

TEXT-GUIDED MASK-FREE LOCAL IMAGE RETOUCHING

Zerun Liu¹, Fan Zhang¹, Jingxuan He², Jin Wang³, Zhangye Wang², Lechao Cheng^{3*}

Zhejiang University of Technology¹,
Zhejiang University²,
Zhejiang Lab³

ABSTRACT

In the realm of multi-modality, text-guided image retouching techniques emerged with the advent of deep learning. Most currently available text-guided methods, however, rely on object-level supervision to confine the region of interest that may be updated. This not only makes it more challenging to develop these algorithms but also limits how widely deep learning can be used for image retouching. In this paper, we offer a text-guided mask-free image retouching approach that yields consistent results to address this concern. Specifically, we propose a two-stage mask-free training paradigm tailored for text-guided image retouching tasks. In the first stage, a unified mask is proposed according to the query description, and then several candidate images are generated with the provided mask and the conditional description based on diffusion model. Extensive experiments have shown that our method can produce high-quality images based on spoken language.

Index Terms— text guided, mask free, image retouching

1. INTRODUCTION

With the advent of the internet and the rise in popularity of numerous photography tools, it is now easier than ever before to obtain enormous photos. In recent years, image retouching has attracted a significant deal of attention due to the need of modifying photographs to accommodate a variety of scenarios. Creating satisfying photographs using conventional techniques, however, often necessitates a substantial amount of talent and labor. Therefore, it would be convenient enough if we could modify photographs with language-specific descriptions. In this work, we focus on retouching images with language guidance.

Approaches for text-guided image retouching tasks can be easily classified into two categories: mask-based and mask-free methods, depending on whether the mask is used. Paint By Word [1] is the first zero-shot solution to solve local image retouching. Following this, several diffusion-based models have manifested themselves with increasing attention, such as Blended Diffusion [2], Glide [3], and Blended Latent Dif-

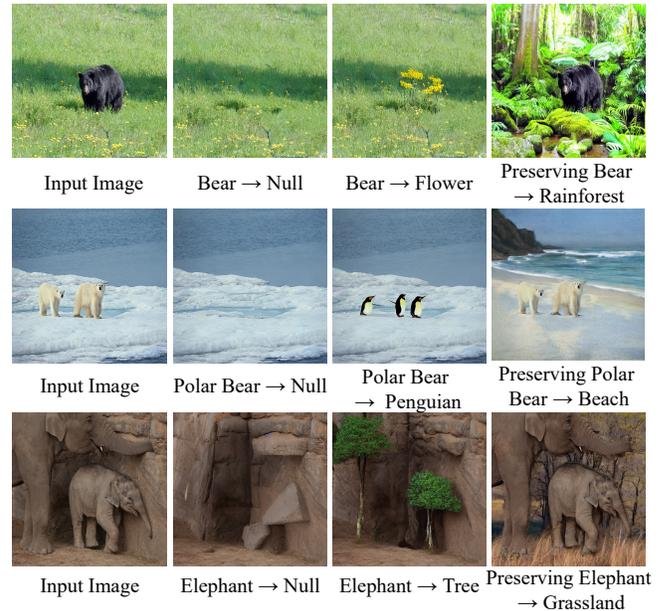


Fig. 1. Applications of our method: from left to right are object removal, object replacement, and background replacement.

fusion [4]. However, it can be challenging to create an appropriate mask for each image that requires manipulation, particularly when the shape of the mask is complex. Therefore, researchers attempt to pay more attention to mask-free techniques, which are more practical. Mask-free methods can address local image retouching without mask guidance, like VQGAN-CLIP [5], CLIP-Guided Diffusion [6], Diffusion-CLIP [7], and DiffuseIT [8]. However, these algorithms often suffer from the drawbacks of skewed object recognition and incomplete preservation of the original backdrop.

