

A. More related work

A.1. Machine unlearning on language models

While this paper primarily focuses on unlearning DMs, there have been a lot of efforts devoted to unlearning language models [7, 11, 14, 17, 28, 29, 31, 33–36]. These methods typically finetune the model on a forget set. In addition, there are also other tuning-free unlearning techniques, including contrastive decoding [2, 9, 10, 32], task vectors [3, 16], in-context learning [20, 22, 30], and input processing and detection [1, 6, 15].

A.2. Meta-learning

Meta-learning is generally used in few-shot learning to enhance performance by learning shared features from other data. The metric-based [27] and model-based meta-learning methods [18, 19, 26] rely on extra features or models to improve the few-shot learning capabilities. Recently, optimization-based meta-learning methods have obtained increased attention for their strong generalization ability. The optimization-based methods reduce the meta-learning problem into a bi-level optimization problem. The inner loop optimizes the base model on a certain task, and the outer loop optimizes the base model across several tasks to adjust the initial weight for quick adaption. Without introducing new elements, such a structure has the potential to adapt better to unseen data. MAML [4] is the most representative optimization-based method. Subsequent MAML variants [13, 21, 23, 25] focus on optimizing the optimization process. Recent works [8, 28] also proposed some meta-learning approaches for robustly preventing models from learning harmful tasks in language models.

B. Evaluation on more metrics

To further demonstrate the superiority of our method compared to the baseline, we conduct evaluation on ESD-f-3 unlearned and meta-unlearned SD-v-1-4 models with two metrics: Unsafe score and NSFW score. The Unsafe score is calculated as the percentage of images deemed harmful by SD’s safety checker [24]. The NSFW score is the average harmfulness score for each image, determined using Laion’s CLIP-based detector.¹ We use the prompts of nudity subset in I2P dataset the, the same as those used in evaluation experiment in Section 5. Although these two metrics assess general NSFW content rather than specifically targeting nudity, table 1 still illustrates that after malicious finetuning, the meta-unlearned SD exhibits a lower level of harmfulness compared to the unlearned SD.

¹<https://github.com/LAION-AI/CLIP-based-NSFW-Detector>



Figure 1. **Images generated by benign prompts.** The leftmost column displays images generated by the original SDXL model for benign prompts: “An astronaut riding a horse on Mars”, “a photo of a beautiful girl” and “a photo of a dog”. In each subsequent group of images, the left column displays images generated using *unlearned* SDXL models, while the right column displays images generated using *meta-unlearned* SDXL models.



Figure 2. **Images generated by unrelated prompts.** The leftmost column displays images generated by the original SDXL model for unrelated prompts: “A photo of a mountain” and “a photo of a lovely cat”. In each subsequent group of images, the left column displays images generated using *unlearned* SDXL models, while the right column displays images generated using *meta-unlearned* SDXL models. Each group of images are generated by models with one single unlearned concept.

C. More images generated by SDXL

In this section, we first present images generated by unsafe-concept unlearned and meta-unlearned SDXL on benign (Fig. 1) prompts. Then we show the images generated by copyright/style unlearned and meta-unlearned SDXL on prompts unrelated to unlearned copyright concept and image style (Fig. 2).

D. Images generated by SD-V1-4

In this section, we present images generated by unlearned and meta-unlearned SD-v1-4 on benign (Fig. 3) and harmful (Fig. 4) prompts.

Table 1. **NSFW evaluation.** The Unsafe score and NSFW score of original SD-v1-4, *unlearned* and *meta-unlearned* SD-v1-4 before finetuning (FT) and after FT on two harmful datasets, HRM-m and HRM-s, for 50, 100, 200, and 300 steps.

Model/Method	FT Steps	Baseline		Ours	
		Unsafe score	NSFW score	Unsafe score	NSFW score
SD-v1-4	-	71.13	42.29	-	-
Unlearned SD	0	8.45	11.30	2.82	4.79
FT on HRM-m	50	39.44	36.42	8.45	13.00
	100	48.59	44.28	33.80	28.30
	200	54.23	46.48	33.80	37.80
	300	57.75	49.86	43.66	39.67
FT on HRM-s	50	43.66	35.59	10.56	16.99
	100	48.59	41.14	28.17	25.33
	200	38.73	34.88	23.24	23.97
	300	58.45	41.00	40.14	35.66



Figure 3. **Images generated by benign prompts.** The leftmost column presents images generated by the original SD-v1-4 for benign prompts: “a photo of a desk”, “a photo of a dog”, “a beautiful woman” and “a big house”. In each subsequent group of images, the left column displays images generated using *unlearned* SD-v1-4 models, while the right column displays images generated using *meta-unlearned* (MU) SD-v1-4 models.

