MasterWeaver: Taming Editability and Face Identity for Personalized Text-to-Image Generation – Supplementary Materials –

Yuxiang Wei^{1,2} Zhilong Ji³ Jinfeng Bai³ Hongzhi Zhang¹ Lei Zhang^{2(\vee)} Wangmeng Zuo^{1,4(\vee)}

¹Harbin Institute of Technology ² The Hong Kong Polytechnic University ³Tomorrow Advancing Life ⁴ Pazhou Lab Huangpu

The following materials are provided in this supplementary file:

- Sec. S1: more implementation details, including the detailed framework, the editing direction loss, and the training details.
- Sec. S2: more dataset details, including the training dataset, face-augmented dataset, and evaluation dataset.
- Sec. S3: more ablation studies, including the layer selection of \mathcal{L}_{edit} , the value of λ , and the number of reference images.
- Sec. S4: more comparisons and results.
- Sec. S5: limitation and social impact.

S1 More Implementation Details

S1.1 Framework of MasterWeaver

Fig. S1 illustrates the detailed framework of MasterWeaver. As shown in the figure, the ID Mapper consists of a stack of cross attention, feed forward layer and self attention layer. To effectively incorporate the identity feature f with text feature to steer the personalized generation, we adopt the dual cross attention mechanism [12, 14]. As illustrated in Fig. S1, for each cross attention block of SD, we further introduce two learnable projection layers W_K^f and W_V^f . The identity feature is then integrated by Attention (Q, K^f, V^f) , where $K^f = W_K^f \cdot f$ and $V^f = W_V^f \cdot f$ represent the projected key and value matrices of the identity information, respectively. To fuse the identity information with the text information, we sum them by:

$$Out = \text{Attention}(Q, K, V) + \lambda \text{Attention}(Q, K^f, V^f), \tag{S1}$$

where λ is a trade-off hyperparameter and set as 1 during training. During training, we optimize the parameters of ID Mapper and the projection layers simultaneously, and keep the parameters of SD fixed.

2 Y. Wei et al.



Fig. S1: Framework of our MasterWeaver.

S1.2 Editing Direction Loss

To calculate the editing direction loss \mathcal{L}_{edit} , we utilize the output features from the cross attention blocks in the first decoder layer. Specifically, given the source prompt y, target prompt y' and reference image x, the editing directions of the l-th cross attention block for SD and MasterWeaver are $\Delta_{\epsilon_{\theta}}^{l}(y, y') \in \mathbb{R}^{H \times W \times D}$ and $\Delta_{\epsilon_{\theta'}}^{l}(y, y', x) \in \mathbb{R}^{H \times W \times D}$, where H, W, and D represent the height, width, and dimension of the output feature, respectively. Then, the editing direction loss is calculated as follows:

$$\mathcal{L}_{edit} = \sum_{l} \mathcal{L}_{edit}^{l}, \tag{S2}$$

$$\mathcal{L}_{edit}^{l} = \frac{1}{H \times W} \sum_{h=1}^{H} \sum_{w=1}^{W} \left[1 - \operatorname{Sim}[(m \cdot \Delta_{\epsilon_{\theta}}^{l})_{h,w}, (m \cdot \Delta_{\epsilon_{\theta'}}^{l})_{h,w}] \right].$$
(S3)

For conciseness, we have omitted the y, y' and x. The source and target prompts employed in our experiments are listed in Table S6. For each training iteration, we randomly select a source and target prompt pair to compute the loss.

S1.3 Training Details

Textual Inversion [1]. We use the official implementation of Textual Inversion¹, training it using SD V1.5. For each identity, the experiment is conducted with a batch size of 1 and a learning rate of 0.005 for 5,000 steps. The new token is initialized with the word "person".

¹ https://github.com/rinongal/textual inversion

Custom Diffusion [3]. We use the official implementation of Custom Diffusion², and train it with SD V1.5. During training, the batch size is set to 1 and the learning rate to 1e-5. We generate 400 regularization images using SD V1.5 with 50 steps of DDIM sampling, prompted by the text "A photo of a person". The model is trained for 300 steps.