To cope with the above issues, we propose a two-stage framework named text-guided mask-free local image retouching, as shown in Fig 1, which converts a specific part of the image that matches the query to a target object described in the text. In the first stage, our method will propose an object mask according to the query. In the second step, given the input image, the suggested mask, and the text description, our

*Corresponding author.

algorithm will generate several candidate images, which will then be assessed to yield the final retouched image. Compared with past approaches, our method can perform local image editing without masking regions of interest. In addition, we also present the multi-modality quality assessment strategy for quantifying the quality of model outputs by picking the best one among candidate ones based on cross-model alignment and image quality assessment.

We summarize the contributions as follows:

- We propose a two-stage mask-free training paradigm tailored for text-guided image retouching tasks. In the first stage, an object mask is proposed according to the query, and several candidate images are generated based on the proposed mask and the text description. Candidate images are then evaluated to produce the final retouched image.
- A location-aware refinement module is introduced in our framework, where users can select one entity from all entities by describing it in the text. Moreover, we also propose a multi-modal quality assessment module in which several candidate images can be assessed to help produce high-quality and semantically consistent retouched images.
- We conducted extensive experiments on ITMSCOCO and ITFlickr which yield promising visual results.

2. RELATED WORK

2.1. Text-guided Image Retouching

The pioneering work [9] introduced a deep generative model that can generate images from natural language descriptions. Later, another work [10] handled image retouching by developing an effective GAN architecture conditioned on text embeddings. Recent advances are fueled by autoregressive transformer models that utilize VQ-VAE [11] to alleviate the problem of an unaffordable amount of computation for large models. Diffusion models have been proven to be able to generate high-quality images, especially when natural language descriptions are given for diversity. Recent works [12, 3] explored diffusion models for text-guided image retouching and achieved promising results with efficient computations. However, these approaches mainly specialize in global image retouching, ignoring the manipulation of a user-defined region of an image, which is a common case in our daily life.

Region-based text-guided image retouching focuses on modifying a specific region of an image while preserving the rest of it. Paint By Word [1] generates a realistic painting of the synthesized image based on user-provided masks and words. Furthermore, Blended Diffusion [2] was proposed for region-based editing of generic natural images instead of synthesized ones, and language descriptions are not restricted to

a specific domain. However, the methods mentioned above require masks to guide where to edit for region-based image retouching. To tackle this problem, our framework utilizes a segmentation model to generate masks based on queries that are more convenient for users to provide.

2.2. Diffusion Models

Diffusion probabilistic models [13] have recently been shown to produce high-quality images while offering sufficient image generation diversity. Ho et al. [14] firstly proposed to synthesize images based on diffusion models. To take more control of the generation process, SDEdit [12] was presented that keeps low-frequency information by terminating the diffusion process early before it turns into a Gaussian distribution. To further improve the quality of generated images, Dhariwal et al. [15] presented a trade-off between quality and diversity by adding attention pooling as a classifier. It beats some famous GAN-based image-generation methods and shows the powerful generative ability of diffusion models. At the same time, multi-modal vision and language models trained on several large-scale datasets, such as CILP, promote the development of diffusion models by adding auxiliary information to the denoising process. For example, Blended Diffusion [2] utilizes CLIP to modify parts of the image based on language and mask guidance. Different from Blended Diffusion, DiffusionCLIP modifies the diffusion model to match up CLIP by introducing a new loss function, which improves the diversity of the generated effect and relieves the performance of the model collapse. Compared with classifier guidance and CLIP guidance, classifier-free guidance [16] can acquire more accurate and better results because it can indicate an unconditional gradient estimation model and a conditional gradient estimation model at the same time in the same model. Benefiting from this, several huge diffusion models like Glide [3], DALLE2 [17], Imagen [18] are proposed and make amazing results.

