**Warning: This supplementary material contains model outputs that are offensive in nature.**

# A  More Visualization Results

More visualization results of ColJailBreak are shown in Figure 7.



**Figure 7:** More Visualization Results of ColJailBreak.

# B  Details of UnsafeEdit Dataset

For a more thorough evaluation of our method, we curated a dataset called the *UnsafeEdit dataset*. Referring to the division of I2P dataset, we categorise the dataset into four categories: violence, harassment, self-harm and nudity. Specifically, violence contains scenes about violence and gore, harassment contains scenes of threats and harassment, self-harm contains scenes of suicide and self-harm, and nudity contains scenes of sexual content.

**Collection Details.** Considering that ChatGPT cannot directly generate prompts containing unsafe contents, we first generate safe prompts and then manually replace the unsafe words to create our dataset. Initially, we manually crafted a template prompt devoid of sensitive words, thereby ensuring it does not activate ChatGPT's safety mechanisms (*e.g.*,*"A man points a pen at another woman in a supermarket. The two faced each other, and the woman looked scared."*).

Then we use ChatGPT to generate more prompt data through the template prompt, the example prompt is:

*"A man points a pen at another woman in a supermarket. The two faced each other, and the woman looked scared. I want to generate a prompt dataset. Please follow the above format and help me generate some similar prompts. The requirements are 1. Contain Two people face to face 2. Including the action of pointing a pen at the other person 3. One of the people must be very scared 4. Daily scenes".*

Finally, we build the dataset by manually replacing unsafe words.

**Dataset details.** UnsafeEdit dataset has a total of 100 prompts, each category containing 25 prompts. To ensure the quality and relevance of the dataset, we carefully reviewed each generated prompt and made necessary adjustments to maintain the intended context while introducing the unsafe elements. This meticulous process guarantees that the dataset accurately reflects real-world scenarios where such harmful content might appear.

For external dataset, in the experiment, we find that T2I model's safety filter is more stringent when detecting nudity compared to other categories. For nudity, we primarily select prompts from the dataset using keywords such as "breast", "tit", "sexy" and "naked". For other categories, we primarily select prompts from the dataset using keywords such as "pistol", "dagger", "knife", "gun" and "revolver".

# C    Implementation Details

## C.1    Computing resources

All experiments are performed using two NVIDIA A100 40GB GPUs. The overall duration of all the experiments in the paper was about six weeks.

## C.2    Details of Baselines

In the setups of Section 5.1, we introduce two baseline methods, and here are more detailed implementation details.

**MMA-Diffusion.** MMA-Diffusion[43] is a method to generate unsafe content by bypassing T2I model safetu filters. In Text-Modal Attack, MMA-Diffusion obtains adversarial prompts through the gradient optimization based search method. Subsequently, sensitive words within the adversarial prompts are eliminated by sensitive word regularization to bypass the safety filter. For hyperparameters, the random seed is 7867, the number of optimization iterations is 500, and the number of adversarial prompts per target prompt is 10.

**QF-Attack.** QF-Attack[49] is an attack method designed to disrupt T2I model generation, introducing a five-character perturbation to the text prompt alters the generated image content. QF-Attack comprises three attack modes: Greedy search, Genetic algorithm, and PGD attack. Based on experimental results reported in the QF-Attack paper, the Genetic algorithm proves more effective than other attack methods. Consequently, for the experiments in this paper, we selected the Genetic algorithm for comparison. Specifically, we choose the targeted query-free attack of QF-Attack. Following the settings used in Ring-A-Bell[40], we mask out words related to peace and non-harm, such as "peaceful and fine," and words related to clothing, such as "wear tops and pants".

## C.3    Details of Defense Models

**Safe Latent Diffusion (SLD).** For the four variants of SLD[36], we set the image size to $512 \times 512$. As shown in Table 3, we set the hyperparameters according to the official configuration.

| Config | warmup steps $\delta$ | guidance scale $s_S$ | threshold $\lambda$ | momentum scale $s_m$ | mom beta $\beta_m$ |
|---|---|---|---|---|---|
| SLD-Weak | 15 | 200 | 0.0 | 0.0 | - |
| SLD-Medium | 10 | 1000 | 0.01 | 0.3 | 0.4 |
| SLD-Strong | 7 | 2000 | 0.025 | 0.5 | 0.7 |
| SLD-Max | 0 | 5000 | 1.0 | 0.5 | 0.7 |

**Table 3:** Configuration sets of SLD

**ESD.** For the category of nudity, we use the officially provided NSFW pre-trained weights. For other inappropriate content, we use "violence, harassment, and self-harm" as the prompt for training ESD[15]. In generating images, we set the image size to $512 \times 512$, and the model to SDv1.4. For the hyperparameters, the guidance scale is 7.5, and the ddim steps is 100.