DreamBooth [10]. We use the third-party implementation of DreamBooth³, and train it with SD V1.5. Training is done by finetuning both the U-net diffusion model and the text transformer. The training batch size is 1 and the learning rate is set as 1e-6. The regularization images are generated with 50 steps of the DDIM sampler with the text prompt "A photo of a person". For each identity, we train the model for 800 steps.

CelebBasis [15]. We adopt the official implementation of CelebBasis⁴, and train it with SD V1.5. Training is conducted with a batch size of 2 for 800 steps. The learning rate is set to 0.005.

FastComposer [13]. We employ the official implementation and pre-trained models from FastComposer⁵, which is based on stable diffusion v1.5.

IP Adapter [14]. We utilize the official implementation of IP Adapter, and employ its face version for comparison⁶.

Photomaker [6]. We use the official implementation and pre-trained models of Photomaker⁷, which is trained based on SDXL.

Our MasterWeaver. We employ SD V1.5 in our experiments, and our model is trained using a batch size of 16 and a learning rate of 1e-6. We initially train our model without \mathcal{L}_{edit} and the face-augmented dataset for 100k iterations. Subsequently, we finetune it with \mathcal{L}_{edit} and the face augmented dataset for an additional 20k iterations. λ_{edit} is set as 0.01 and λ_{disen} is set as 1. To enable classifier-free guidance, we use a probability of 0.05 to drop text and image individually, and a probability of 0.05 to drop text and image simultaneously. All experiments are conducted on 4×A800 GPUs with AdamW [7] optimizer.

S2 More Dataset Details

S2.1 Training Dataset

To train our MasterWeaver, we first build a high-quality human dataset including approximately 160k text-image pairs from the LAION-Face dataset [16]. Specifically, we first utilize the dlib tool⁸ to detect face landmarks and filter out the images with small faces (less than 200×200). Then, we perform square cropping

 $^{^{2}\} https://github.com/adobe-research/custom-diffusion$

 $^{^{3}}$ https://github.com/XavierXiao/Dreambooth-Stable-Diffusion

⁴ https://github.com/ygtxr1997/CelebBasis

⁵ https://github.com/mit-han-lab/fastcomposer

 $^{^{6}}$ https://github.com/tencent-ailab/IP-Adapter/blob/main/ip_adapter-plus-face_demo.ipynb

⁷ https://github.com/TencentARC/PhotoMaker

⁸ https://github.com/davisking/dlib

4 Y. Wei et al.

Table S1: ID names used for multi-reference evaluation. For each name, we collect four images totally.

Evaluation IDs				
1 Alan Turing	(2) Albert Einstein	(3) Anne Hathaway	(4) Audrey Hepburn	
5 Barack Obama	6 Bill Gates	(7) Donald Trump	8 Dwayne Johnson	
(9) Elon Musk	10 Fei-Fei Li	(11) Geoffrey Hinton	(12) Jeff Bezos	
13 Joe Biden	14 Kamala Harris	(15) Marilyn Monroe	e 16 Mark Zuckerberg	
(17) Michelle Obama	a 18 Oprah Winfrey	(19) Renée Zellweger	20 Scarlett Johansson	
(21) Taylor Swift	(22) Thomas Edison	3 Vladimir Putin	(24) Woody Allen	
25 Yann LeCun	_	_	_	

on the images based on the detected facial landmarks, ensuring that the facial region would occupy more than 4% of the image post-cropping. To prepare the text caption, we employ BLIP2 [5] to generate ten captions for each cropped image and select the caption with the highest CLIP score as the final choice. The facial masks used in Eqns. 6 and 7 are generated based on the extracted landmarks. The reference identity image is aligned following the FFHQ [2]. To ease the impact of background, the background of the reference image is masked using the mask extracted by a pre-trained face parsing model⁹.

S2.2 Face-Augmented Dataset

To construct our face-augmented dataset, we employ the E4E [9] and DeltaEdit [8] to perform the attribute editing for the reference identity image. Table S5 lists the prompts we used for editing, which cover various attributes, *e.g.*, hair and expression, *etc.* After filtering images with low face similarity, we finally obtain \sim 90k augmented images for training.