3. METHOD

3.1. Overview

Let $I \in \mathbb{R}^{H \times W \times 3}$ be the input image, \mathbf{q} be the query description and \mathbf{v} be the conditional text, respectively. The goal of the proposed framework attempts to retouch I at the local region of interest based on \mathbf{q} and \mathbf{v} to produce a high-quality natural image.

To achieve this, we draw inspiration from entity segmentation [19] to segment all visual entities $I_i, i \in \{0, \dots, n-1\}$ in an image based on mask potentials M_i without considering semantic category labels. Each visual entity I_i is generated as follows:

$$I_i = I \odot M_i \quad (1)$$

Where \odot denotes an element-wise multiplication operation.

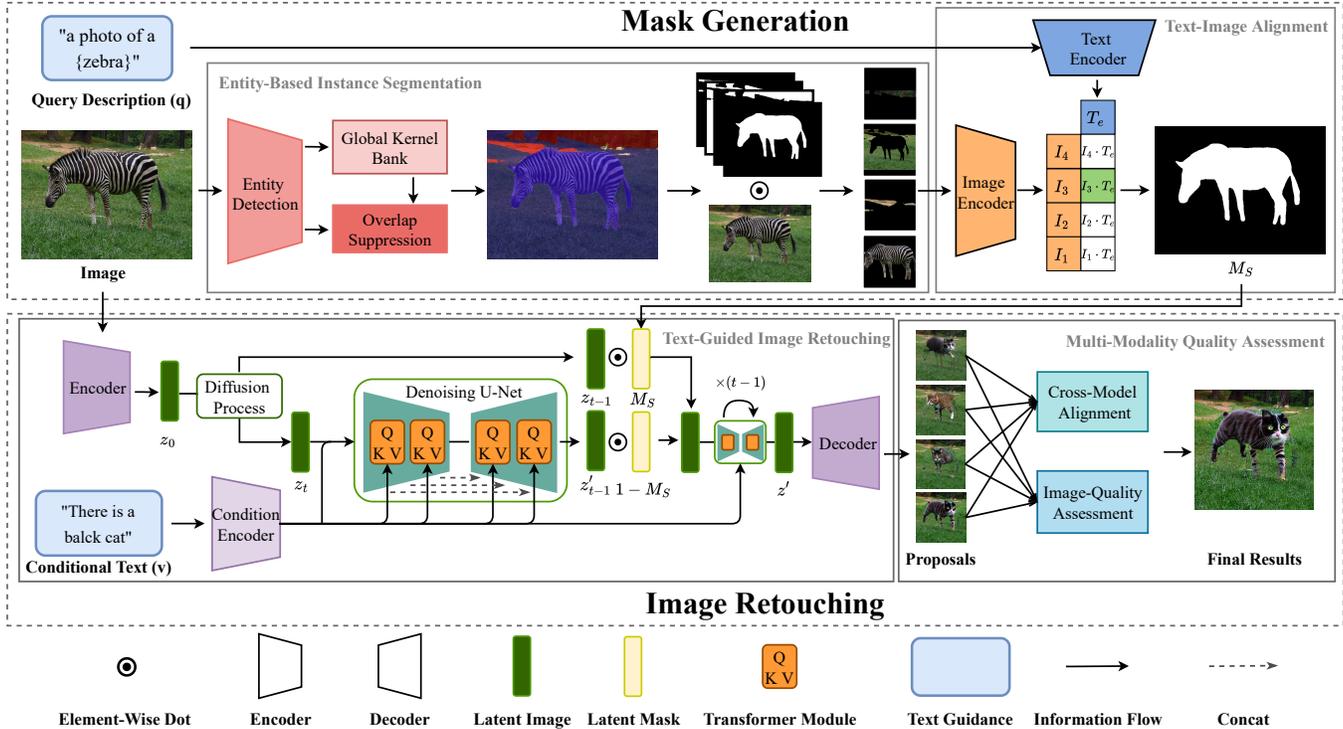


Fig. 2. The overview of our framework. It comprises two parts: image retouching and mask generation. Based on the query, mask generation tries to create the mask based on the query description. Image retouching focuses on altering the image based on the provided conditional text and picking the best of several proposals via the multi-modality quality assessment module.

