

Supplementary Materials for LoMOE: Localized Multi-Object Editing via Multi-Diffusion

Anonymous Authors

CONTENTS

Contents	1
1 Introduction	1
2 Method Details	1
2.1 Regularised Inversion	1
2.2 Temperature Scaling	1
3 Ablation Study	1
3.1 Temperature Scaling	1
3.2 Bootstrap	1
3.3 Inversion	2
3.4 Inference Time	3
4 Experimental Protocol	4
4.1 Datasets	4
4.2 Baselines	4
4.3 Additional Results	5
5 Limitations	5
6 Broader Impact	6
References	6

1 INTRODUCTION

To keep the overall manuscript self-contained, we include additional details in the supplementary material. The source code for LoMOE along with the LoMOE-Bench dataset will be released in due time.

2 METHOD DETAILS

Specific aspects of the framework, including regularized inversion and temperature scaling, are described below.

2.1 Regularised Inversion

To softly enforce gaussianity on the inverted noise maps generated during the DDIM Inversion, we use a pairwise regularization $\mathcal{L}_{\text{pair}}$ [19] and a divergence loss \mathcal{L}_{KL} [12] weighted by λ (refer Sec. 3.1 of the main paper). These losses ensure that there is (1) no correlation between any pair of random locations and (2) zero mean, unit variance at each spatial location, respectively. Mathematically, the pairwise regularization loss is given by:

$$\mathcal{L}_{\text{pair}} = \sum_p \frac{1}{S_p^2} \sum_{\delta=1}^{S_p-1} \sum_{x,y,c} \eta_{x,y,c}^p (\eta_{x-\delta,y,c}^p + \eta_{x,y-\delta,c}^p) \quad (1)$$

where $\{\eta^0, \eta^1, \dots, \eta^p\}$ denote the noise maps with size S_p at the p^{th} pyramid level, δ denotes the offset which helps propagate long-range information [11, 19], and $\{x, y, c\}$ denotes a spatial location. Here, we set $p = 4$ and $\eta^0 = \epsilon_\theta \in \mathbb{R}^{64 \times 64 \times 4}$, where the subsequent noise maps are obtained via max-pooling.

The divergence loss is given by:

$$\mathcal{L}_{KL} = \sigma_{\epsilon_\theta}^2 + \mu_{\epsilon_\theta}^2 - 1 - \log(\sigma_{\epsilon_\theta}^2 + \varepsilon) \quad (2)$$

where μ_{ϵ_θ} and $\sigma_{\epsilon_\theta}^2$ denotes the mean and variance of ϵ_θ and ε is a stabilization constant.

2.2 Temperature Scaling

Given a vector $z = (z_1, \dots, z_n) \in \mathbb{R}^n$, it can be transformed into a probability vector via

$$\text{Softmax}(z|\tau)_i = \frac{e^{z_i/\tau}}{\sum_{j=1}^n e^{z_j/\tau}} \quad (3)$$

where τ is a temperature parameter [8] which varies the smoothness of the output distribution. In general, lower values of τ result in a sharp distribution, and increasing τ softens the distribution. This method has been used in applications such as model calibration [9], image restoration [20] and image inpainting [23]. In this work, we use a constant temperature scale to ensure the distributional smoothness of the cross-attention maps, setting $\tau = 1.25$. Further ablation on τ is discussed in Sec. 3.

3 ABLATION STUDY

In addition to the quantitative ablation of λ_{xa} and λ_b , we further study the impact of varying the temperature scaling parameter τ and bootstrap T_b . Specifically, we experiment for $\tau \in \{1.00, 1.25, 1.50, 1.75, 2.0\}$ and $T_b \in \{5, 10, 20, 30, 35\}$ and report the results in Table 1. We also do a visual ablation for the effect of the tightness of the mask on the performance of LoMOE. We also explore varying the tightness of the mask. We find that as long as the object to be edited is fully masked, the model performs well, as depicted in Figure 3.

3.1 Temperature Scaling

The results for variation in τ are summarized in Tab. 1 and have been depicted visually in Fig. 1. We observe that the edited image tends to go towards the source image with an increase in τ , which can be attributed to over-smoothing the distribution. This is also indicated by the *neural* metrics in Tab. 1, where an increase in τ results in increasing source CS and a decreasing target CS. This is further exemplified by the *background* metrics and Structural Distance, which are the best for $\tau = 2.00$. In this work, we set $\tau = 1.25$ as mentioned in Sec. 3.4 of the main paper. This choice of τ resulted in visually pleasing edits and we observed semantically coherent outputs for $\tau \in [1, 1.5]$.

3.2 Bootstrap

Upon analyzing the findings presented in Tab. 1, we opt for $T_b = 10$ based on the observation that the general structure and overall layout of the image is established within the first 10 denoising steps. Subsequently, the diffusion model manifests the finer details of the

<i>τ</i>	<i>T_b</i>	Source CLIP Score (\uparrow)	Background LPIPS (\downarrow)	Background PSNR (\uparrow)	Background SSIM (\uparrow)	Structural Distance (\downarrow)	Target CLIP Score (\uparrow)
1.00	-	23.4216	0.0586	30.1023	0.8822	0.0728	25.9163
1.25	-	23.7507	0.0522	30.4707	0.8849	0.0715	26.0902
1.50	-	24.1785	0.0497	30.7565	0.8863	0.0708	25.7919
1.75	-	25.0428	0.0466	31.1206	0.8875	0.0709	24.9769
2.00	-	25.4275	0.0409	31.5829	0.8896	0.0652	24.1544
-	05	23.5422	0.0562	30.1123	0.8838	0.0782	25.9403
-	10	23.5445	0.0546	30.3154	0.8847	0.0710	26.0740
-	20	23.4344	0.0587	30.0937	0.8822	0.0723	25.8746
-	30	23.4494	0.0618	29.8495	0.8792	0.0757	25.9404
-	35	23.2644	0.0621	29.8123	0.8792	0.0774	25.8089

Table 1: Further Ablation: We experiment with different values of the temperature parameter (τ) and bootstrap (T_b) parameters. From the *neural* and *background* metrics, we observe that the similarity between the edited and the input image increases for higher values of τ and that $T_b = 10$ is the optimal value for the bootstrap parameter.