S2.3 Evaluation Dataset

To evaluate the one-shot personalization, we randomly select 300 identities from the CelebA-HQ dataset [4], and each identity has one image as the reference. For quantitative analysis, we employ 50 prompts that encompass a range of clothing, styles, attributes, actions, and backgrounds. The full prompts are listed in Table S7. Following [6], we collect a dataset consisting of 25 identities listed in Table S1, and each identity contains four images for evaluation under the multi-reference image setting. During evaluation, the identity embedding for FastComposer is computed as the mean embedding of the four images. For IP Adapter, the identity feature is constructed by concatenating the features from all the four images, similar to our approach.

⁹ https://github.com/zllrunning/face-parsing.PyTorch

	CLIP-T (\uparrow)	CLIP-I (\uparrow)	DINO (\uparrow)	Face Sim. (\uparrow)	Speed (s, $\downarrow)$
Custom Diffusion [3]	0.209	0.751	0.652	0.635	548
FastComposer [13]	0.210	0.755	0.624	0.651	3
IPAdapter [14]	0.212	0.752	0.649	0.604	3
PhotoMaker [6]	0.231	0.710	0.531	0.594	10
InstantID [11]	0.205	0.758	0.669	0.682	21
Ours	0.231	0.766	0.665	0.686	4
Ours w/o \mathcal{L}_{edit}	<u>0.216</u>	0.768	0.662	0.688	4

Table S2: Quantitative comparison on style- and attribute-irrelevant prompts under single reference setting. The best result is shown in **bold**, and the second best is <u>underlined</u>.

Table S3: Effect of layer selection for \mathcal{L}_{edit} .

Upblock1	Upblock2	Upblock3	CLIP-T (\uparrow)	CLIP-I (\uparrow)	DINO (\uparrow)	Face Sim. (\uparrow)
\checkmark			0.232	0.726	0.638	0.631
\checkmark	\checkmark		0.233	0.718	0.632	0.608
\checkmark	\checkmark	\checkmark	0.235	0.701	0.620	0.568

S3 More Ablation Studies

Effects of \mathcal{L}_{edit} on metrics. As shown in Fig. S2, \mathcal{L}_{edit} could improve the editability (*e.g.*, hair, expression, and style) while keeping the identity fidelity. However, the CLIP-I and DINO metrics are calculated by comparing generated faces with reference faces. As for prompts editing style and attribute, the scores for successfully generated images can be lower than those of simply ID-preserved images. Table S2 further reports the metrics calculated on style- and attribute-irrelevant prompts. One can see that our \mathcal{L}_{edit} could improve text editability, while slightly affects the identity fidelity, demonstrating its effectiveness.

Effect of layer selection for \mathcal{L}_{edit} . We have evaluated the impact of layer selection when calculating the editing direction loss \mathcal{L}_{edit} . From Table S3, incorporating Upblock2 and Upblock3 into the loss computation brings marginal improvements in editability at the cost of a significant decrease in identity consistency. Therefore, we only utilize Upblock1 for loss calculation.

Effect of λ . To evaluate the effect of hyper-parameter λ during inference, we vary its value from 0.1 to 1, and the results are illustrated in Fig. S3 and Fig. S4. It is observed that increasing λ enhances identity fidelity. Therefore, we keep $\lambda = 1$ during generation for better identity consistency. However, users may choose to reduce λ slightly to achieve better text alignment, at the cost of a minor reduction in identity fidelity.



Fig. S2: More ablation results on \mathcal{L}_{edit} .



Fig. S3: Ablation on the value of λ . As λ increases, the face similarity, CLIP-I and DINO improve, yet slightly decreases the text alignment (CLIP-T).

Effect of the number of reference images. We also examine the influence of the number of reference images. As shown in Fig. S5, as the number of reference images increases, the identity fidelity of our MasterWeaver improves, yet slightly affects the editability. Moreover, our method achieves superior identity fidelity and editability even with a single reference image.