Then, we adopt pre-trained CLIP models to predict the correct pairings of a batch of (image, text) examples. Assuming that \mathbf{E}_T and \mathbf{E}_I are text encoder and image encoder, respectively. For each entry in multi-modal embedding space spanning by visual entities and query descriptions, we compute the correlation C_i as:

$$C_i = \mathbf{E}_I(I_i) \cdot \mathbf{E}_T(\mathbf{q}) \quad (2)$$

Notably that we consider all query words as one instance to form single embedding $T_e = \mathbf{E}_T(\mathbf{q})$, as illustrated in Fig 2. To determine the regions that ought to be edited, the mask M_S with high correlation confidence is selected. Formally, $M_S = \bigcup_i M_i$, *st.* $C_i \geq \tau$. In practice, fixing threshold τ does not work well in real sceneries, and we introduce an adaptive thresholding strategy inspired by [20]. More details can be referred at Sec 3.2.2.

Once we obtain the accurate regions of interest, the input image I , and the conditional text \mathbf{v} as well as detected areas M_S are fed into a diffusion model to generate retouched proposals \mathbf{P}_k , $k \in \{0, \dots, m-1\}$. m is the number of proposals. We finally attach a multi-modality quality assessment module to produce high-quality and semantic-aware results based on proposed Cross-Model Alignment (CMA) and Image-Quality Assessment (IQA) strategies. The subsequent sections will

detail the above components.

3.2. Mask Generation

Unlike existing popular local text-guided image retouching techniques, our model will construct the mask based on the query to determine where to modify. In particular, it can generate accurate masks in relatively simple scenarios, however, it fails when it comes to more complex situations. In this section, we attempt to generate more precise masks based on the query description. We will outline them in detail below.

3.2.1. Entity-Based Instance Segmentation

The entity-based instance segmentation module is built upon ES [19], a novel open-world entity segmentation task to investigate the feasibility of convolutional center-based representation to segment things and stuff in a unified manner. The introduced global kernel bank and overlap suppression aim to exploit the requirements of ES in order to improve the segmentation quality of predicted entity masks. As a result, the mask potentials M_i and visual entities I_i are obtained for subsequent processing.

3.2.2. Text-Image Alignment

We leverage the CLIP [21] strategy to model cross-modality relations for query description \mathbf{q} and visual entities I_i . As stated in eq 3.1, τ plays an important role in obtaining precise regions. Thus, we adopt an adaptive thresholding strategy based on sampling over cumulative correlation confidence. Specifically, we first sort the correlation values in $\mathcal{C}_i, i \in \{0, \dots, n-1\}$ from high to low and calculate the cumulative confidence. Then the adaptive threshold τ is obtained based on the inverse transform sampling over the maximal contribution step.

3.2.3. Location-Aware Refinement

The above text-image alignment procedure can work well in most cases, but it may fail when descriptions contain location constraints. To tackle this issue, we propose a location-aware refinement module that amends the mask region as expected. We simply split the image space into nine regions evenly and bag each detected visual entity into the corresponding location. Recall that we separate the regions and background based on the correlation confidence \mathcal{C} , and we proceed to perform a location-aware refinement operation by restricting the orientation.