Figure 1: Ablation on Temperature Scaling: Impact of an increasing temperature parameter, τ 's on the edits. We observe that an increase in the value of τ results in the edited image moving towards the input image. Empirically, we see visually appealing edits are achieved at $\tau = 1.25$.

image, in accordance with [2]. We also observe using a higher value of bootstrap aids in *addition-based edits*.

3.3 Inversion

As mentioned in Sec. 3.1 of the main paper, *inversion* helps initiate the editing procedure and ensures a coherent and controlled edit. To understand the impact of *inversion*, we compare two different initializations for the *edit* process (refer Sec. 3.2 of the main paper), namely (1) $x_T = x_{inv}$ and (2) $x_T = \zeta$. Here, $\zeta \in \mathbb{R}^{64 \times 64 \times 4}$ denotes

a random latent with elements sampled from $\mathcal{N}(0, 1)$. Specifically, we choose to showcase this impact on *style transfer* based edits.

From Fig. 2, we observe that the images with *inversion* are structurally much closer to the input image compared to the ones generated using a *random latent*, which is also indicated by the Structural Distance metric. In most cases, although using a random latent generates a faithful edit to the given prompt, it changes the content of the image, resulting in undesirable outputs. Therefore, using *inversion* is crucial for faithful image editing.

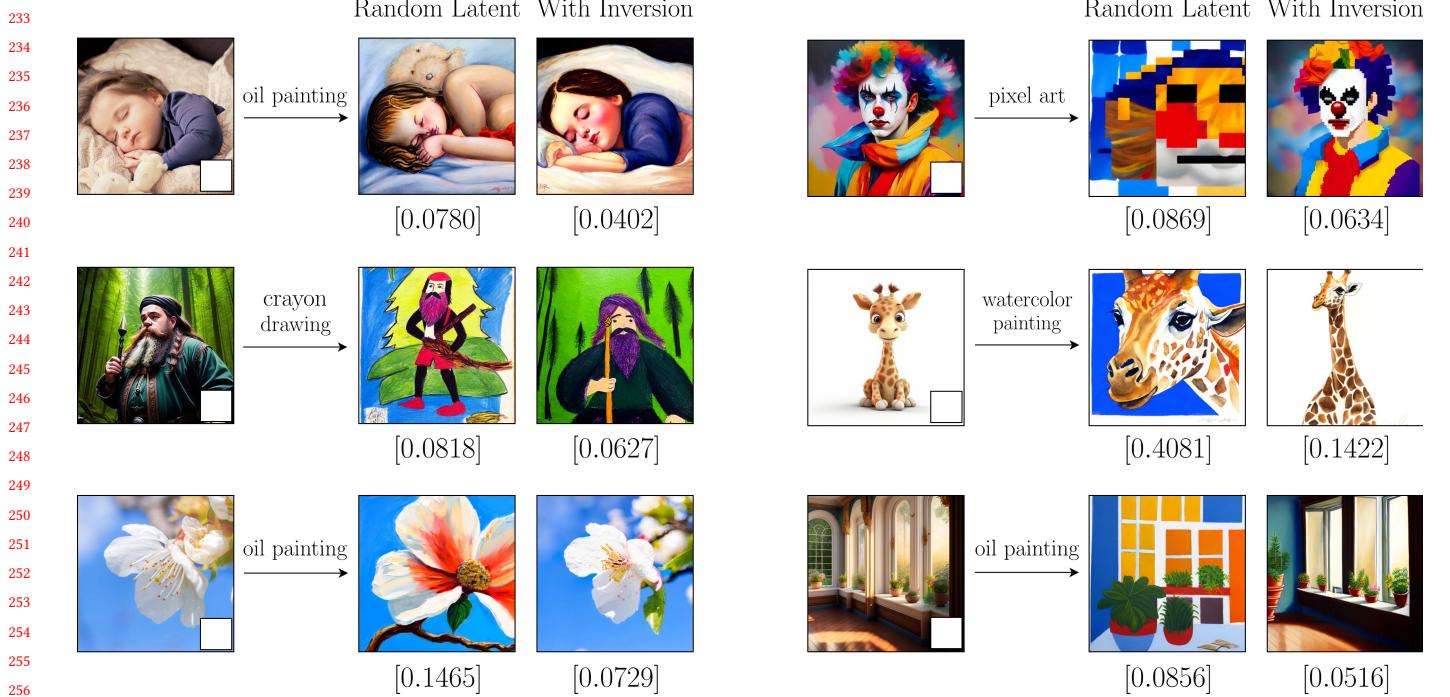


Figure 2: Ablation on Inversion: We study the impact of editing with a *random latent* compared to initiating the editing process via *inversion*. The outputs from LoMOE for both cases are captioned with the [Structural Distance (\downarrow)]. We observe that the structural similarity is preserved when using *inversion* instead of a *random latent* to initiate the editing process.