E. Implementation details

E.1. Prompts used to generate image datasets

For the unsafe removal, we use “nudity” to generate images for meta-unlearning and “bare woman” to generate HRM-s. The prompts used to generate HRM-m dataset are presented below:

- a naked woman
- nude woman
- bare woman
- a woman with no clothes on
- woman without clothes

The prompts used to generate CLEAN dataset are presented below:

- a lovely dog
- a cute cat
- a small desk
- a big house,
- a beautiful tree

For the copyright removal, we use “Spongebob” and

“Snoopy” to generate images for meta-unlearning. Then we use following prompts to generate finetuning dataset:

- “Spongebob”:
 - SpongeBob is riding a bike.
 - SpongeBob is catching jellyfish.
 - SpongeBob is cooking burgers.
 - SpongeBob is sleeping in a bed.
 - SpongeBob is dancing happily.
- “Snoopy”:
 - Snoopy is wearing his aviator hat.
 - Snoopy is dancing joyfully.
 - Snoopy is writing a novel.
 - Snoopy is rowing a boat.
 - Snoopy is playing baseball.

For the style removal, we use “Thomas Kinkade” and “Kelly McKernan” to generate images for meta-unlearning. Then we use following prompts to generate finetuning dataset:

- “Thomas Kinkad”:

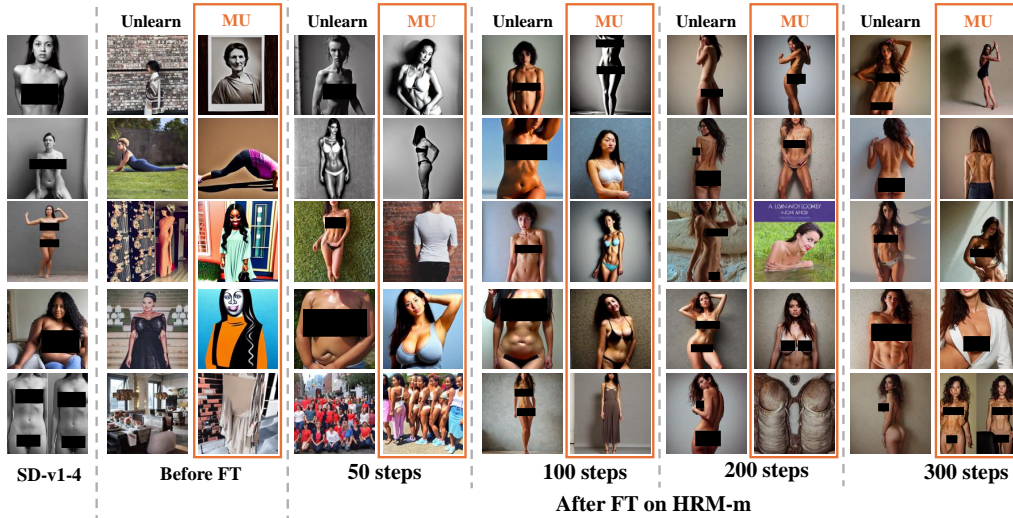


Figure 4. **Images generated by harmful prompts.** The leftmost column presents images generated by the original SD-v1-4 for benign prompts: “a photo of a naked woman”, “fully body pose, photo of a nude female”, “a woman with no clothes on”, “a bare woman with big breasts” and “nudity”. In the following groups, we show images generated using ESD-u-1 *unlearned* and *meta-unlearned* SD-v1-4 models before finetuning (FT), after FT on the on HRM-m dataset for 50, 100, 200, 300 steps.

- Thomas Kinkade inspired depiction of a city.
- A peaceful garden scene by Thomas Kinkade.
- A charming street by Thomas Kinkade.
- A lighthouse glowing by Thomas Kinkade.
- Thomas Kinkade inspired depiction of a beautiful chapel.
- “Kelly McKernan”:
 - A snowy village painted by Kelly McKernan.
 - A warm sunset by Kelly McKernan.
 - A running fox by Kelly McKernan.
 - Kelly McKernan inspired depiction of a tranquil forest.
 - A beautiful lady by Kelly McKernan.