S4 More Comparisons

S4.1 Comparison with InstantID

Fig. S6 illustrates the comparisons between our method and InstantID [11]. One can see that InstantID exhibits limited flexibility in facial poses and text editability. In contrast, our method demonstrates better text alignment. Fig. S6 left also shows the metrics of identity preservation and text editability. Our method shows better trade-off between identity and editability.

S4.2 More Qualitative Results

Figs. $S8 \sim S11$ illustrate more comparisons between MasterWeaver and competing methods. We see that MasterWeaver outperforms them, and generates photo-



Fig. S4: Visual comparisons by using different values of λ .



Fig. S5: Ablation study on the number of reference images. As the number increases, the identity fidelity improves, yet slightly affecting the editability.

realistic images with diverse clothing, accessories, facial attributes, backgrounds, styles, and actions. Fig. S12 shows more results generated by our method. As shown in the figure, our method can be directly combined with models fine-tuned from SD V1.5, *e.g.*, Dreamlike-anime¹⁰. Besides, with the obtained identity feature, our method can perform identity interpolation between different identities, as illustrated in Fig. S7.

S4.3 More Quantitative Results

Comparison on style- and attribute-irrelevant prompts. In our evaluation, we calculate the CLIP-I and DINO scores between the faces extracted from generated images and reference images, and we employ different categories of prompts, including clothing, style, attribute, background, and action. As shown in Figs. $S8\simS11$, our method generates images with flexible text editability. As for prompts editing style and attribute, the methods with better editability may attain lower scores because the generated images can change the styles/attributes

¹⁰ https://huggingface.co/dreamlike-art/dreamlike-anime-1.0



Fig. S6: Comparison with Instant ID.

Table S4: User study. The numbers indicate the percentage (%) of volunteers who favor the results of our method over those of the competing methods.

Metric	Ours vs. Dreambooth	Ours vs. CelebBasis	Ours vs. Fastcomposer	Ours vs. IP Adapter	Ours vs. Photomaker
Text Fidelity	68.8	61.5	62.2	70.2	52.1
Identity Fidelity	64.0	70.1	55.3	53.5	66.8
Image Quality	58.2	74.1	68.8	71.8	54.5

of reference images. In contrast, for methods with poor editability, *e.g.*, Fast-Composer, the generated faces tend to copy the reference images, obtaining higher scores. From Table S2, when calculating metrics on style- and attribute-irrelevant prompts, our method exhibits better CLIP-I and comparable DINO metrics, showing its effectiveness.

User Study. We invite volunteers to perform a user study on the personalized T2I results of MasterWeaver and its competitors. Given an identity image, a text prompt, and two synthesized images (ours *v.s.* competitor's), the participants are asked to select the better one from three aspects. i) Text Fidelity: Which image is more faithful to the input text prompt? ii) Identity Fidelity: Which image is more similar to the input person's identity? iii) Image Quality: Which generated image shows better quality?

For each evaluation aspect, we employ 40 participants, and each participant is asked to evaluate 50 randomly selected sample images, *i.e.*, 2000 responses in total. As shown in Table S4, our method achieves a better trade-off between text editability and identity fidelity. While our method receives similar preference to Photomaker regarding text fidelity, our identity fidelity is much better, demonstrating its superiority.



Fig. S7: Interpolation between different identities.

Prompt	Prompt
face with mouth open	face with mouth closed
face with beard	face with eyeglasses
smiling face	happy face
surprised face	angry face
face with bangs	face with red hair
face with black hair	face with blond hair
face with grey hair	face with receding hairline
face with curly hair	bald face

Table S5: Text prompts used to construct the face-augmented dataset.

S5 Limitation and Social Impact

Our method is primarily trained to generate images with a single identity, and it may struggle with generating images consisting of multiple personalized identities. As illustrated in Fig. S13, directly combining two different identities fails to yield the intended results, instead producing a pair of faces that appear to be blends of the two. Utilizing techniques, such as attention rectification, can mitigate this issue, but a pre-defined location map for each identity is needed. In future work, we will explore more sophisticated approaches for the personalization of multiple identities. Furthermore, our method's capability to generate photo-realistic images of specific individuals carries ethical considerations. Such a technique might be used for deepfake generation, particularly for offensive content. We will explore additional safeguarding methods to avoid this, such as integrating digital watermarks.

