNzpEo> (3999x3041)
- 255 4. Selfie by author (3456x4608)
- 256 5. "white wooden dresser with mirror photo", by Minh Pham, via
257 https://unsplash.com/photos/7pCFUybP_P8 (3902x5853)
- 258 6. "Lofi Cyberpunk" by David Revoy <https://www.davidrevoy.com/article867/lofi-cyberpunk>
259 (2431x1930)
- 260 7. Anime-style image of river generated with stable diffusion by the authors (2048x2048)