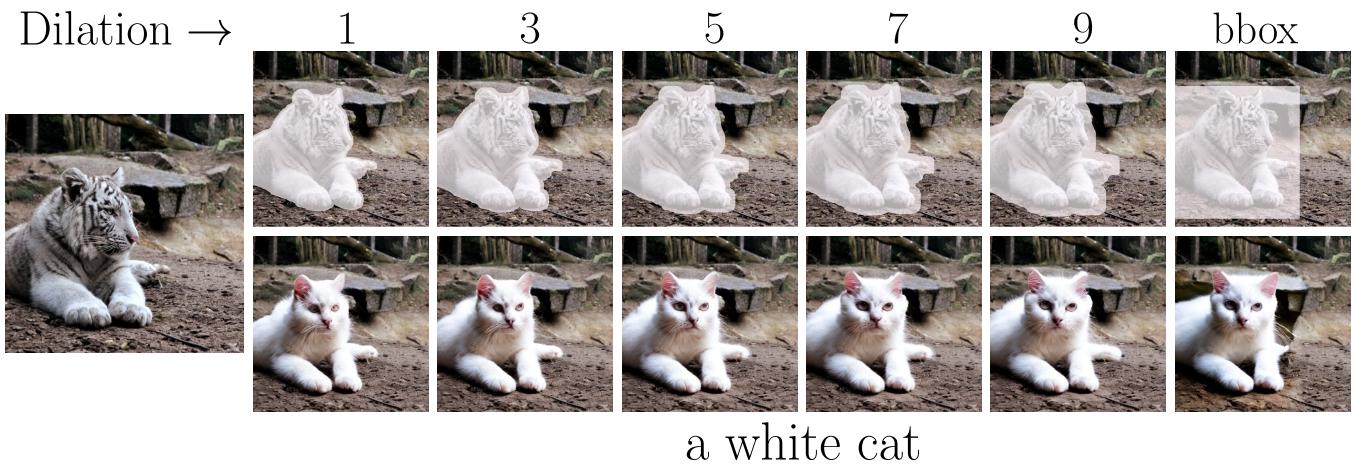


Figure 3: Illustrating the impact of enlarging the mask beyond the target object in LoMOE. As the mask undergoes progressive dilation, culminating in the use of a bounding box as the mask, it becomes evident that such transformations have minimal effect.

3.4 Inference Time

In a multi-object scenario, LoMOE separates itself by executing all edits in a single pass, resulting in substantial time savings compared to iterative methods. This is highlighted in Table 2, where our approach proves particularly advantageous in scenarios involving multiple objects, demonstrating a notable decrease in edit

time. Unlike other methods that run iteratively to generate multi-object edits, LoMOE's streamlined approach minimizes the need for repeated computations, enhancing overall efficiency. The gains in edit time underscore LoMOE's practical applicability in real-world editing tasks, showcasing its potential to streamline and expedite complex multi-object editing processes. The peak memory utilization (PMU) using the setup in Sec. 3.4 of the main paper has been

349 detailed in Table 3. Additionally, the time taken to optimize each
 350 objective is as follows (Ref. Main Paper Sec. 3):
 351

$$\begin{aligned} \text{start} &\rightarrow \mathcal{L}_{xa}(\text{Eqn. (9)}) : 0.165025 \text{ sec} \\ \mathcal{L}_{xa}(\text{Eqn. (9)}) &\rightarrow \mathcal{L}_b(\text{Eqn. (10)}) : 0.001386 \text{ sec} \\ \mathcal{L}_b(\text{Eqn. (10)}) &\rightarrow y_{t-1}^*(\text{Eqn. (13)}) : 0.185235 \text{ sec} \end{aligned}$$

Method	Inference Time for N masks (sec)				
	1	2	3	5	7
GLIDE	22.10	41.10	63.76	106.99	153.11
DiffEdit	22.25	41.30	65.91	110.85	152.60
BLD	27.20	51.60	80.40	135.24	185.37
SDInpaint	29.43	49.02	71.91	116.34	158.40
Iterative	24.97	45.75	70.50	117.35	162.37
LoMOE	23.19 (7.1)	31.3 (31.6)	39.35 (44.2)	55.47 (52.7)	76.15 (53.1)

360 **Table 2:** In a multi-object setting, we report the inference
 361 time of all the methods for varying number of masks N .
 362 **Iterative** denotes the average runtime of GLIDE, DiffEdit and
 363 BLD. We report the percentage improvement by LoMOE over
 364 **Iterative** (in green)

Method	Peak Memory Utilization for N masks (MiB)				
	1	2	3	5	7
<i>inversion</i>	-----	11029	-----	-----	-----
<i>reconstruction</i>	-----	6813	-----	-----	-----
<i>edit</i>	12609	17127	21367	30235	38979

380 **Table 3: Peak Memory Utilization (PMU) by LoMOE.**

4 EXPERIMENTAL PROTOCOL

4.1 Datasets

387 To facilitate a comparison between various baselines on *single-*
 388 *object* edits, we employ a modified subset of the PIE-Bench [10]
 389 dataset supplemented with images from AFHQ [5], COCO [14],
 390 and Imagen [22]. Overall, the benchmark consists of 300 images,
 391 covering editing types such as changing objects, adding objects,
 392 changing object content, changing object color, changing object
 393 material, changing background, and changing image style. Sample
 394 images for each edit type are shown in Fig. 14.

395 The newly proposed *multi-object* editing benchmark LoMOE-Bench
 396 consists of 64 images, covering various editing types, with each im-
 397 age featuring 2 to 7 masks paired with corresponding text prompts.
 398 The masks for the images in LoMOE-Bench and the supplemental
 399 images in the *single-object* dataset are generated using SAM [13]. In
 400 practice, the user is required to provide a bounding box around the
 401 object via a GUI interface, which then automatically saves the seg-
 402 mented mask, as mentioned on Sec. 4.1 of the main paper. Sample
 403 images from LoMOE-Bench are depicted in Fig. 18. We further ob-
 404 serve the variety of objects in LoMOE-Bench images and their masks
 405 in Fig. 4 and Fig. 5 shows the name of the objects being edited in

Method		Image	Mask	TIP	SMP	TMP	EIn
SDEdit	[16]	✓	✗	✓	✗	✗	✗
I-P2P	[3]	✓	✗	✗	✗	✗	✓
NTI (w/ P2P)	[17]	✓	✗	✓	✗	✗	✗
MasaCtrl	[4]	✓	✗	✓	✗	✗	✗
DiffEdit	[6]	✓	✓	✗	✓	✓	✗
GLIDE	[18]	✓	✓	✗	✗	✓	✗
BLD	[1]	✓	✓	✗	✗	✓	✗
SDInpaint	[21]	✓	✓	✗	✗	✓	✗
LoMOE		✓	✓	✗	✗	✓	✗

407 **Table 4: Annotations required by various baseline meth-
 408 ods included in the modified *single-object* dataset and LoMOE-
 409 Bench.**

410 the form of a word cloud. The images are also supplemented with
 411 various text-based annotations used by different baselines (refer
 412 Table 4) via a JSON file, including