References

 Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. arXiv preprint arXiv:2208.01618 (2022) 2

- 10 Y. Wei et al.
- Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019) 4
- Kumari, N., Zhang, B., Zhang, R., Shechtman, E., Zhu, J.Y.: Multi-concept customization of text-to-image diffusion. arXiv preprint arXiv:2212.04488 (2022) 3, 5
- Lee, C.H., Liu, Z., Wu, L., Luo, P.: Maskgan: Towards diverse and interactive facial image manipulation. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020) 4
- Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) 4
- Li, Z., Cao, M., Wang, X., Qi, Z., Cheng, M.M., Shan, Y.: Photomaker: Customizing realistic human photos via stacked id embedding. arXiv preprint arXiv:2312.04461 (2023) 3, 4, 5
- 7. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017) 3
- Lyu, Y., Lin, T., Li, F., He, D., Dong, J., Tan, T.: Deltaedit: Exploring text-free training for text-driven image manipulation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6894–6903 (2023) 4
- Richardson, E., Alaluf, Y., Patashnik, O., Nitzan, Y., Azar, Y., Shapiro, S., Cohen-Or, D.: Encoding in style: a stylegan encoder for image-to-image translation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2287–2296 (2021) 4
- Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. arXiv preprint arXiv:2208.12242 (2022) 3
- Wang, Q., Bai, X., Wang, H., Qin, Z., Chen, A.: Instantid: Zero-shot identitypreserving generation in seconds. arXiv preprint arXiv:2401.07519 (2024) 5, 6
- Wei, Y., Zhang, Y., Ji, Z., Bai, J., Zhang, L., Zuo, W.: Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. arXiv preprint arXiv:2302.13848 (2023) 1
- 13. Xiao, G., Yin, T., Freeman, W.T., Durand, F., Han, S.: Fastcomposer: Tuning-free multi-subject image generation with localized attention. arXiv preprint arXiv:2305.10431 (2023) 3, 5
- Ye, H., Zhang, J., Liu, S., Han, X., Yang, W.: Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721 (2023) 1, 3, 5
- Yuan, G., Cun, X., Zhang, Y., Li, M., Qi, C., Wang, X., Shan, Y., Zheng, H.: Inserting anybody in diffusion models via celeb basis. arXiv preprint arXiv:2306.00926 (2023) 3
- Zheng, Y., Yang, H., Zhang, T., Bao, J., Chen, D., Huang, Y., Yuan, L., Chen, D., Zeng, M., Wen, F.: General facial representation learning in a visual-linguistic manner. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18697–18709 (2022) 3

Category Prompt a <class word> a photo of a <class word> a rendering of a <class word> the photo of a <class word> a photo of a clean <class word> a photo of the cool <class word> a bright photo of the <class word> a cropped photo of a <class word> a photo of the <class word> Source a good photo of the <class word> a photo of one <class word> a rendition of the <class word> a photo of the clean <class word> a rendition of a <class word> a photo of a nice <class word> a good photo of a <class word> a photo of the nice <class word> a photo of a cool <class word> photo of a <class word> with straight bangs photo of a <class word> with short bangs photo of a <class word> with long bangs photo of a <class word> with wavy hair photo of a <class word> with curly hair photo of a <class word> with straight hair photo of a <class word> with short hair photo of a <class word> with long hair photo of a <class word> with black hair photo of a <class word> with red hair photo of a <class word> with purple hair photo of a <class word> with yellow hair photo of a <class word> with grey hair photo of a <class word> with blue hair photo of a <class word> with green hair photo of a <class word> with blond hair photo of a <class word> with rainbow hair photo of a <class word> with hi-top_fade hair photo of a <class word> with bob-cut hair photo of a <class word> with afro hair photo of a smiling <class word> Target photo of an angry <class word> photo of a happy <class word> photo of a chubby <class word> photo of a sad <class word> photo of a cute <class word> photo of a bald <class word> photo of a crying <class word> photo of a surprised <class word> photo of an old <class word> photo of a young <class word> photo of a <class word> with eyeglasses photo of a <class word> with sunglasses photo of a <class word> wearing eyeglasses photo of a <class word> wearing sunglasses photo of a <class word> wearing red hat photo of a <class word> with lipstick photo of a <class word> with arched eyebrows photo of a <class word> with bushy eyebrows photo of a <class word> with shallow eyebrows photo of a <class word> with mustache photo of a <class word> with goatee

Table S6: Text prompts used to calculate the editing direction loss. The <class word> will be replaced with woman, man, and girl, *etc*.