3.3. Text-Guided Image Retouching

After obtaining regions \mathbf{M}_S correlated with query description \mathbf{q} , we further retouch the image I with conditional text \mathbf{v} with provided mask \mathbf{M}_S . The input image I and the conditional text \mathbf{v} will be compressed into the same latent space, respectively. Formally, denoting z_0 as the initial latent vector, we utilize the denoising network $\epsilon_\theta(z_t, t, \mathbf{v})$ to reconstruct z from the Gaussian distribution $\mathcal{N}(0, I)$. Thus, we can simplify the corresponding objective as:

$$\mathcal{L} := \mathbb{E}_{z, \epsilon \sim \mathcal{N}(0, I), t} [\|\epsilon - \epsilon_\theta(z_t, t, \mathbf{v})\|] \quad (3)$$

For timestep t , the denoising result z'_t is further updated with mask \mathbf{M}_S as:

$$z'_t = z_t \odot (1 - \mathbf{M}_S) + z'_t \odot \mathbf{M}_S \quad (4)$$

Eq 4 indicates that we keep the original background while learning text-guided visual entities based on the conditional text \mathbf{v} . The last z'_t is put forward into a decoder network to generate plausible proposals $\mathbf{P}_k, k \in \{0, \dots, m-1\}$.

3.4. Multi-Modality Quality Assessment

Randomization is an essential factor of the denoising procedure, which assures the variety of the diffusion model, yet it may also lead to undesirable outcomes. To address this issue, we present a multi-modality quality assessment module for

selecting the best result from several candidates. In this module, we examine the quality of the output from two perspectives: cross-model assessment and image quality assessment.

3.4.1. Cross-Model Alignment

In practice, we employ CLIP to measure the consistency between proposals \mathbf{P}_k and the conditional text \mathbf{v} as described in section 3.2.2. The cross-model confidence \mathcal{C}'_k for proposal \mathbf{P}_k is then defined as:

$$\mathcal{C}'_k = \mathbf{E}_I(\mathbf{P}_k) \cdot \mathbf{E}_T(\mathbf{v}) \quad (5)$$

Noted that \mathcal{C}'_k is normalized to $[0, 1]$. Recall that \mathbf{E}_T and \mathbf{E}_I are text encoder and image encoder before.

3.4.2. Image Quality Assessment

In addition to cross-modality alignment, we shall also pay attention to the visual quality of the generated result. Therefore, we introduce an image quality assessment module to enable high-quality output images. Formally, we define the image quality confidence score \mathcal{S}_k as follows:

$$\mathcal{S}_k = \frac{1}{H * W} \sum_i \sum_j \|I(i, j) - \mathbf{P}_k(i, j)\| \quad (6)$$

Finally, we choose the best suitable proposal $\mathbf{P}_{k'}$ that fulfills:

$$k' := \arg \max_k (\mathcal{C}'_k - \alpha \cdot \mathcal{S}_k) \quad (7)$$

Where α is the hyperparameter, and we empirically set $\alpha = 5.0$. Experiments indicate they can increase the overall quality of image retouching results with semantic consistency.

4. EXPERIMENT

4.1. Experimental Settings

We follow [22] to generate ITMSCOCO and ITFlickr by selecting 5000 text-image pairs from MSCOCO [23] and Fliker30K [24], respectively. Furthermore, SSIM, PSNR, FID, MSE, and LPIPS are utilized as metrics to quantify the quality of retouching results.

4.2. Compared to Existing Methods

To evaluate the performance of our method, we compare our results against the following state-of-the-art baseline models: Glide, Blended Latent Diffusion, VQGAN-CLIP, CLIP-Guided Diffusion, and DiffuseIT. Experiments show that our model outperforms the above methods qualitatively and quantitatively.