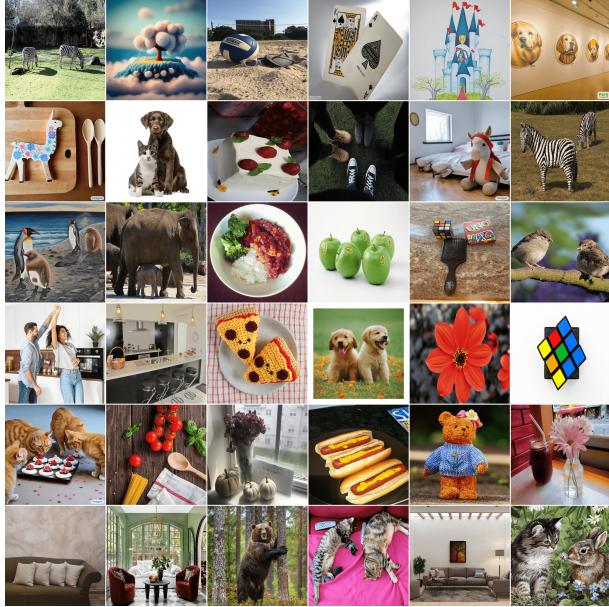
- **Target Image Prompt (TIP):** A complex prompt describing the complete image after the edit.
- **Source Mask Prompt (SMP):** A simple text prompt describing the object inside the masked region of the input image.
- **Target Mask Prompt (TMP):** A simple text prompt that describes the edited object inside the masked region.
- **Edit Instruction (EIn):** Edit instruction for I-P2P [3].

413 The dataset statistics are the following. On average, masks cover
 414 8.34% of the image (standard deviation: 10.1%), with each image
 415 containing an average of 3 masks (standard deviation: 1.17). In
 416 addition, we see the spatial distribution of the masks in a heat map
 417 in Fig. 6, demonstrating the frequency with which each pixel is
 418 masked.

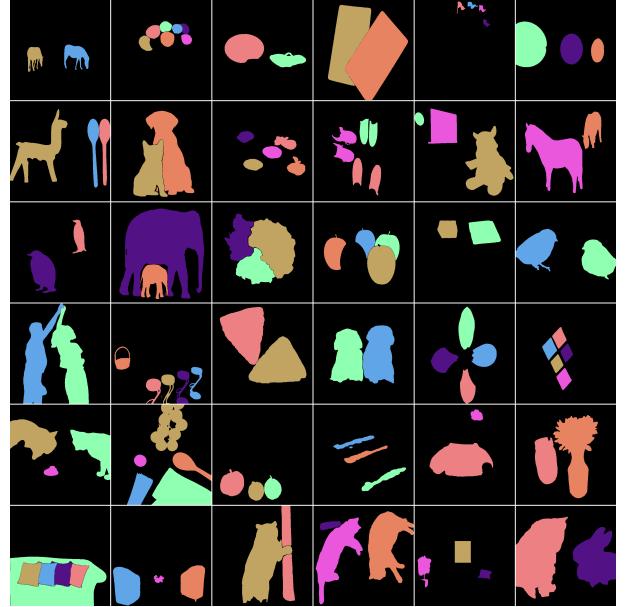
4.2 Baselines

419 We use the official implementation for all baseline methods us-
 420 ing PyTorch, except for DiffEdit as the code has not been made
 421 public. SDEdit uses the target prompt for text-guided image edit-
 422 ing and does not require any other input. DiffEdit by construction
 423 uses the DDIM solver, but the unofficial implementation uses DPM
 424 solver [15] for better sample efficiency. The method also generates
 425 noisy masks based on the source and target mask prompts, thus we
 426 choose to use the masks in the dataset (as mentioned in Sec. 4.2 of
 427 the main paper).

428 I-P2P requires an edit instruction along with the image and does
 429 not need any other inputs. For example, the edit instruction for
 430 the first image in Fig. 13 would look like: “*change the shape of*
 431 *the cake to a square*”. It is also important to note that although all
 432 other methods use the pre-trained Stable Diffusion model directly,
 433 Instruct-P2P is trained by finetuning this model. Finally, GLIDE
 434 and BLD are similar to LoMOE in that they only require the target
 435 mask prompt as additional inputs.



(a) Images to be edited



(b) Masks corresponding to the real images

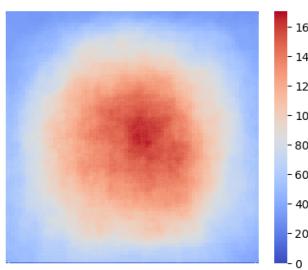
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522

Figure 4: LoMOE-Bench Dataset: Dataset comprising images with multiple edit targets along with their corresponding masks (a) and (b), and a word cloud



523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580

Figure 5: Word Cloud depicting the variety of objects present in LoMOE-Bench



523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580

Figure 6: Heatmap showing the number of masks that cover a certain pixel in 512 × 512 images.

4.3 Additional Results

We supplement the qualitative results provided in the main paper (Ref. Main Paper Sec. 5) by comparing LoMOE against baselines on more single-object edits, depicted in Fig. 13. Furthermore, we showcase single-object and multi-object edits with LoMOE in Figs. 14 and 18 for samples from LoMOE-Bench and the *single-object* benchmark for multiple masks and various edit types, respectively.

5 LIMITATIONS

In Figures 8, 9 and 7, we analyze the effect of lighting, shadows and reflections in the outputs produced by LoMOE. Additionally, Figures 10 and 11 delve into object coherence and some failure cases of LoMOE.

Lighting: In Fig. 8, we observe that the edited image has similar lighting properties to those of the source image.

Reflections: In Fig. 9, we observe that LoMOE adds in realistic reflections for *addition* based examples, based on the mask used. For example, there is no reflection in the “blue fish” example due to the mask constraint. This opens up avenues for a mask-free version of LoMOE. Additionally, in the example of a “sunny day”, the LoMOE uses the reflections in the water as a prior to aptly fill in the masked region.

Shadows: In Fig. 7, we observe that LoMOE attempts to add realistic shadows.

Realism: In Fig. 10, we demonstrate object-background coherence using a *background* change and an *addition* example. In both cases, LoMOE produces realistic and coherent results.