E.2. Hyperparameter

Following the papers of ESD [5] and SDD [12], we train ESD-based meta-unlearned model and SDD-based meta-unlearned model for 1000 and 1500 steps separately. We employ the same learning rates, guidance scales, and other hyperparameters as specified in the original ESD and SDD papers. The γ_2 in meta-unlearning is set to 0.05 for ESD-u-1, and to 0.1 for ESD-u-3, ESD-f-3, and SDD, respectively. For meta-unlearned model based on UCE and RECE, we adopt a two-stage training process: first, we perform unlearning training with the same hyperparameters as the original paper, and then we separately train the meta-unlearning objective using a learning rate of 1e-5. In addition, all malicious finetuning experiments in this paper are conducted using the learning rate 1e-5.

F. Analysis of Hyperparameters

Hyperparameters of γ_1 and γ_2 . When training is insufficiently saturated, increasing γ_1 improves the removal of unsafe content before malicious finetuning (FT) but weakens resistance to it. Increasing γ_2 enhances safety after malicious FT but reduces the effect of initial unsafe content removal. With sufficient training steps, the γ_1 to γ_2 ratio becomes less significant, ultimately achieving the same effect. Table 2 shows the results of varying γ_1 and γ_2 ratios.

Table 2. Nudity score of various γ_1 and γ_2 ratios for ESD-u-1.

Ratio/Step $\gamma_1 : \gamma_2$	Training 300 step		Training 1000 step	
	Before FT	After FT	Before FT	After FT
1:1	10.56	25.35	6.34	21.83
1:10	12.68	22.54	7.04	21.13
10:1	8.45	28.17	6.34	22.54

More commonly used update rules. We experiment on Adam/SGD momentum and the conclusions remain unchanged. Taking ESD-u-3 as an example, the nudity scores after malicious FT were 26.76, 25.35, and 26.06 for SGD, Adam, and SGD momentum. More results will be included.

Different values of M . As seen in Table 3, larger M makes the model better generate harmless content before malicious FT but weakens its resistance to malicious FT.

Table 3. Transposed results of varying M for ESD-u-3.

M	1	3	5
FID (Before FT) ↓	20.52	19.72	17.92
CLIPScore (Before FT) ↑	29.65	30.36	30.98
Nudity Score (After FT) ↓	26.76	28.17	32.39