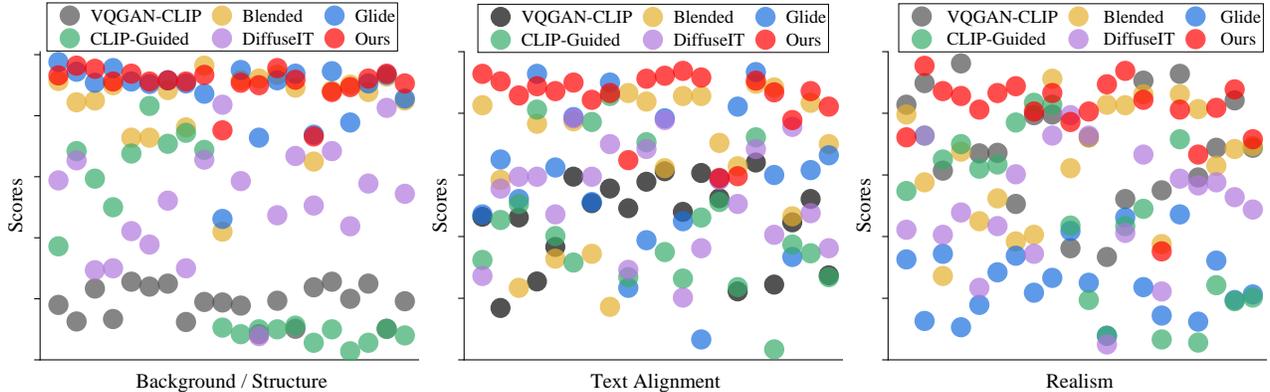


Fig. 3. After each image was assessed by 20 participants, we computed the mean scores of **Background / Structure**, **Text Alignment** and **Realism** for 20 generated images.

4.2.1. Qualitative Comparison

As shown in Fig 7, we compare our results with two mask-based methods and three mask-free methods. The masks for the mask-based methods are generated by our method, which shows our superiority that we needn't provide the mask for the regions of interest. As for mask-free methods, the main drawback of them is that they couldn't preserve the whole background well, which may lead to the failure of retouching, besides, these methods may fail when handling complex situations like Column 7 of Fig 7. As a result, our method brings the best of both worlds and obtains persuasive results.

4.2.2. Quantitative Comparison

In the realm of text-guided image retouching, it's difficult to measure the performance of models by metrics quantitatively. Therefore, we introduce a user study, as depicted in Table 1 and Fig. 3. The study participants were instructed to evaluate each result based on (1) the preservation of the background and structure relative to the source image, (2) alignment with the corresponding text, and (3) the realism of the output. A total of 20 participants took part in the study and evaluated 20 randomly selected examples for each method. The examples were presented in random order. For each example, the participants rated the quality on a scale of 0 to 5, with higher ratings indicating better performance.

Table 1. The quantitative statistical rating of generated results. \uparrow means that the higher the value, the better. **Red** indicates the best scores, while **blue** indicates the second-best results.

Methods	Background / Structure (\uparrow)	Text Alignment (\uparrow)	Realism (\uparrow)
Glide [3]	4.36 \pm 0.59	2.97 \pm 1.16	3.47 \pm 0.9
Blended [4]	4.19 \pm 0.64	3.29 \pm 1.16	3.23 \pm 0.95
VQGAN [5]	0.93 \pm 0.28	2.25 \pm 0.73	1.27 \pm 0.58
Clip-guided [6]	1.68 \pm 1.5	2.25 \pm 1.11	2.35 \pm 1.35
DiffuseIT [8]	2.58 \pm 0.94	2.66 \pm 0.87	2.5 \pm 0.96
Ours	4.52\pm0.3	4.21\pm0.53	4.03\pm0.65

4.3. Ablation Studies

In this section, we conduct ablation studies to verify the effectiveness of the mask generation and image retouching modules.

4.3.1. Mask generation

Adaptive Threshold As illustrated in Fig 4, the empirically rough threshold often cause the loss of visual entities when complicated scenarios are encountered. While our adaptive threshold exhibits much more characteristics that capture the regions of interest well.