Failures: In Fig. 11, in cases where reflections fall outside the mask, LoMOE cannot remove them due to the mask-constraint, resulting in inconsistent lighting. Furthermore, LoMOE can perform global edits like style-transfer (Ref. Fig. 14) but cannot facilitate spatial manipulation like swapping or translation while preserving the identity and style of the object. While translation can be achieved through deletion and addition, as demonstrated in Fig. 11, it doesn't maintain the object's identity. Textual inversion [7] might be a potential solution to preserve identity for future works.

6 BROADER IMPACT

Generative image editing models are powerful tools that can create realistic and diverse images from text or other inputs. They have many potential applications in domains such as art, entertainment, education, medicine, and security. However, they also pose significant ethical and social challenges that need to be addressed. Some of these challenges include:

- The risk of generating harmful or offensive images that may violate human dignity, privacy, or rights.
- The possibility of manipulating or deceiving people with fake or altered images that may affect their beliefs, emotions, or behaviours.
- The difficulty of verifying the authenticity or provenance of images that may have legal or moral implications.
- The impact of replacing or reducing human creativity and agency with automated or algorithmic processes.

These challenges require careful consideration and regulation from various perspectives, such as technical, legal, ethical, and social. However, we believe that despite these drawbacks, better content creation methods will produce a net positive for society. Furthermore, we advocate for conducting such research in the public domain, emphasizing transparency and collaborative efforts to ensure responsible and beneficial outcomes for the broader community.

REFERENCES

- | | |
|---|-----|
| [1] Omri Avrahami, Ohad Fried, and Dani Lischinski. 2023. Blended Latent Diffusion. <i>ACM Trans. Graph.</i> 42, 4, Article 149 (jul 2023), 11 pages. https://doi.org/10.1145/3592450 | 639 |
| [2] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. 2023. MultiDiffusion: Fusing Diffusion Paths for Controlled Image Generation. In <i>ICML</i> . PMLR. | 640 |
| [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2023. Instructpix2pix: Learning to follow image editing instructions. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> . 18392–18402. | 641 |
| [4] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yingqiang Zheng. 2023. MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image Synthesis and Editing. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> . 22560–22570. | 642 |
| [5] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. 2020. StarGAN v2: Diverse Image Synthesis for Multiple Domains. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> . | 643 |
| [6] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. 2023. DiffEdit: Diffusion-based semantic image editing with mask guidance. In <i>The Eleventh International Conference on Learning Representations</i> . https://openreview.net/forum?id=3lge0p5o-M | 644 |
| [7] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. <i>arXiv preprint arXiv:2208.01618</i> (2022). | 645 |
| [8] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2019. On calibration of modern neural networks. (2019). | 646 |
| [9] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. 2018. Why ReLU Networks Yield High-Confidence Predictions Far Away From the Training Data | 647 |
| and How to Mitigate the Problem. <i>2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> (2018), 41–50. https://api.semanticscholar.org/CorpusID:55700923 | 648 |
| [10] Xuan Ju, Ailing Zeng, Yuxuan Bian, Shaoteng Liu, and Qiang Xu. 2024. PnP Inversion: Boosting Diffusion-based Editing with 3 Lines of Code. <i>International Conference on Learning Representations (ICLR)</i> (2024). | 649 |
| [11] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. 2020. Analyzing and Improving the Image Quality of StyleGAN. In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> . IEEE Computer Society, Los Alamitos, CA, USA, 8107–8116. https://doi.org/10.1109/CVPR42600.2020.00813 | 650 |
| [12] Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In <i>International Conference on Learning Representations, ICLR 2014</i> , Yoshua Bengio and Yann LeCun (Eds.). | 651 |
| [13] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment Anything. <i>arXiv:2304.02643</i> (2023). | 652 |
| [14] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In <i>Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13</i> . Springer, 740–755. | 653 |
| [15] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2022. DPM-Solver: A Fast ODE Solver for Diffusion Probabilistic Model Sampling in Around 10 Steps. <i>arXiv preprint arXiv:2206.00927</i> (2022). | 654 |
| [16] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2021. Sdedit: Guided image synthesis and editing with stochastic differential equations. <i>arXiv preprint arXiv:2108.01073</i> (2021). | 655 |
| [17] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2023. Null-text inversion for editing real images using guided diffusion models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> . 6038–6047. | 656 |
| [18] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. In <i>International Conference on Machine Learning</i> . PMLR, 16784–16804. | 657 |
| [19] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. 2023. Zero-shot image-to-image translation. In <i>ACM SIGGRAPH 2023 Conference Proceedings</i> . 1–11. | 658 |
| [20] Tobias Plötz and Stefan Roth. 2018. Neural Nearest Neighbors Networks. In <i>Proceedings of the 32nd International Conference on Neural Information Processing Systems (Montreal, Canada) (NIPS'18)</i> . 1095–1106. | 659 |
| [21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> . 10684–10695. | 660 |
| [22] Si Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. 2023. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> . 18359–18369. | 661 |
| [23] Xiang Zhou, Yuan Zeng, and Yi Gong. 2023. Learning to Scale Temperature in Masked Self-Attention for Image Inpainting. <i>ArXiv abs/2302.06130</i> (2023). https://api.semanticscholar.org/CorpusID:256827540 | 662 |

697
698
699
700
701
702
703
704
705
706

INPUT



MASK



OUTPUT



CAKES STACKED ON EACH OTHER

755
756
757
758
759
760
761
762
763
764

A SHINY WHITE CAR

765
766
767
768
769
770
771
772
773
774
775

ORANGE CURTAINS

776
777
778
779
780
781
782
783
784
785
786**Figure 7: Shadows:** The images illustrate that the shadows are consistent with the original image before object editing.732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

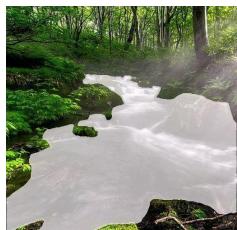
INPUT

MASK

OUTPUT



A SQUARE CAKE



A ROAD



A PIECE OF SALMON



RED BERRIES

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

Figure 8: Consistency in illumination of object edit with background: Note the initially illuminated portion of the cake situated on the right, with the relatively dimmer region on the left. Output depiction of the square cake also adheres to the luminosity characteristics inherent in the initial input, thereby ensuring coherence in visual representation. Observe a similar coherence in salmon, road, and red berries.