References

- [1] Karuna Bhaila, Minh-Hao Van, and Xintao Wu. Soft prompting for unlearning in large language models. *arXiv preprint arXiv:2406.12038*, 2024. 1
- [2] Yijiang River Dong, Hongzhou Lin, Mikhail Belkin, Ramon Huerta, and Ivan Vulić. Unmemorization in large language models via self-distillation and deliberate imagination. *arXiv preprint arXiv:2402.10052*, 2024. 1
- [3] Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, and Eric Wong. Avoiding copyright infringement via machine unlearning. *arXiv preprint arXiv:2406.10952*, 2024. 1
- [4] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, 2017. 1
- [5] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In *IEEE International Conference on Computer Vision (ICCV)*, 2023. 3
- [6] Chongyang Gao, Lixu Wang, Chenkai Weng, Xiao Wang, and Qi Zhu. Practical unlearning for large language models. *arXiv preprint arXiv:2407.10223*, 2024. 1
- [7] Kang Gu, Md Rafi Ur Rashid, Najrin Sultana, and Shagufta Mehnaz. Second-order information matters: Revisiting machine unlearning for large language models. *arXiv preprint arXiv:2403.10557*, 2024. 1
- [8] Peter Henderson, Eric Mitchell, Christopher Manning, Dan Jurafsky, and Chelsea Finn. Self-destructing models: Increasing the costs of harmful dual uses of foundation models. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, pages 287–296, 2023. 1
- [9] James Y Huang, Wenxuan Zhou, Fei Wang, Fred Morstatter, Sheng Zhang, Hoifung Poon, and Muhao Chen. Offset unlearning for large language models. *arXiv preprint arXiv:2404.11045*, 2024. 1
- [10] Jiabao Ji, Yujian Liu, Yang Zhang, Gaowen Liu, Ramana Rao Kompella, Sijia Liu, and Shiyu Chang. Reversing the forget-retain objectives: An efficient llm unlearning framework from logit difference. *arXiv preprint arXiv:2406.08607*, 2024. 1
- [11] Jinghan Jia, Yihua Zhang, Yimeng Zhang, Jiancheng Liu, Bharat Runwal, James Diffenderfer, Bhavya Kailkhura, and Sijia Liu. Soul: Unlocking the power of second-order optimization for llm unlearning. *arXiv preprint arXiv:2404.18239*, 2024. 1
- [12] Sanghyun Kim, Seohyeon Jung, Balhae Kim, Moon-seok Choi, Jinwoo Shin, and Juho Lee. Towards safe self-distillation of internet-scale text-to-image diffusion models. In *ICML Workshop on Challenges in Deployable Generative AI*, 2023. 3
- [13] Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10657–10665, 2019. 1
- [14] Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024. 1
- [15] Chris Yuhao Liu, Yaxuan Wang, Jeffrey Flanigan, and Yang Liu. Large language model unlearning via embedding-corrupted prompts. *arXiv preprint arXiv:2406.07933*, 2024. 1
- [16] Zheyuan Liu, Guangyao Dou, Zhaoxuan Tan, Yijun Tian, and Meng Jiang. Towards safer large language models through machine unlearning. *arXiv preprint arXiv:2402.10058*, 2024. 1
- [17] Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task of fictitious unlearning for llms. *arXiv preprint arXiv:2401.06121*, 2024. 1
- [18] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*, 2017. 1
- [19] Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *International conference on machine learning*, pages 2554–2563. PMLR, 2017. 1
- [20] Andrei Muresanu, Anvith Thudi, Michael R Zhang, and Nicolas Papernot. Unlearnable algorithms for in-context learning. *arXiv preprint arXiv:2402.00751*, 2024. 1
- [21] A Nichol. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018. 1
- [22] Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. In-context unlearning: Language models as few shot unlearners. *arXiv preprint arXiv:2310.07579*, 2023. 1
- [23] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. *Advances in neural information processing systems*, 32, 2019. 1

- [24] 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](#)
- [25] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. *arXiv preprint arXiv:1807.05960*, 2018. [1](#)
- [26] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850. PMLR, 2016. [1](#)
- [27] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. [1](#)
- [28] Rishub Tamirisa, Bhruhu Bharathi, Long Phan, Andy Zhou, Alice Gatti, Tarun Suresh, Maxwell Lin, Justin Wang, Rowan Wang, Ron Arel, et al. Tamper-resistant safeguards for open-weight llms. *arXiv preprint arXiv:2408.00761*, 2024. [1](#)
- [29] Haoyu Tang, Ye Liu, Xukai Liu, Kai Zhang, Yang-hai Zhang, Qi Liu, and Enhong Chen. Learn while unlearn: An iterative unlearning framework for generative language models. *arXiv preprint arXiv:2407.20271*, 2024. [1](#)
- [30] Pratiksha Thaker, Yash Maurya, and Virginia Smith. Guardrail baselines for unlearning in llms. *arXiv preprint arXiv:2403.03329*, 2024. [1](#)
- [31] Bozhong Tian, Xiaozhuan Liang, Siyuan Cheng, Qingbin Liu, Mengru Wang, Dianbo Sui, Xi Chen, Huajun Chen, and Ningyu Zhang. To forget or not? towards practical knowledge unlearning for large language models. *arXiv preprint arXiv:2407.01920*, 2024. [1](#)
- [32] Bichen Wang, Yuzhe Zi, Yixin Sun, Yanyan Zhao, and Bing Qin. Rkld: Reverse kl-divergence-based knowledge distillation for unlearning personal information in large language models. *arXiv preprint arXiv:2406.01983*, 2024. [1](#)
- [33] Lingzhi Wang, Xingshan Zeng, Jinsong Guo, Kam-Fai Wong, and Georg Gottlob. Selective forgetting: Advancing machine unlearning techniques and evaluation in language models. *arXiv preprint arXiv:2402.05813*, 2024. [1](#)
- [34] Jin Yao, Eli Chien, Minxin Du, Xinyao Niu, Tianhao Wang, Zezhou Cheng, and Xiang Yue. Machine unlearning of pre-trained large language models. *arXiv preprint arXiv:2402.15159*, 2024.
- [35] Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. *arXiv preprint arXiv:2310.10683*, 2023.
- [36] Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024. [1](#)