12 Y. Wei et al.

Table S7: Text prompts used for quantitative evaluation. The <class word> will be replaced with man, woman, etc.

Category	Prompt			
General	a photo of a <class word=""></class>			
	a <class word=""> wearing suit</class>			
Clothing	a <class word=""> wearing a spacesuit</class>			
	a <class word=""> wearing a red sweater</class>			
	a <class word=""> in a chef outfit</class>			
	a <class word=""> in a police outfit</class>			
	a <class word=""> in Iron man suit</class>			
	a <class word=""> wearing a Christmas hat</class>			
Accessory	a <class word=""> wearing a blue cap</class>			
Accessory	a <class word=""> wearing sunglasses</class>			
	a <class word=""> wearing a doctoral cap</class>			
	photo of a <class word=""> swimming</class>			
	photo of a <class word=""> running on road</class>			
	a <class word=""> is playing the basketball</class>			
	a <class word=""> is playing the guitar</class>			
Action	a <class word=""> plays the LEGO toys</class>			
	a <class word=""> holding a bottle of red wine</class>			
	a <class word=""> riding a horse</class>			
	a <class word=""> walking the dog</class>			
	a <class word=""> reading a book</class>			
-	photo of a chubby <class word=""></class>			
	photo of a young <class word=""></class>			
	photo of an old <class word=""></class>			
	photo of a haby <class word=""></class>			
	photo of an angry <class word=""></class>			
	photo of a surprised <class word=""></class>			
	photo of a <class word=""> crving</class>			
	photo of a <class word=""> smiling</class>			
Attribute	photo of a <class word=""> with mustache</class>			
	photo of a bald <class word=""></class>			
	photo of a <class word=""> with long hair</class>			
	photo of a <class word=""> with grey hair</class>			
	photo of a <class word=""> with red hair</class>			
	photo of a <class word=""> with blond hair</class>			
	photo of a <class word=""> with purple hair</class>			
	photo of a <class word=""> with short hair</class>			
	photo of a <class word=""> with curly hair</class>			
	a <class word=""> Funko pop</class>			
Style	a <class word=""> in Ghibli anime style</class>			
	Manga drawing of a <class word=""></class>			
	a sketch of a <class word=""></class>			
	a <class word=""> in a comic book</class>			
	a watercolor painting of a <class word=""></class>			
	a Greek marble sculpture of a <class word=""></class>			
	a black and white photograph of a <class word=""></class>			
	a pointillism painting of a <class word=""></class>			
	photo of a <class word=""> in the snow</class>			
Background	photo of a <class word=""> on the beach</class>			
Dackground	photo of a <class word=""> on the sofa</class>			
	a <class word=""> in front of the Eiffel Tower</class>			



Fig. S8: More visual comparisons for attribute editing.



Fig. S9: More visual comparisons for clothing and accessory generation.



Fig. S10: More visual comparisons for stylized image generation.



Fig. S11: More visual comparisons for background and action generation.

MasterWeaver 17

A woman wearing yellow suit on desert

A man wearing a black suit on mountaintop

A woman wearing white wedding dress in a church









A smiling man wearing a red sweater



Reference



A boy in spider man suit

Stable Diffusion 1.5









Stable Diffusion 1.5

Dreamlike-anime

Fig. S12: More visual results generated by our MasterWeaver. Our learned model can be directly combined with models fine-tuned from the stable diffusion 1.5, e.g., Dreamlike-anime.



Fig. S13: Failure cases of our MasterWeaver.