929
930
931
932
933
934
935
936
937
938
939

INPUT



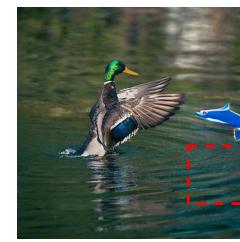
MASK



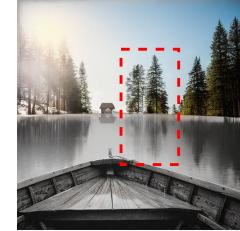
OUTPUT

940
941942
943
944
945
946
947
948
949
950951
952953
954
955
956
957
958
959
960
961962
963964
965
966
967968
969
970
971
972973
974975
976977
978979
980
981
982983
984
985
986

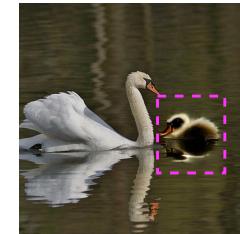
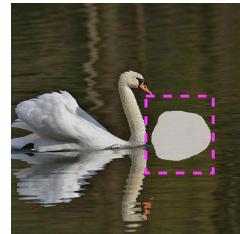
A GLASS BOWL

987
988
989
990
991
992
993
994
995
996
997998
9991000
1001
1002
1003
1004
1005
1006
1007
10081009
10101011
1012
1013
1014
1015
1016
1017
1018
10191020
10211022
1023
10241025
1026
10271028
1029
10301031
1032
10331034
1035
10361037
1038
10391040
1041
1042
1043
1044

A BLUE FISH



A SUNNY DAY AT A FOREST



A BABY SWAN

Figure 9: Observe the consistency in reflections of the edited object with respect to the background as highlighted by a dashed rectangle for emphasis.

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

Figure 10: Realism: Note the meticulous attention to detail in the editing of the forest scene, which imbues it with a semblance of reality. Likewise, the dog collar in the output exhibits a naturalistic deformation along the neck region, thereby enhancing realism.

INPUT



MASK



OUTPUT

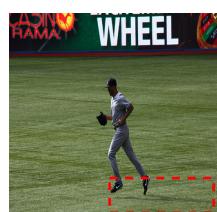


A FOREST

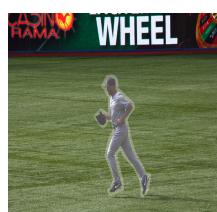


A RED DOG COLLAR

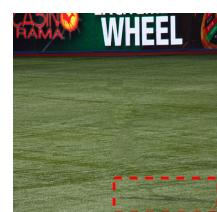
INPUT



MASK



OUTPUT



REMOVE



MOVE FOOTBALL FROM A TO B

Figure 11: Deletion: LoMOE’s limitations about shadow handling. Specifically, upon the removal of an individual, their associated shadows persist post-editing due to our approach’s failure to address elements outside the designated mask area. Furthermore, it is noteworthy that despite the presence of motion, the edited football image retains clarity without exhibiting blurriness.

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160



Figure 12: Additional Comparison among Contemporary Methods for Single Object Edits: We present a qualitative comparison of LoMOE against other baseline methods on additional single-object edits. The observations stand similar to Fig. 3 in the main paper, where our proposed method LoMOE makes the intended edit, preserves the unmasked region and avoids unintended attribute edits.