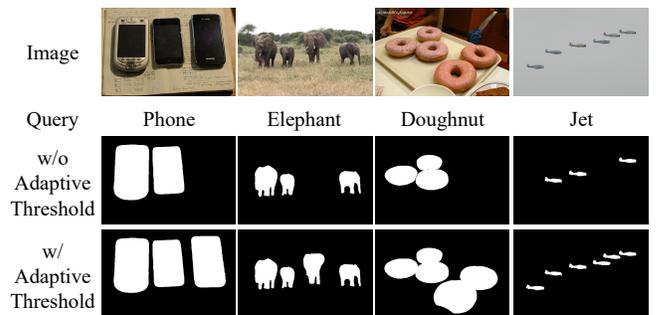


Fig. 4. Adaptive threshold. We empirically set the threshold to τ to 0.2 for the baseline.

Location-Aware Refinement We demonstrate the effectiveness of the proposed location-aware refinement in Fig 5 and most common positional restraints are well retained in real cases.

4.3.2. Multi-Modality Quality Assessment

CMA vs. IQA.

As can be seen in Tab 2, we evaluate the CMA and IQA modules with the mean values of all validation images. The quantitative results just reveal that the IQA component plays an essential role in the production of high-quality images,

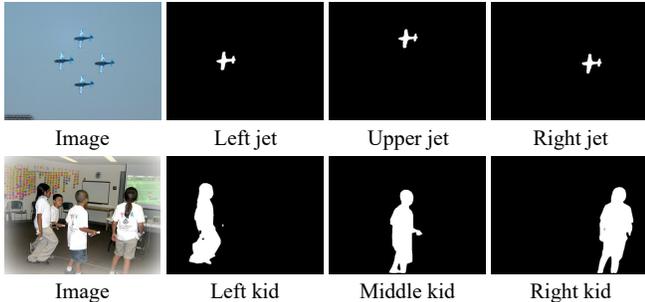


Fig. 5. Visual entities with orientation.

while the qualitative visual results in Fig 6 show us the significant involvement of CMA to yield semantic-aware targets.

Table 2. CMA vs. IQA. We carry out comparative experiments on the ITMSCOCO datasets. The max timestep is 200, and the proposal number m equals 4.

CMA	IQA	SSIM(\uparrow)	PSNR(\uparrow)	FID(\downarrow)	MSE(\downarrow)	LPIPS(\downarrow)
×	×	0.7042	17.7389	16.7268	0.0218	0.2673
√	×	0.7121	17.7823	14.5790	0.0191	0.2608
×	√	0.7218	18.6803	12.7913	0.0160	0.2499
√	√	0.7125	17.8141	14.5512	0.0190	0.2605

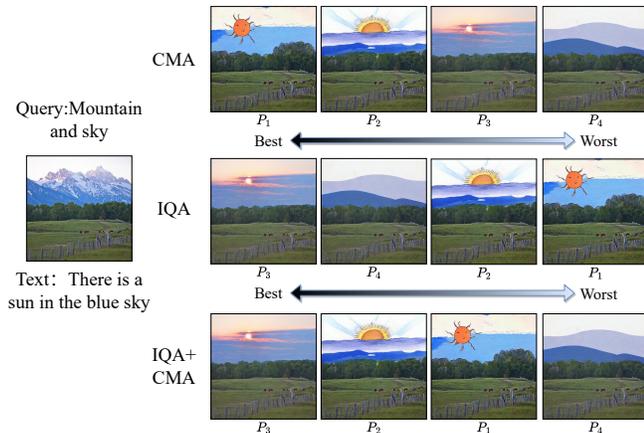


Fig. 6. The effectiveness of CMA and IQA. P_1, P_2, P_3, P_4 are ranked only based on CMA strategy, while P_3, P_4, P_2, P_1 are ranked by the IQA module. After combining the CMA and IQA module, the final result changed to P_3, P_2, P_1, P_4 .

5. CONCLUSION

This study offers a two-stage framework for text-guided mask-free local image retouching. This framework transforms a specified region of the image that fits the query into a text-described target object. Extensive experiments demon-

strate the superiority of the proposed approach against existing methods.

Acknowledgments

This work is partially supported by the National Natural Science Foundation of China (Grant No. 62106235), by the Exploratory Research Project of Zhejiang Lab(2022PG0AN01), by the Zhejiang Provincial Natural Science Foundation of China (LQ21F020003).

6. REFERENCES

- [1] David Bau, Alex Andonian, Audrey Cui, YeonHwan Park, Ali Jahanian, Aude Oliva, and Antonio Torralba, “Paint by word,” *arXiv preprint arXiv:2103.10951*, 2021.
- [2] Omri Avrahami, Dani Lischinski, and Ohad Fried, “Blended diffusion for text-driven editing of natural images,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18208–18218.
- [3] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen, “Glide: Towards photorealistic image generation and editing with text-guided diffusion models,” *arXiv preprint arXiv:2112.10741*, 2021.
- [4] Omri Avrahami, Ohad Fried, and Dani Lischinski, “Blended latent diffusion,” *arXiv preprint arXiv:2206.02779*, 2022.
- [5] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff, “Vqgan-clip: Open domain image generation and editing with natural language guidance,” in *European Conference on Computer Vision*. Springer, 2022, pp. 88–105.
- [6] Katherine Crowson, “Clip-guided diffusion,” <https://github.com/afiaka87/clip-guided-diffusion>.
- [7] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye, “Diffusionclip: Text-guided diffusion models for robust image manipulation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2426–2435.
- [8] Gihyun Kwon and Jong Chul Ye, “Diffusion-based image translation using disentangled style and content representation,” *arXiv preprint arXiv:2209.15264*, 2022.
- [9] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov, “Generating images from captions with attention,” *arXiv preprint arXiv:1511.02793*, 2015.
- [10] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee, “Gener-



Fig. 7. Qualitative comparison of text-guided image retouching on MSCOCO dataset. We choose two mask-based methods and three mask-free methods. Our approach generates realistic samples based on the conditional text with better perceptual quality than the baselines.

ative adversarial text to image synthesis,” in *International conference on machine learning*. PMLR, 2016, pp. 1060–1069.

- [11] Aaron Van Den Oord, Oriol Vinyals, et al., “Neural discrete representation learning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [12] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon, “Sdedit: Guided image synthesis and editing with stochastic differential equations,” in *International Conference on Learning Representations*, 2021.

- [13] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” in *International Conference on Machine Learning*. PMLR, 2015, pp. 2256–2265.
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel, “Denoising diffusion probabilistic models,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [15] Prafulla Dhariwal and Alexander Nichol, “Diffusion models beat gans on image synthesis,” *Advances in Neu-*

ral Information Processing Systems, vol. 34, pp. 8780–8794, 2021.

- [16] Jonathan Ho and Tim Salimans, “Classifier-free diffusion guidance,” *arXiv preprint arXiv:2207.12598*, 2022.
- [17] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen, “Hierarchical text-conditional image generation with clip latents,” *arXiv preprint arXiv:2204.06125*, 2022.
- [18] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al., “Photorealistic text-to-image diffusion models with deep language understanding,” *arXiv preprint arXiv:2205.11487*, 2022.
- [19] Lu Qi, Jason Kuen, Yi Wang, Jiuxiang Gu, Hengshuang Zhao, Zhe Lin, Philip Torr, and Jiaya Jia, “Open-world entity segmentation,” *arXiv preprint arXiv:2107.14228*, 2021.
- [20] Hui Su, Yue Ye, Zhiwei Chen, Mingli Song, and Lechao Cheng, “Re-attention transformer for weakly supervised object localization,” *arXiv preprint arXiv:2208.01838*, 2022.
- [21] 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 *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [22] Ruotian Luo, Brian Price, Scott Cohen, and Gregory Shakhnarovich, “Discriminability objective for training descriptive captions,” *arXiv preprint arXiv:1803.04376*, 2018.
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [24] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier, “From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions,” *Transactions of the Association for Computational Linguistics*, vol. 2, pp. 67–78, 2014.