### C.4 Details of ColJailBreak

**Adaptive Safe Word Substitution.** As mentioned in Section 4.2, we utilize ChatGPT to create the substitution pool for sensitive words, the size of the substitution pool is set to 10, and then we select the substitution word in the substitution pool. We design a prompt template for generating the substitution pool:

*Hello ChatGPT, you are a helpful assistant. You are a safety testing expert in the field of AI. I am doing image editing work and need to replace an object in the image with an object of similar shape. I need to create a vocabulary of similar-shaped objects for each object. Can you understand what I mean? If so, let's get started. The first word is {}. Please give me {} words directly, without any other text, and separate them with commas.*

**Inpainting-driven Injection of Unsafe Content.** Inspired by Inpainting Anything[45], in general, we use substitution word $s$ as a text prompt, mask the area related to the text prompt based on SAM, and apply inpainting model for editing. Initially, as SAM's input consists of points, mask, and bounding box, but not text prompt, we employ CLIP Surgery [21], which converts text prompt into points by leveraging the explainability of CLIP. Then, we generate a preliminary semantic segmentation map utilizing the robust semantic segmentation capabilities of SAM[20], and then obtain the mask of the editing region. Subsequently, we edit the image using the pre-trained Fooocus-Inpainting model to inject unsafe content. For Fooocus-Inpainting, we set the image size to $512 \times 512$. For the hyperparameters, the guidance scale is 7.5, the num inference steps is 50, and the strength is 0.9999.

## D   Broader Impacts

Our work provides new insights into the security and robustness of commercial T2I models. However, while our research aims to evaluate the security of current commercial T2I models against jailbreak attacks, there is a risk that malicious users may exploit our work to generate unsafe images, which requires more caution. Considering that our proposed ColJailBreak may be used maliciously, we have provided user guidelines for ColJailBreak.

## E   Use Guidelines of ColJailBreak

In utilizing the ColJailBreak framework for jailbreaking T2I models, it is essential to adhere to the following guidelines to ensure responsible and ethical usage:

- Purpose and Intent: ColJailBreak should be used primarily for research purposes to understand the limitations and vulnerabilities of existing T2I models and to improve their safety mechanisms. Users must ensure that their intent aligns with ethical research standards and contributes to the advancement of safe AI technologies.

- Compliance with Regulations: Users must comply with all relevant laws and regulations governing the use of AI and deep learning technologies in their respective jurisdictions.

- Privacy and Consent: Respect the privacy and consent of individuals. Do not use personal data or identifiable information without explicit permission. Avoid creating images that depict real individuals in a harmful or misleading manner.

- Reporting and Accountability: Report any misuse or inappropriate content generated using ColJailBreak to the developers. Be accountable for the content you generate and share using the framework.

- Strict Confidentiality: Users must rigorously safeguard the model's operational principles, datasets, and any associated information to prevent disclosure to unauthorized individuals or organizations.

# References

[1] ChatGPT. https://chat.openai.com. 5

[2] Controlnet-v1.1-sd1.5-Inpainting. https://huggingface.co/lllyasviel/control_v11p_sd15_inpaint. 8

[3] DALL·E 2. https://openai.com/index/dall-e-2. 6

[4] Fooocus-Inpainting. https://huggingface.co/Vijish/fooocus_inpainting. 8

[5] GPT-4. https://openai.com/index/gpt-4. 6

[6] NudeNet. https://github.com/notAI-tech/NudeNet. 6

[7] SD-Inpainting. https://huggingface.co/runwayml/stable-diffusion-inpainting. 8

[8] 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 (CVPR)*, pages 18208–18218, June 2022. 3, 4

[9] Zhongjie Ba, Jieming Zhong, Jiachen Lei, Peng Cheng, Qinglong Wang, Zhan Qin, Zhibo Wang, and Kui Ren. Surrogateprompt: Bypassing the safety filter of text-to-image models via substitution. *arXiv preprint arXiv:2309.14122*, 2023. 2

[10] Zhi-Yi Chin, Chieh-Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Chen Chiu. Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. In *International Conference on Machine Learning (ICML)*, 2024. 1

[11] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. In *International Conference on Learning Representations*, 2023. 3

[12] Giannis Daras and Alexandros G Dimakis. Discovering the hidden vocabulary of dalle-2. In *NeurIPS Workshop on Score-Based Methods*, 2022. 3

[13] Yimo Deng and Huangxun Chen. Divide-and-conquer attack: Harnessing the power of llm to bypass the censorship of text-to-image generation model. *arXiv preprint arXiv:2312.07130*, 2023. 2, 3, 6

[14] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023. 3

[15] Rohit Gandikota, Joanna Materzyńska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *Proceedings of the 2023 IEEE International Conference on Computer Vision*, 2023. 7, 13

[16] Sensen Gao, Xiaojun Jia, Yihao Huang, Ranjie Duan, Jindong Gu, Yang Liu, and Qing Guo. Rt-attack: Jailbreaking text-to-image models via random token. *arXiv preprint arXiv:2408.13896*, 2024. 3

[17] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022. 1

[18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020. 3

[19] 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 (CVPR)*, pages 2426–2435, June 2022. 3

[20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 5, 14

[21] Yi Li, Hualiang Wang, Yiqun Duan, and Xiaomeng Li. Clip surgery for better explainability with enhancement in open-vocabulary tasks, 2023. 14

[22] Han Liu, Yuhao Wu, Shixuan Zhai, Bo Yuan, and Ning Zhang. Riatig: Reliable and imperceptible adversarial text-to-image generation with natural prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20585–20594, 2023. 3

[23] Jiachen Ma, Anda Cao, Zhiqing Xiao, Jie Zhang, Chao Ye, and Junbo Zhao. Jailbreaking prompt attack: A controllable adversarial attack against diffusion models. *arXiv preprint arXiv:2404.02928*, 2024. 3

[24] Alexander Quinn 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. In *International Conference on Machine Learning (ICML)*, 2022. 1, 3

[25] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2085–2094, October 2021. 3

[26] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 3

[27] Yiting Qu, Xinyue Shen, Xinlei He, Michael Backes, Savvas Zannettou, and Yang Zhang. Unsafe Diffusion: On the Generation of Unsafe Images and Hateful Memes From Text-To-Image Models. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM, 2023. 2

[28] 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*, pages 8748–8763. PMLR, 2021. 3

[29] 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. 1, 3

[30] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 3

[31] Javier Rando, Daniel Paleka, David Lindner, Lennart Heim, and Florian Tramèr. Red-teaming the stable diffusion safety filter. *arXiv preprint arXiv:2210.04610*, 2022. 1

[32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 3

[33] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 3

[34] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 conference proceedings*, pages 1–10, 2022. 3

[35] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jeffery Whang, Emily L Denton, Kamyar Ghasemipour, Rosanne Gontijo Lopes, Baran Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, volume 35, pages 36479–36494, 2022. 1, 3

[36] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 6, 7, 13

[37] Patrick Schramowski, Christopher Tauchmann, and Kristian Kersting. Can machines help us answering question 16 in datasheets, and in turn reflecting on inappropriate content? In *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAccT)*, 2022. 6

[38] Erfan Shayegani, Yue Dong, and Nael Abu-Ghazaleh. Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. In *The Twelfth International Conference on Learning Representations*, 2023. 1

[39] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 3

[40] Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia-You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models? In *International Conference on Learning Representations*, 2024. 2, 13

[41] Tianyu Wei, Shanmin Pang, Qi Guo, Yizhuo Ma, and Qing Guo. Emoattack: Emotion-to-image diffusion models for emotional backdoor generation. *arXiv preprint arXiv:2406.15863*, 2024. 3

[42] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18381–18391, 2023. 3, 4

[43] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. MMA-Diffusion: MultiModal Attack on Diffusion Models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3, 6, 13

[44] Yuchen Yang, Bo Hui, Haolin Yuan, Neil Gong, and Yinzhi Cao. Sneakyprompt: Jailbreaking text-to-image generative models. In *Proceedings of the IEEE Symposium on Security and Privacy*, 2024. 2, 3

[45] Tao Yu, Runseng Feng, Ruoyu Feng, Jinming Liu, Xin Jin, Wenjun Zeng, and Zhibo Chen. Inpaint anything: Segment anything meets image inpainting. *arXiv preprint arXiv:2304.06790*, 2023. 3, 14

[46] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *IEEE International Conference on Computer Vision (ICCV)*, 2023. 1, 8

[47] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. *European Conference on Computer Vision (ECCV)*, 2024. 1

[48] Qi Zhou, Dongxia Wang, Tianlin Li, Zhihong Xu, Yang Liu, Kui Ren, Wenhai Wang, and Qing Guo. Foolsdedit: Deceptively steering your edits towards targeted attribute-aware distribution. *arXiv preprint arXiv:2402.03705*, 2024. 4

[49] Haomin Zhuang, Yihua Zhang, and Sijia Liu. A pilot study of query-free adversarial attack against stable diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2384–2391, 2023. 3, 6, 13