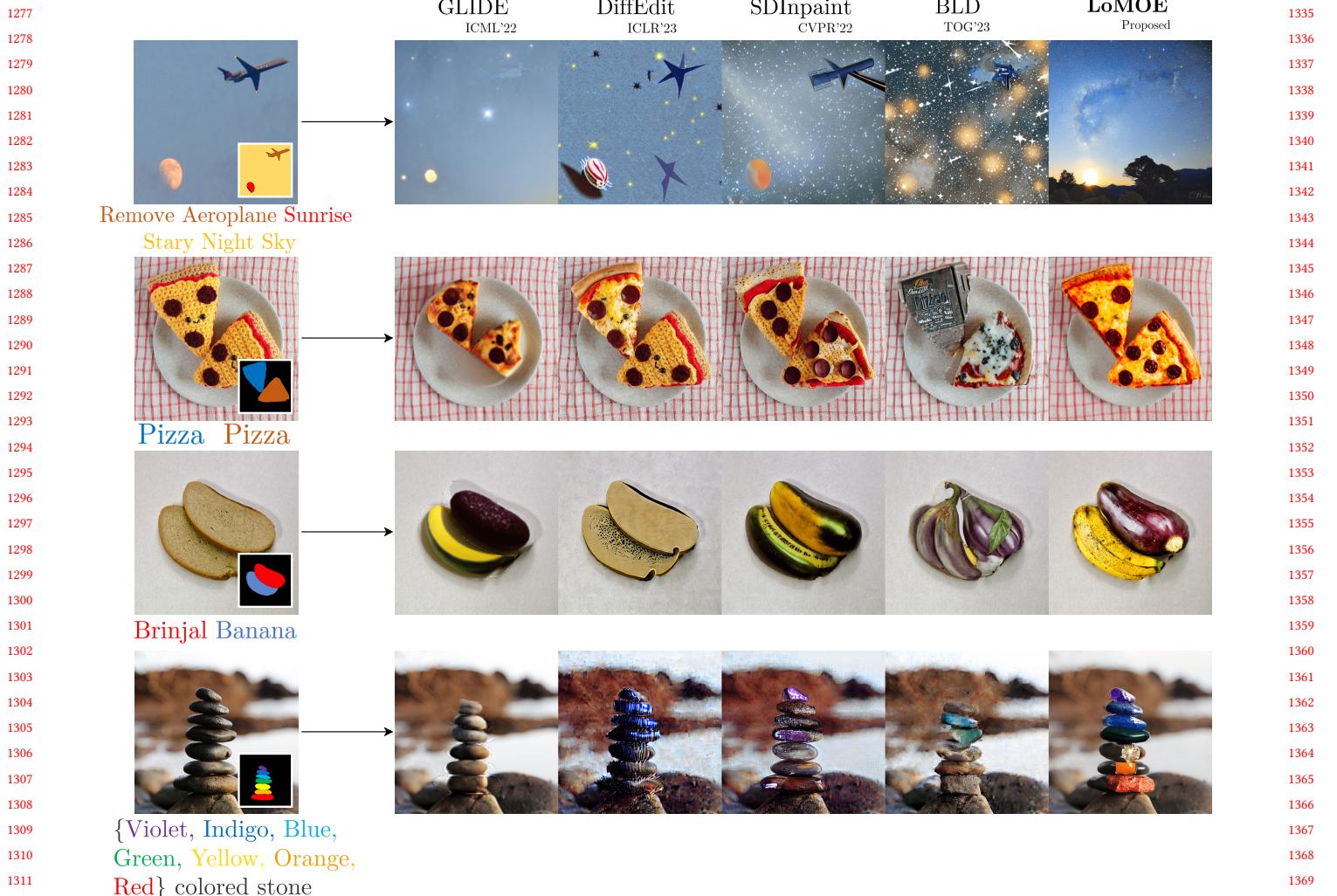


Figure 13: Additional Comparison among Contemporary Methods for Multi Object Edits: We present a qualitative comparison of LoMOE against other baseline methods on additional multi-object edits. The observations stand similar to Fig. 4 in the main paper, where our proposed method LoMOE makes the intended edit, preserves the unmasked region, and avoids unintended attribute edits.

1393	Edit Type	Input Image	Mask	Annotation	Edited image	1451
1394						1452
1395						1453
1396	Change Object			{ TIP: a plate with roasted salmon on it SMP: a steak TMP: a roasted salmon EIn: change the steak to roasted salmon }		1454
1397						1455
1398						1456
1399						1457
1400	Change Object			{ TIP: a roaring tiger wearing a pink hat SMP: a cat TMP: a roaring tiger EIn: Change the animal from a cat to a roaring tiger }		1458
1401						1459
1402	Adding Object			{ TIP: a dog with a red dog collar looking at the camera SMP: no collar TMP: red dog collar EIn: put a red dog collar on the dogs neck }		1460
1403						1461
1404	Adding Object			{ TIP: a small brown bird sitting on top of a pine tree SMP: background TMP: a bird flying EIn: Add a small bird flying }		1462
1405						1463
1406						1464
1407						1465
1408	Adding Object			{ TIP: a bulb with red glowing filament hangs from a string SMP: filament TMP: red glowing filament EIn: Add red color to the lights }		1466
1409						1467
1410	Change Content			{ TIP: a glass of hot chocolate SMP: glass of cold coffee TMP: glass of hot chocolate EIn: change cold coffee to hot chocolate }		1468
1411						1469
1412	Change Content			{ TIP: a woman with blue hair SMP: brown hair TMP: blue hair EIn: change brown to blue hair }		1470
1413						1471
1414	Change Color			{ TIP: a red curtain is hanging over a stone wall SMP: red curtain TMP: orange curtain EIn: change the color of the curtain from red to orange }		1472
1415						1473
1416	Change Color			{ TIP: a shiny metal car SMP: rusted metal car TMP: shiny metal car EIn: make the car shiny metal }		1474
1417						1475
1418	Change Material			{ TIP: a fabric ladybug with black spots on its back is sitting on a leaf SMP: a ladybug TMP: a fabric ladybug EIn: make the ladybug fabric }		1476
1419						1477
1420	Change Material			{ TIP: "a bride standing in a garden" SMP: on a mountain TMP: a garden EIn: substitute the mountain with a garden }		1478
1421						1479
1422	Change Background			{ TIP: a barn owl with a black eye and spots in a blurry forest background SMP: a grey background TMP: a blurry forest background EIn: add a blurry forest background }		1480
1423						1481
1424	Change Style			{ TIP: a spring road lined with trees and leaves SMP: an autumn road lined with trees and leaves TMP: a spring road lined with trees and leaves EIn: change the season from autumn to spring }		1482
1425						1483
1426	Change Style			{ TIP: an oil painting of a cherry blossom with blue sky SMP: cherry blossom with blue sky TMP: an oil painting of a cherry blossom with blue sky EIn: add an oil painting effect to the cherry blossom }		1484
1427						1485
1428						1486
1429						1487
1430						1488
1431						1489
1432						1490
1433						1491
1434						1492
1435						1493
1436						1494
1437						1495
1438						1496
1439						1497
1440						1498
1441						1499
1442						1500
1443						1501
1444						1502
1445						1503
1446	Figure 14: Single Object Benchmark: Examples from Single-Object dataset. The columns are (1) Edit type (2) The input image on which the editing is done, (3) The mask used for localizing the edit, (4) JSON annotation containing the Target Image Prompt (TIP), Source Mask Prompt (SMP), Target Mask Prompt (TMP), and the Edit Instruction (EIn), and (5) The edited images produced by LoMOE.					1504
1447						1505
1448						1506
1449						1507
1450						1508

1509	Input Image	Mask	Annotation	Edited image	1567
1510			{ SMP: "a brown bird", "a brown bird" TMP: "a crochet bird", "a origami bird" }		1568
1511					1569
1512					1570
1513					1571
1514					1572
1515					1573
1516					1574
1517					1575
1518					1576
1519					1577
1520					1578
1521					1579
1522			{ SMP: "pink cake frosting", "a real jar of candy", "chocolate cake truffles" TMP: "pink cake frosting", "a real jar of candy", "chocolate cake truffles" }		1580
1523					1581
1524					1582
1525					1583
1526					1584
1527					1585
1528					1586
1529					1587
1530					1588
1531					1589
1532					1590
1533					1591
1534					1592
1535					1593
1536					1594
1537					1595
1538					1596
1539					1597
1540					1598
1541					1599
1542					1600
1543					1601
1544					1602
1545					1603
1546					1604
1547					1605
1548					1606
1549					1607
1550					1608
1551					1609
1552					1610
1553	Figure 15: LoMOE-Bench: Examples from Multi-Object Dataset. The columns are (1) The input image on which the editing is done, (2) The masks used for localizing the edit, (3) JSON annotation containing the Source Mask Prompts (SMP) and Target Mask Prompts (TMP), and (4) The edited images produced by LoMOE.				
1554					1611
1555					1612
1556					1613
1557					1614
1558					1615
1559					1616
1560					1617
1561					1618
1562					1619
1563					1620
1564					1621
1565					1622
1566					1623
					1624

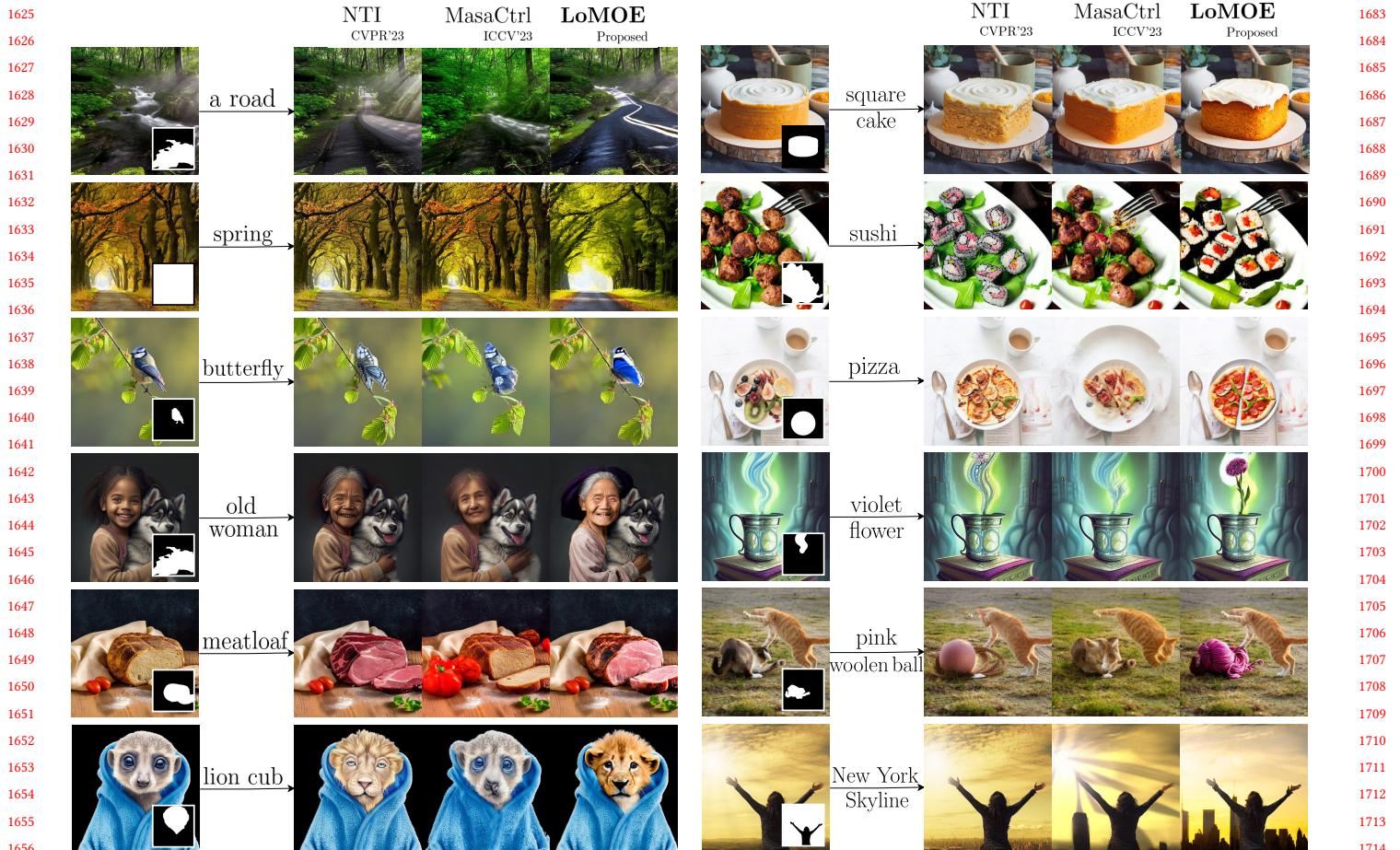


Figure 16: Additional Comparison among Contemporary Methods for Single Object Edits:

1741

1742

1743

1744

1745

1746

1747

1748

1749

1750

1751

1752

1753

1754

1755

1756

1757

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1778

1779

1780

1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

SDEdit

ICLR'22

InstructP2P

CVPR'23

NTI

CVPR'23

MasaCtrl

ICCV'23

LoMOE

Proposed

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

1836

1837

1838

1839

1840

1841

1842

1843

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

1856



Figure 17: Additional Comparison among Contemporary Methods for Multi Object Edits:

1857

1858

1859

1860

1861

1862

1863

1864

1865

1866

Remove Aeroplane Sunrise
Stary Night SkySDEdit
ICLR'22InstructP2P
CVPR'23NTI
CVPR'23MasaCtrl
ICCV'23LoMOE
Proposed

1867



Pizza Pizza



1868

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

Brinjal Banana



1884

1885

1886

1887

1888

1889

1890

1891

{Violet, Indigo, Blue,
Green, Yellow, Orange,
Red} colored stone

1892

1893

1894

1895

1896

1897

1898

1899

1900

1901

1902

1903

1904

1905

1906

1907

1908

1909

1910

1911

1912

1913

1914

1915

1916

1917

1918

1919

1920

1921

1922

1923

1924

1925

1926

1927

1928

1929

1930

1931

1932

1933

1934

1935

1936

1937

1938

1939

1940

1941

1942

1943

1944

1945

1946

1947

1948

1949

1950

1951

1952

1953

1954

1955

1956

1957

1958

1959

1960

1961

1962

1963

1964

1965

1966

1967

1968

1969

1970

1971

1972

Figure 18: Additional Comparison among Contemporary Methods for Multi Object Edits: