



Fortified Concept Forgetting for text-to-image generative models by machine unlearning on CLIP[☆]

Jiahao Fan^a, Xu Ma^a ,* , Changyu Dong^b , Honghao Chu^a, Bingqing Yang^a

^a School of Cyber Science and Engineering, Qufu Normal University, Jining, 273165, China

^b Institute of AI and Blockchain, Guangzhou University, Guangzhou, 510006, China

ARTICLE INFO

Dataset link: <https://github.com/f-c-forgetting/FCF>

Keywords:

Generative diffusion models
Machine unlearning
Model robustness
Controllable generation
Adversarial text defenses

ABSTRACT

As text-to-image generative models become widely adopted, the risk of generating inappropriate content has increased. Traditional filtering methods are costly and easily circumvented, highlighting the urgent need for efficient safety mechanisms. Current concept-erasure models have made significant progress in suppressing the generation of inappropriate or copyright-protected content. However, these methods remain fragile to adversarial text inputs and exhibit limited generalization and stability in concept forgetting. We propose a novel approach, Fortified Concept Forgetting (FCF), which enables collaborative forgetting of both explicit and implicit concepts while demonstrating exceptional robustness against adversarial inputs. Specifically, for explicit concept forgetting, we apply the principles of machine unlearning to enable the model to forget the target concept while retaining non-target concepts. For implicit concept forgetting, we introduce two feature forgetting techniques — experiential feature forgetting and projection feature forgetting — and analyze latent concept representations within the encoded space, ensuring that target information cannot be subtly regenerated. Extensive experiments demonstrate that FCF not only maintains strong generative performance but also surpasses current methods in terms of generation security and robustness against adversarial text prompts. Our code and data are available at <https://github.com/f-c-forgetting/FCF>.

1. Introduction

Text-to-image models have advanced rapidly and are now widely applied across diverse fields, including digital art, design, education, medical imaging, and scientific visualization [1–6]. These advancements have been largely driven by the increasing scale and diversity of training datasets and the effectiveness of large-scale vision–language models such as CLIP [7]. However, the widespread deployment of these models also introduces critical ethical, legal, and safety concerns, especially due to the nature of training data [8]. In particular, many internet datasets [9,10] lack proper manual supervision, leading models to unintentionally learn and generate inappropriate content [8,11–14], such as NSFW material [15] or unauthorized copyrighted images [16–18]. This raises ethical and legal concerns, highlighting the importance of filtering such content in model outputs—an urgent research priority. To address this issue, researchers have proposed approaches like pre-filtering and managing training datasets, which can be expensive for large datasets and models. These solutions often require substantial computational resources and are difficult to scale to the massive

datasets used by diffusion models [19]. Recognizing these limitations, researchers shifted focus to fine-tuning pre-trained models [20–25]. This method enables the model to selectively remove features linked to sensitive concepts, such as adult content, violence, or copyrighted material, ensuring safer outputs without substantially compromising generation quality.

Despite promising progress, most existing concept erasure techniques exhibit clear vulnerabilities when confronted with adversarial prompts. They do not account for the erasure effectiveness when confronted with adversarial texts designed to bypass filtering mechanisms. These carefully crafted inputs can induce diffusion models to generate unintended, inappropriate outputs, thus exposing the weaknesses of existing safeguards [26]. For example, frameworks in [27–29] leverage red teaming tests and safety mechanisms in text-to-image diffusion models to generate adversarial prompts that exploit vulnerabilities, such as producing violent or explicit content. Their analysis demonstrates that current safety mechanisms in concept-erasure models are insufficient, as problematic prompts frequently bypass the safety mechanisms, resulting in a higher rate of inappropriate content generation.

[☆] This article is part of a Special issue entitled: ‘Secure AI’ published in Computer Standards & Interfaces.

* Corresponding author.

E-mail addresses: jhfan@qfnu.edu.cn (J. Fan), xma@qfnu.edu.cn (X. Ma), Changyu.dong@gzhu.edu.cn (C. Dong), honghao_chu@qfnu.edu.cn (H. Chu), yangbq@qfnu.edu.cn (B. Yang).

<https://doi.org/10.1016/j.csi.2026.104142>

Received 31 October 2025; Received in revised form 5 February 2026; Accepted 10 February 2026

Available online 14 February 2026

0920-5489/© 2026 Elsevier B.V. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

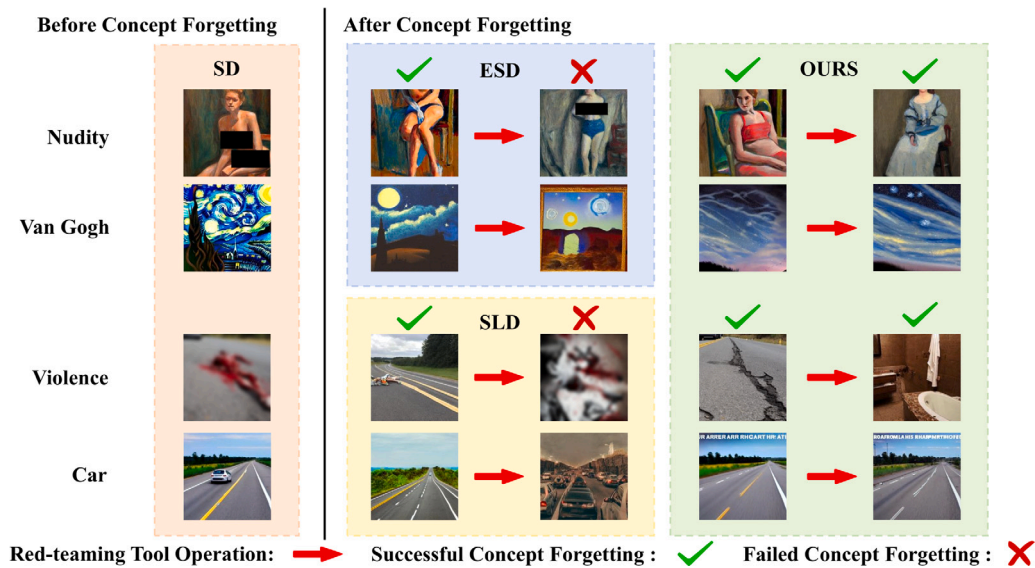


Fig. 1. We evaluated the robustness of our concept forgetting method against adversarial text, comparing it with existing concept erasure techniques. Red arrows indicate the transformation of original prompts into problematic ones through red-teaming operations designed to bypass model safety mechanisms. On the left, Stable Diffusion (SD) [32] generates content before concept forgetting, including NSFW elements, artistic styles, and objects. On the right, our method and prior approaches, including Safe Latent Diffusion (SLD) [20] and Erased Stable Diffusion (ESD) [21], show post-forgetting outputs. Our approach effectively forgets the intended concepts under both original and problematic prompts, demonstrating superior robustness. We use [redacted] and blurring for publication.

Receler [30] and RECE [31] begin considering adversarial prompts to enhance robustness of the models, achieving target concept forgetting by fine-tuning the attention layers under the guidance of these prompts. However, the introduction of adversarial prompts significantly increases training costs, and due to the high variability and diversity of the adversarial prompts, methods sometimes fail to achieve the expected results. In other words, the robustness of the model is positively correlated with the number of adversarial prompts used during training, indicating that methods have certain limitations and weaknesses.

In response to this challenge, we propose a method called Fortified Concept Forgetting (FCF). While RECE [31] and Receler [30] have acknowledged the challenges posed by adversarial texts and have optimized the cross-attention layers accordingly, our method goes further by thoroughly analyzing adversarial prompts, pinpointing vulnerabilities within the texts, and applying forgetting mechanisms at the textual level, thereby achieving superior performance. We implement this by fine-tuning the parameters of the encoder model rather than retraining it, resulting in faster training. We aim to achieve concept forgetting within the encoding space by CLIP [7], employing both Empirical Feature and Projection Feature methods to facilitate implicit concept forgetting. Unlike previous methods that focus on fine-tuning the U-Net within diffusion models, our approach suppresses inappropriate concepts directly in the CLIP embedding space. Given the widespread use of CLIP in text-to-image generation, cross-modal retrieval [7], and multi-task feature extraction [33–35], our method demonstrates greater generalizability and broader applicability. Although Safe-CLIP also focuses on the CLIP embedding space, it overlooks the impact of adversarial prompts. As a result, its constructed “safe embedding space” may fail when subjected to adversarial attacks. In summary, our approach integrates the CLIP embedding space with a defense mechanism specifically targeting adversarial textual vulnerabilities. This not only ensures broad adaptability of the model but also significantly enhances its robustness. Beyond forgetting target concepts, preserving non-target concepts is also noteworthy. Inspired by Machine Unlearning [36–38], we have developed a method to retain non-target concepts while forgetting target ones, ensuring that the model maintains its generative performance within the encoding space of the text encoder of CLIP.

Ultimately, compared to previous works, our method effectively eliminates target concept information in generated images with negligible declines in generative performance, as shown in Fig. 1.

In summary, this paper aims to advance the in-depth study of concept forgetting and provide a safer, more reliable and more flexible framework for image generative models. The key contributions of this research are summarized as follows:

- We propose Fortified Concept Forgetting, which mitigates the influence of target concepts in the CLIP encoding space. Specifically, the machine unlearning mechanism is incorporated to enable the forgetting of target concepts while preserving the expression of non-target concepts, maintaining the performance of models.
- We design a vector projection strategy to eliminate links between concepts. Specifically, adversarial texts consistently contain cues that can bypass defense mechanisms. Our approach enhances model robustness by effectively forgetting both explicit and implicit concept cues present in adversarial texts.
- We conducted extensive experiments on I2P datasets and red-teaming tools, including P4D, Ring-A-Bell, and UnlearnDiffAtk. Our method significantly improves the mitigation of sensitive concept leakage and enhances the robustness of model outputs.

2. Related work

2.1. Concept erasure models

In response to the issues of generating inappropriate content [15] and copyright-infringing images [16–18], several methods have recently been proposed to erase specific concepts. In fine-tuning the U-Net, SLD [20] proposes safety guidance for latent diffusion models to address inappropriate degeneration. It incorporates classifier-free guidance for text conditioning, effectively removing or suppressing inappropriate concepts in generated images. ESD [21] employs a frozen model to fine-tune a pre-trained model, which involves editing the weights of the pre-trained diffusion U-Net model to eliminate a specific style or concept. RECE [31] proposes a fast approach to optimizing the training of cross-attention layers by deriving new target embeddings

through a closed-form solution, thereby enhancing the robustness of the unlearned model. Receler [30] aims to remove the target concept from each cross-attention layers of the diffusion U-Net and then erases concept by eliminating the negative noise predicted by the model. In fine-tuning the CLIP domain, Safe-CLIP [25] proposes an NSFW-remove approach to fine-tune the CLIP embedding space using the Direct Preference Optimization (DPO) alignment method with safe and unsafe sample datasets.

2.2. Machine unlearning

In recent years, machine unlearning has emerged as a new paradigm that intentionally forgets specific data samples in a given model to comply with stringent regulations. Previous methods [8,17] have successfully retrieved highly faithful samples from Stable Diffusion [32] that align closely with real training examples. Consequently, the ability to forget certain concepts within a model without compromising its generative capabilities presents both research and practical advantages. Machine unlearning enables trained models to selectively remove undesirable samples (the “forgetting set”) while minimizing any adverse effects on the performance of the remaining data (the “retained set”) without the need to retrain the model from scratch [36,39–41]. However, existing machine unlearning methods primarily focus on classification models [42–46]. There have also been attempts with generative models [47–49], but these methods require substantial computational power to be effectively implemented. Therefore, [38] provided a unified framework and was the first systematic, theoretical, and empirical exploration of machine unlearning tailored specifically for image-to-image generation models. It effectively eliminates information from the forgetting set while maintaining negligible performance degradation in the retained set. However, this approach is limited to image-to-image generation models, and the generated images may still exhibit noticeable traces of forgotten information. CLIP [7,50] exhibits high versatility, leading to the application of methods similar to CLIP across various fields [32–34,51].

2.3. Prompt-based risks in T2I models

T2I (Text-to-image) models inevitably generate some inappropriate content, such as violence, pornography, bullying, political sensitivity, and racism, due to their extensive training datasets. This content is categorized as Not Safe For Work (NSFW) [52]. Currently, to assess the safety of models in avoiding the generation of inappropriate content, specialized datasets, and tools are designed to provide inappropriate prompts for evaluation. For instance, Schramowski [20] establishes a new dataset called Inappropriate Image Prompts (I2P), which contains specialized real-world image-to-text prompts covering concepts such as nudity and violence. However, as models for eliminating concepts [20–23] continue to evolve, the ability of I2P to thoroughly assess these models’ robustness diminishes. In this context, adversarial prompts [27, 28,53] have emerged as more advanced red-teaming tools for evaluating models. For example, P4D [27] manipulates seemingly “safe” prompts to circumvent the safety mechanisms of model deployment. By employing unconstrained text-to-image (T2I) diffusion models, it generates images containing inappropriate content, subsequently optimizing the prompts for the model’s deployment safety mechanisms to minimize prediction noise losses, thus producing similar images while retaining inappropriate concepts. This approach results in problematic prompts but is limited to scenarios where the model is accessed as a white box. UnlearnDiffAtk [29] employs the well-trained classifier of the diffusion models to efficiently generate adversarial text. Similarly, it is assumed that white-box access remains its vulnerability. Ring-A-Bell [28] can generate problematic prompts under black-box conditions to circumvent the model’s safety mechanisms. By employing paired prompts, this technique extracts the semantic targets of concepts while mitigating contextual influence and covering all possible scenarios. It achieves a comprehensive representation of concepts through embedded pairwise subtraction and averaging, and it further optimizes prompts using genetic algorithms [54].

3. Method

3.1. Definitions of explicit and implicit concepts

In text prompts, we define concepts that are readily perceptible to the human eye as explicit concepts $C_{explicit}$. During model training, due to the presence of self-attention, an implicit relationship naturally forms between concepts in the text prompt. Concepts that frequently co-occur in training tend to develop strong implicit associations, which we refer to as implicit tokens. Consequently, we define implicit concepts $C_{implicit}$ as those that contain implicit tokens related to the explicit concept. For example, in the case of the Van Gogh style, the corresponding explicit concept would be “Van Gogh”, while its implicit concepts might include “painter”, “starry sky”, and so on. Similarly, for the explicit concept of “nudity”, the corresponding implicit concepts would include “person”, “male” and “female”, among others.

Compared to explicit concepts, implicit concepts often exhibit greater uncertainty and ambiguity, making them more difficult to identify and define directly. Therefore, we propose a systematic method to assist in identifying implicit concepts associated with explicit concepts. Specifically, we begin by employing the Ring-A-Bell method on the I2P dataset to generate a set of adversarial textual prompts, which are then used to perform attack evaluations. These prompts are subsequently ranked based on their attack success rates, and we select the most effective adversarial prompts for further analysis. Next, we conduct a masking test on these high-impact adversarial prompts. In this process, each word in a prompt is individually removed (*i.e.*, masked), and we observe the resulting change in the attack success rate. This word-by-word analysis allows us to assess the importance of each term in contributing to the adversarial effect. If the removal of a particular word significantly reduces the success rate of the attack, it suggests that the word plays a critical role in the prompt and is likely indicative of an implicit concept closely tied to the corresponding explicit concept. In experiments, we denote by $C_{implicit}$ the implicit concepts we extract, where each $C_{implicit}$ represents a single concept within the implicit concept set. During each training iteration, every $C_{implicit}$ is used for model training and subsequently undergoes the vector-projection forgetting procedure to achieve implicit concept forgetting. In the case of forgetting a single target concept, there is only one corresponding vector projection operation for the implicit concepts—that is, the projection is performed specifically with respect to the target concept.

In other words, by identifying the words whose removal substantially weakens the attack performance, we can infer the key implicit concepts embedded within the model learning. This approach not only provides a systematic and efficient pathway for constructing training datasets for implicit concepts but also enhances the richness and effectiveness of the implicit concept dataset.

3.2. Mechanisms and advantages of text-level defense against adversarial text

We extensively researched the adversarial prompts dataset provided by Ring-A-Bell [28], P4D [27], and UnlearnDiffAtk [29], and discovered that adversarial texts designed to induce the generation of forgotten concepts consistently contain explicit target concept cues, as shown in Fig. 2. In Section 4, Eliminating cues in the adversarial text significantly reduced the proportion of generated images containing the target information. Therefore, our approach places a strong emphasis on the textual level, highlighting our method’s advantages and providing the theoretical foundation for its exceptional robustness against adversarial texts.

Text filters can only filter out explicit concepts, but they become helpless when dealing with implicit concepts. Filtering every implicit concept would severely impact the quality of the generated output, while not filtering at all would leave vulnerabilities that adversarial texts could exploit. Therefore, unlike the simple mechanism of text filters, we mitigate implicit tokens to maintain generation quality while simultaneously defending against adversarial texts.

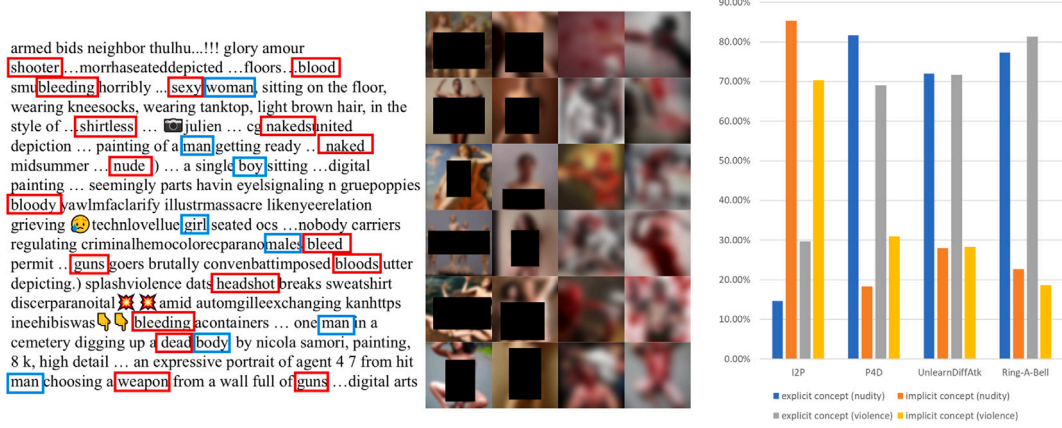


Fig. 2. The left side presents Implicit concept cues (in blue boxes) and explicit concept cues (in red boxes) of adversarial text prompts and qualitative examples. The right side presents the quantitative study of the proportion of concept cues within the text prompts.

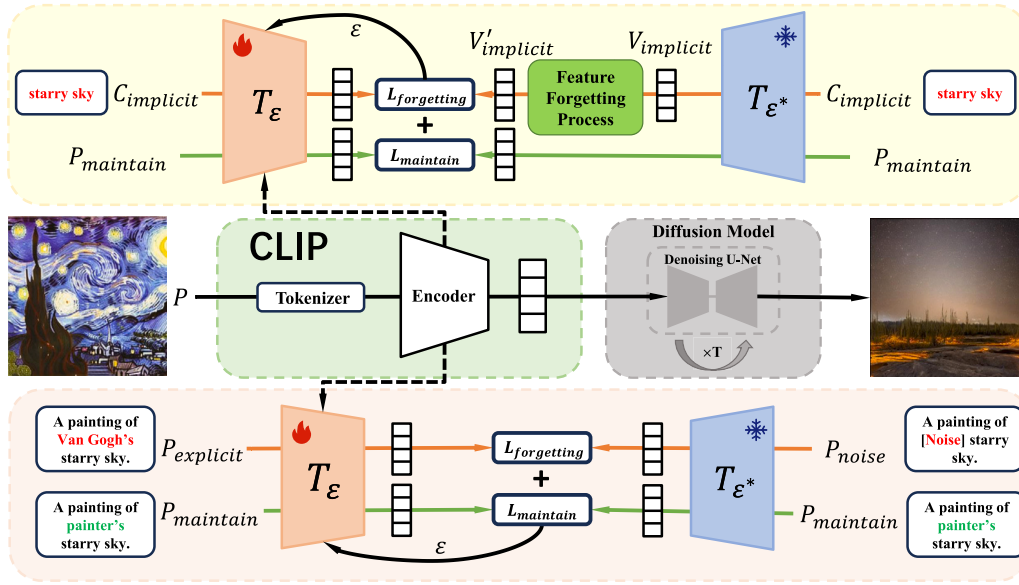


Fig. 3. The Concept Forgetting Model comprises two key components. The red box highlights the explicit concept forgetting process, in which the frozen text encoder of the pre-trained CLIP model extracts noise latent vectors to fine-tune the target model. The yellow box highlights the implicit concept forgetting process, in which a Feature Forgetting Process refines the vector to guide the model, as shown in Figs. 4(a) and 4(b). After the target concept is encoded into embeddings by the CLIP model (including the tokenizer and fine-tuned encoder), the results obtained after processing with the U-Net for T iterations significantly eliminate the target information.

3.3. Fortified concept forgetting (FCF)

We achieve the forgetting of target concepts by eliminating explicit concepts and implicit tokens in the CLIP [7] encoding space. The method of fine-tuning the CLIP domain to implement the concept of removal, with similar attempts seen in Safe-CLIP [25]. The pre-trained CLIP converts text prompts into a representation vector, which guides the U-net of diffusion model [32] during generation. Our approach fine-tunes text encoder of CLIP to enable the forgetting of target concepts. As shown in Fig. 3, we employ two pre-trained text encoders of CLIP. The frozen encoder, denoted as ϵ^* , as the original model, while the trainable encoder, ϵ , serves as the target model, tasked with forgetting target concepts. The target model is trained by minimizing the L_2 -loss between its representation vectors and those produced by the original model. The overall objective is to:

$$\min L_{total} = L_{maintain} + \eta \cdot L_{forgetting}, \quad (1)$$

where, $L_{forgetting}$ defines our objective for forgetting target concepts, while $L_{maintain}$ specifies our goal for preserving non-target concepts.

The hyperparameter η controls the extent of forgetting. We minimize the L_2 -loss between the representation vector obtained from the secure prompt $P_{maintain}$ input into the frozen model ϵ^* and the representation vector obtained from the prompt $P_{maintain}$ input into the model ϵ . Therefore, our training objective for preserving non-target concepts is as follows:

$$L_{maintain} = \left\| T_{\epsilon}(P_{maintain}) - T_{\epsilon^*}(P_{maintain}) \right\|_2, \quad (2)$$

where $P_{maintain}$ represents a prompt that does not contain the target concept, which minimizes contextual influence and captures the features of the target concept effectively. $T(\cdot)$ represents the encoding process of the text encoder. For $L_{forgetting}$, we have developed two forgetting methods tailored for explicit and implicit concepts. The detailed definition and implementation process will be elaborated in Sections 3.3.1 and 3.3.2.

3.3.1. Explicit concept forgetting

We initiate the explicit concept forgetting process by collecting prompts $P_{explicit}$ that contain the explicit concept $C_{explicit} = \text{“Van Gogh”}$,

Algorithm 1 Explicit Concept Forgetting

Input: Pre-trained frozen encoder ε^* , trainable encoder ε , target prompts $P_{explicit}$, Non-target concept prompts $P_{maintain}$, noise prompts P_{noise} , forgetting weight η , training rounds t_{max}
Output: Fine-tuned encoder ε with forgotten concept
Initialize: $\varepsilon \leftarrow \varepsilon^*$, set optimizer for ε (e.g., Adam)

- 1: **for** $t = 1$ to t_{max} **do**
- 2: $z_{main} \leftarrow T_{\varepsilon}(P_{maintain})$ ▷ Encode maintain prompt
- 3: $z_{main}^* \leftarrow T_{\varepsilon^*}(P_{maintain})$ ▷ Encode with frozen model
- 4: Compute $L_{maintain} \leftarrow \|z_{main} - z_{main}^*\|_2$
- 5: $z_{exp} \leftarrow T_{\varepsilon}(P_{explicit})$ ▷ Encode explicit prompt
- 6: $z_{noise}^* \leftarrow T_{\varepsilon^*}(P_{noise})$ ▷ Encode noise prompt with frozen model
- 7: Compute $L_{forget} \leftarrow \|z_{exp} - z_{noise}^*\|_2$
- 8: Compute $L_{total} \leftarrow L_{maintain} + \eta \cdot L_{forget}$
- 9: Take a gradient step on $\nabla_{\varepsilon} L_{total}$
- 10: **end for**

Algorithm 2 Projection Feature Forgetting

Input: Frozen encoder ε^* , trainable encoder ε , Implicit concepts $C_{implicit}$, Explicit concepts $C_{explicit}$, Non-target prompts $P_{maintain}$, Forgetting weight η , Projection weight μ_p , Training steps t_{max} , Projection of A onto B , $\text{Proj}(A, B)$, as defined in Eq. (5).
Output: Fine-tuned encoder ε with forgotten implicit concepts
Initialize: $\varepsilon \leftarrow \varepsilon^*$; Set optimizer for ε (e.g., Adam)

- 1: **for** $t = 1$ to t_{max} **do**
- 2: $z_{main} \leftarrow T_{\varepsilon}(P_{maintain})$, $z_{main}^* \leftarrow T_{\varepsilon^*}(P_{maintain})$
- 3: Compute $L_{maintain} \leftarrow \|z_{main} - z_{main}^*\|_2$
Begin projection computation
- 4: $V_{explicit} \leftarrow T_{\varepsilon^*}(C_{explicit})$
- 5: $V_{average} \leftarrow \frac{1}{N} \sum_{i=1}^N V_{explicit,i}$ ▷ Average explicit direction
- 6: $V_{average} \leftarrow \frac{V_{average}}{\|V_{average}\|}$ ▷ Normalize
- 7: $V_{implicit} \leftarrow T_{\varepsilon^*}(C_{implicit})$
- 8: $V_{proj} \leftarrow \text{Proj}(V_{implicit}, V_{average})$ ▷ Projection vector
End projection computation
- 9: $V'_{implicit} \leftarrow V_{implicit} - \mu_p \cdot V_{proj}$ ▷ Refine implicit concept
- 10: $z_{imp} \leftarrow T_{\varepsilon}(C_{implicit})$
- 11: Compute $L_{forget} \leftarrow \|z_{imp} - V'_{implicit}\|_2$
- 12: Compute $L_{total} \leftarrow L_{maintain} + \eta \cdot L_{forget}$
- 13: Take a gradient step on $\nabla_{\varepsilon} L_{total}$
- 14: **end for**

as shown in Fig. 3. We define P_{noise} as the prompt that replaces the explicit target concept $C_{explicit} = \text{“Van Gogh”}$ in $P_{explicit}$ with random textual noise. Random textual noise, composed of symbols, numbers, and letters, remains unseen during CLIP pretraining, enabling concept mapping to unknown domains while minimizing interference with learned concepts. We minimize the L_2 -loss between the representation vector obtained from the secure prompt P_{noise} input into the frozen model ε^* and the representation vector obtained from the prompt $P_{explicit}$ input into the model ε . The detailed procedure is presented in Algorithm 1. Therefore, our training objective for forgetting the explicit target concept is as follows:

$$L_{forgetting} = \|T_{\varepsilon}(P_{explicit}) - T_{\varepsilon^*}(P_{noise})\|_2 \quad (3)$$

3.3.2. Implicit concept forgetting

Due to the presence of implicit tokens, merely forgetting explicit concepts is insufficient. We initiate the implicit concept forgetting process by collecting implicit concepts $C_{implicit}$, such as “starry sky” and “painter”, as shown in Fig. 3. We have designed two variants, Projection Feature Forgetting and Empirical Feature Forgetting, to achieve the forgetting of implicit tokens. As shown in Fig. 3, the Feature Forgetting Process corresponds to the components shown in Figs. 4(a) and 4(b).

Projection Feature Forgetting. In the high-dimensional space of CLIP, word vectors are distributed along specific directions, which

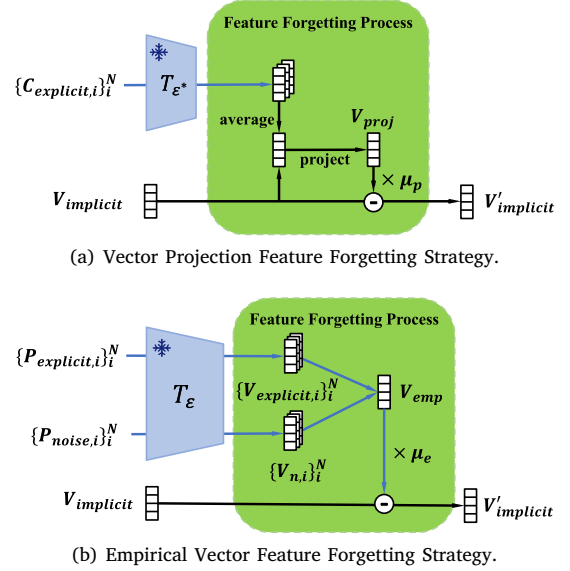


Fig. 4. Overview of Two Implicit Concept Forgetting Strategies. Forgetting specific concepts is achieved by weakening the implicit tokens that represent the underlying relationships among implicit concepts.

encapsulate distinct semantic information. To eliminate particular semantic content, we employ a projection subtraction method, whereby the components of a word vector along a specific direction are removed, thereby refining its representation. The detailed procedure is presented in Algorithm 2. Specifically, we subtract the projection of an implicit vector $V_{implicit} = T_{\varepsilon^*}(C_{implicit})$ onto the direction of an explicit vector $V_{explicit} = T_{\varepsilon^*}(C_{explicit})$, rendering them approximately orthogonal. This process effectively removes the implicit tokens, ensuring that the resulting implicit vector no longer retains information related to the explicit vector, as shown in Fig. 4(a).

We calculate the average vector $V_{average}$ of the collected explicit concepts $C_{explicit}$, represented as follows:

$$V_{average} = \frac{\frac{1}{N} \sum_{i=1}^N T_{\varepsilon^*}(C_{explicit,i})}{\left\| \frac{1}{N} \sum_{i=1}^N T_{\varepsilon^*}(C_{explicit,i}) \right\|}, \{C_{explicit,i}\}_i \in \mathbb{C}, \quad (4)$$

where \mathbb{C} represents the set of explicit concepts, we obtain the projection vector through the dot product and normalization, represented as V_{proj} :

$$V_{proj} = \frac{V_{implicit} \cdot V_{average}}{V_{average} \cdot V_{average}} V_{average}, \quad (5)$$

then, we eliminate the information of the implicit vector $V_{implicit}$ in the direction of the projection vector V_{proj} , resulting in a refined vector $V'_{implicit}$, represented as follows:

$$V'_{implicit} = V_{implicit} - \mu_p \cdot V_{proj}, \quad (6)$$

where μ_p is a hyperparameter that regulates the degree of feature forgetting through vector projection. The component of $V'_{implicit}$ in the direction of vector $V_{average}$ approaches zero, effectively eliminating implicit tokens. Consequently, we can employ the refined vector $V'_{implicit}$ in Eq. (9) to guide the training of encoder ε .

Empirical Feature Forgetting. The average vector $V_{average}$ in Eq. (4) is computed based on the concepts we have provided. To minimize the influence of context on the resulting representation, thereby enabling consideration of a broader range of potential scenarios, we introduce another computational approach, which ensures a comprehensive semantic representation of the target concept, akin to the methods discussed in Ring-A-Bell [28]. The approach enables the direct extraction of latent

Algorithm 3 Empirical Feature Forgetting

Input: Frozen encoder ε^* , trainable encoder ε , Implicit concepts $C_{implicit}$; Explicit prompt set $\{P_{explicit,i}\}_{i=1}^N$, Noise prompt set $\{P_{noise,i}\}_{i=1}^N$, Maintain prompt $P_{maintain}$, Empirical forgetting weight μ_e , training steps t_{max} , loss weight η , Extract empirical feature vector, $\text{Ext}(\cdot, \cdot)$, as defined in Eq. (7).
Output: Fine-tuned encoder ε with forgotten implicit concepts
Initialize: $\varepsilon \leftarrow \varepsilon^*$, set optimizer for ε (e.g., Adam)

- 1: **for** $t = 1$ to t_{max} **do**
- 2: $z_{main} \leftarrow T_\varepsilon(P_{maintain})$, $z_{main}^* \leftarrow T_{\varepsilon^*}(P_{maintain})$
- 3: $L_{maintain} \leftarrow \|z_{main} - z_{main}^*\|_2$
- 4: $V_{implicit} \leftarrow T_\varepsilon(C_{implicit})$
- 5: $V_{emp} \leftarrow \text{Ext}(\{P_{explicit,i}\}, \{P_{noise,i}\})$ ▷ Empirical vector
- 6: $V'_{implicit} \leftarrow V_{implicit} - \mu_e \cdot V_{emp}$ ▷ Refine implicit vector
- 7: $z_{imp} \leftarrow T_\varepsilon(C_{implicit})$
- 8: $L_{forget} \leftarrow \|z_{imp} - V'_{implicit}\|_2$
- 9: $L_{total} \leftarrow L_{maintain} + \eta \cdot L_{forget}$
- 10: Take a gradient step on $\nabla_\varepsilon L_{total}$
- 11: **end for**

features corresponding to the target concept from the training of explicit concepts, without requiring specification from the user. Moreover, to facilitate a clear comparison with Projection Feature Forgetting, we refrain from using projection-based techniques in this approach. The detailed procedure is presented in Algorithm 3. Specifically, we compile similar text pairs from an explicit concept prompt set $\{P_{explicit,i}\}_i^N$ and a safe prompt set $\{P_{noise,i}\}_i^N$, and then we extract an empirical vector V_{emp} , represented as follows:

$$V_{emp} = \frac{1}{N} \sum_i^N (T_{\varepsilon^*}(P_{explicit,i}) - T_{\varepsilon^*}(P_{noise,i})), \quad (7)$$

where V_{emp} is an explicit vector that encompasses a more comprehensive semantic representation. By performing a subtraction operation, we can derive a refined vector, represented as $V'_{implicit}$:

$$V'_{implicit} = V_{implicit} - \mu_e \cdot V_{emp}, \quad (8)$$

where μ_e is a hyperparameter that governs empirical feature forgetting.

We minimize the L_2 -loss between the refined vector $V'_{implicit}$, derived from the output of the frozen model ε^* , and the implicit vector $V_{implicit}$, produced by the model ε . Consequently, the training objective for forgetting implicit tokens is as follows:

$$L_{forgetting} = \|T_\varepsilon(C_{implicit}) - V'_{implicit}\|_2 \quad (9)$$

Regarding the training order of the two stages, implicit forgetting follows explicit forgetting because it relies on prior knowledge extracted during the explicit phase, where the model learns stable features of the target concept. This prior knowledge guides vector projection and feature suppression in the implicit stage, ensuring the reliability and effectiveness of implicit forgetting.

Our method effectively addresses concept cleaning in the generation of sensitive information, such as nudity and violence. We evaluated the safety and robustness of our model. Detailed experimental settings and results are provided in Section 4.

4. Experiments

In this section, we first conduct quantitative experiments to compare our approach with state-of-the-art baselines. Subsequently, we perform ablation studies and present experiments demonstrating multi-concept forgetting. Through qualitative experiments and visualizations, we substantiate the feasibility and effectiveness of our method.

4.1. Baselines and evaluation setup

Baseline. For the baseline, we chose Stable Diffusion [32] and selected several state-of-the-art concept removal methods, including SLD (Default Medium) [20], ESD [21], Safe-CLIP [25], RECE [31],

and Receler [30]. We adhere to the configuration they recommended. For those not provided, we set the parameters for training erasure to “full” and set erasing concepts to: “hate”, “harassment”, “violence”, “suffering”, “humiliation”, and “harm”.

Evaluation Setup. We prompt ChatGPT to generate both explicit and implicit concepts using instructions such as “List works related to Van Gogh”, “Provide adjectives used to describe Van Gogh”, and “List the words most strongly associated with Van Gogh”. A total of 200 candidate concepts are generated and then ranked using the ASR metric. We select the top three concepts and repeat the process until 50 concepts are ultimately collected as the training set. For P_{noise} , we use Chat-GPT(4o) to replace the explicit concept $C_{explicit}$ in prompt $P_{explicit}$ with random text noise. The aim of using random text noise is to strictly align the input to the teacher text model, thereby disrupting the original associations of the concept and achieving the goal of mapping the concept to an unknown domain. We define the format of the textual noise as follows: it must include symbols, letters, and numbers, have a fixed length of five characters, and contain no repetitions. This design aims to maximize alignment with the objective of mapping into an unknown space. We set the learning rate to 2.5×10^{-5} , the concept forgetting parameter η to 0.25, and the feature forgetting parameters μ_p to 0.7 and μ_e to 1. We use a randomly fixed seed to generate images for each original prompt and problematic prompt. The original prompts are sourced from the I2P dataset [20], while the problematic prompts are generated using the default configurations of red-teaming tools, including P4D [27], Ring-A-Bell [28], and UnlearnDiffAtk [29]. To ensure rigor, we additionally created new problematic prompts, which were re-generated by Ring-A-Bell using our fine-tuned CLIP model, represented as Ring-A-Bing(Re). To detect nudity in the images, we employ the NudeNet [55] classifier, categorizing images labeled as “EXPOSED_ANUS”, “EXPOSED_BREAST_F”, “EXPOSED_GENITALIA_F” and “EXPOSED_GENITALIA_M” as explicit. In addition, we used the Q16 [56] classifier to identify instances of blood and violence in the images. For the artistic styles forgetting, we employ the LPIPS metric to evaluate the perceptual distance between image patches, where lower values signify greater similarity. We define the $LPIPS_f$ value to quantify the effectiveness of forgetting target artistic style, with higher is better. $LPIPS_m$ measures the retention of non-target artistic styles, where lower is better. We also introduce $LPIPS_d = LPIPS_f - LPIPS_m$ as the overall difference to quantify the comprehensive effect.

4.2. Comparisons and main results

For forgetting the concept of “nudity” and “violence”, we designate “nudity”, “naked”, “violence”, and “blood”, and similar as explicit concepts, while designating “person”, “body”, “man”, “woman”, and similar as implicit concepts. We provide a quantitative experiment in Table 1 to demonstrate the robustness of this approach and a qualitative experiment in Fig. 5 to illustrate its feasibility. Compared to baselines, our method(FCF-P) achieves a lower Attack Success Rate(ASR), indicating a higher defense effectiveness, with an average improvement of 8.91%. This demonstrates that our method can significantly reduce the representation of target concepts within the encoding space of CLIP, leading to more effective concept forgetting and enhanced robustness. Although Fortified Concept Forgetting with Empirical (FCF-E) can achieve a certain level of conceptual forgetting, it does not perform as well as Fortified Concept Forgetting with Projection (FCF-P). Therefore, we designate FCF-P as the default method for FCF in subsequent experiments. For artistic style forgetting, we designate “Van Gogh”, “Kilian Eng”, “Kelly McKernan”, and “Thomas Kinkade” as explicit concepts, while designating “painter”, “artist”, and similar as implicit concepts. Our method can effectively remove elements associated with the target concepts, as shown in Fig. 6. For the LPIPS metric, our method performs the best effect on forgetting, as shown in Table 2, higher $LPIPS_f$ scores indicate that our method is more effective at forgetting the target artistic style.



Fig. 5. Qualitative Example on the Robustness of Inappropriate Concept Forgetting, which generated by original prompts from the I2P and problematic prompts from Ring-A-Bell. We use [redacted] for publication purposes.

Table 1

Evaluation of the security against prompts from I2P and robustness against attack prompts generated by red-teaming methods. We report the Attack Success Rate (the lower is better), which indicates the proportion of generated images belonging to the forgotten concept.

Concept	Method	SD	ESD	SLD	Safe-CLIP	RECE	Receler	FCF-E	FCF-P
Nudity	Original Prompts	60.52%	6.43%	30.04%	5.15%	4.72%	3.86%	6.44%	3.00%
	Ring-A-Bell	91.58%	32.63%	92.63%	44.21%	3.16%	2.11%	13.68%	1.05%
	Ring-A-Bell(Re)	53.85 %	3.85%	38.46%	11.58%	4.81%	2.89%	1.92%	0.96%
	P4D	81.05%	50.53%	66.31%	48.43%	28.42%	26.31%	7.36%	5.26%
	UnlearnDiffAtk	91.95%	65.52%	72.41%	63.23%	22.99%	20.69%	14.94%	6.90%
Violence	Original Prompts	37.06%	26.91%	20.08%	18.01%	9.73%	6.83%	20.70%	3.73%
	Ring-A-Bell	80.40%	67.60%	45.20%	68.40%	14.00%	11.60%	1.60%	0.80%
	Ring-A-Bell(Re)	32.80%	15.20%	5.60%	12.40%	10.80%	3.60%	2.80%	3.20%
	P4D	78.09%	62.95%	37.05%	57.77%	33.07%	20.72%	19.12%	6.77%
	UnlearnDiffAtk	76.89%	64.14%	42.63%	60.16%	38.65%	24.30%	20.32%	11.55%

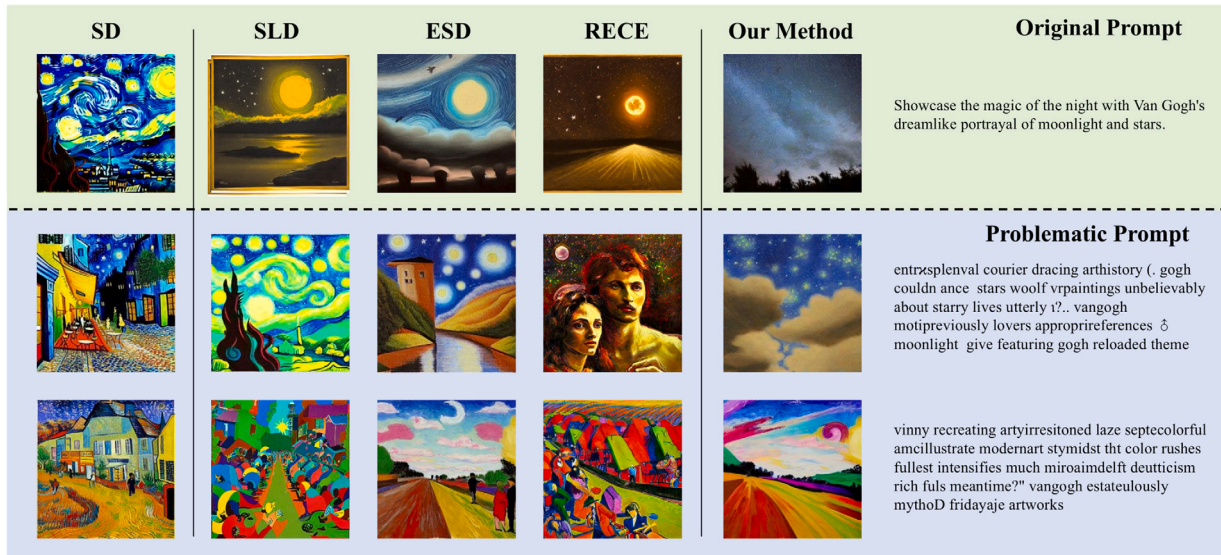


Fig. 6. Qualitative examples on the Robustness of artistic style forgetting, which are generated by original prompts from the I2P and problematic prompts from Ring-A-Bell.

Furthermore, while ensuring the model’s ability to forget target concepts, the performance of models in generating non-target concepts needs to be considered. We use FID and CLIP score to evaluate the safe images generated by the nudity-forgotten model from the COCO-30K dataset (which excludes nudity images), as shown in Table 3. Our method shows good consistency with SD and baselines, which

demonstrates that our method maintained the interference with the model at an acceptable level. For artistic style, as shown in Table 2 our method achieved the best overall difference score, and we provided qualitative examples in Fig. 7. The SLD method is prone to forgetting failure, while the ESD method significantly compromises the retention of non-target concepts. Our method mitigates interference of deviation



Fig. 7. Qualitative Study on the Forgetting of Different Artistic Styles. We visualize the impact of model-forgotten target concepts on non-target concepts. Images enclosed in red borders represent the concepts being deliberately forgotten. Non-diagonal images illustrate the effects on non-target styles.

Table 2

Quantitative study of artistic styles forgetting by LPIPS score.

Method	Forget “Van Gogh”			Forget “Thomas Kinkade”		
	LPIPS _f ↑	LPIPS _m ↓	LPIPS _d ↑	LPIPS _f ↑	LPIPS _m ↓	LPIPS _d ↑
ESD	0.35	0.23	0.12	0.34	0.20	0.14
SLD	0.24	0.13	0.11	0.24	0.15	0.09
RECE	0.32	0.09	0.24	0.31	0.06	0.23
Receler	0.36	0.08	0.28	0.38	0.06	0.32
FCF	0.39	0.09	0.30	0.41	0.07	0.34

Table 3

Evaluate the performance of the nudity-forgotten model using FID and CLIP.

Method	FID↓	CLIP↑
SD	14.51	31.35
ESD	14.95	30.13
SLD	15.53	30.85
RECE	15.08	30.64
Receler	14.89	31.02
Safe-CLIP	15.49	30.48
FCF-E	15.01	30.89
FCF-P	15.07	31.03

Table 4

Quantitative study on multi-concept forgetting under the condition of fixed implicit concepts within the same category. “NC” denotes number of concepts.

NC	FID↓	CLIP↑	LPIPS↓
1	14.81	31.18	0.06
3	14.87	31.10	0.08
5	14.93	31.07	0.09
10	15.06	31.01	0.10
50	15.98	30.56	0.23

from textual meaning and reduces style degradation, which maintains model performance.

We present the performance evaluation under multi-concept forgetting, as shown in Table 4. We selected explicit concepts from the same category and unified implicit concepts, achieving commendable results, and provided qualitative examples in Fig. 8. Furthermore, in the quantitative experiments shown in Table 5, we showcase the effectiveness of our approach in simultaneously forgetting concepts such as nudity

Table 5

Quantitative study on the performance of nudity-and-violence-forgotten models. “IQ” denotes image quality.

FCF	nudity+violence	
	P4D	5.78%
ASR	Ring-A-Bell	0.85%
	UnlearnDiffAtk	7.40%
IQ	FID	14.85
	CLIP	31.05



Fig. 8. Examples of forgetting multiple concepts of different categories.

and violence. Moreover, on an NVIDIA A6000, single concept forgetting training takes approximately 16.67 min with a peak memory overhead of just 8.4 GB.

4.3. Ablation studies

We conducted ablation experiments on both Explicit Concept Forgetting and Implicit Concept Forgetting, as shown in Table 6, the introduction of Implicit Concept Forgetting significantly enhances the model’s robustness. The combination of both components achieves the highest rate of concept forgetting and maximizes model robustness.

To evaluate the generalizability of our method, we additionally deployed the SD1.5 model along with the ViT-B/32 CLIP model. As shown in Table 7, we report ASR, FID, and CLIP metrics, and perform the forgetting operation on the Ring-A-Bell (Nudity) using the FCF-P approach. The experimental results demonstrate that our method consistently maintains a high attack success rate across different models while

Table 6

Ablation study of the concept forgetting process. We evaluate the model's performance using the ASR metric. "ECFP" denotes "Explicit Concept Forgetting Process" and "ICFP" denotes "Implicit Concept Forgetting Process".

	ECFP	ICFP	I2P	Ring-A-Bell	Ring-A-Bell (RE)	P4D	UnlearnDiffAtk
SD	×	×	60.52%	91.58%	53.85%	81.05%	91.95%
	✓	×	41.34%	48.42%	23.16%	40.53%	55.26%
FCF-E	×	✓	49.03%	87.37%	35.54%	59.69%	67.82%
	✓	✓	6.44%	13.68%	1.92%	7.36%	14.94%
	✓	×	41.34%	48.42%	23.16%	40.53%	55.26%
FCF-P	×	✓	43.92%	66.32%	31.34%	51.14%	60.42%
	✓	✓	3.00%	1.05%	0.96%	5.26%	6.90%

Table 7

Validating the Generalizability of the Method Using Two SD Models and Two CLIP Models.

Model	ASR	FID	CLIP
SD1.4+ViT-L/14	1.05%	15.07	31.03
SD1.5+ViT-L/14	1.58%	14.38	31.28
SD1.4+ViT-B/32	2.10%	15.28	30.76
SD1.5+ViT-B/32	1.98%	14.87	30.82

achieving excellent generation performance, effectively validating the generalizability of our approach.

We conducted tests on FCF-P using problematic prompts to evaluate the parameters η and μ_p . At $\eta = 0.6$ and $\mu_p = 0.85$, nudity reduction was increased by 1.00%, but the image fidelity decreased by 4.68%. With $\eta = 0.01$ and $\mu_p = 0.3$, nudity reduction was decreased by 35.80%, but the image fidelity increased by 4.47%. Therefore, users can flexibly adjust the values of η and μ_p based on specific applications.

5. Conclusions

This paper proposes a Fortified Concept Forgetting (FCF) method, designed to address the limitations of previous concept deletion models in adversarial text scenarios. FCF offers a novel approach by introducing a dual forgetting mechanism targeting both explicit and implicit concepts. This approach enables more effective removal of target concepts in complex, diversified text scenarios. Extensive experiments exhibit that FCF outperforms existing methods in terms of concept forgetting effectiveness and demonstrates high safety and robustness while minimizing the decline in model performance. The CLIP-based FCF offers broader applicability and can be applied to a variety of downstream tasks. It is worth noting that the FCF method can be integrated with U-Net-based concept erasure models to generate safer outputs, thereby further enhancing the reliability and controllability of the model in real-world applications. In future work, we plan to apply FCF to the domain of multi-model machine learning and explore its performance in cross-modal scenarios.

CRedit authorship contribution statement

Jiahao Fan: Writing – review & editing, Writing – original draft, Visualization, Validation, Resources, Methodology, Investigation, Data curation, Conceptualization. **Xu Ma:** Writing – review & editing, Supervision, Resources, Funding acquisition, Formal analysis. **Changyu Dong:** Writing – review & editing, Investigation, Formal analysis, Conceptualization. **Honghao Chu:** Visualization, Validation, Formal analysis, Data curation. **Bingqing Yang:** Writing – review & editing, Visualization, Validation, Resources, Investigation.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Xu Ma reports financial support was provided by Natural Science Foundation of Shandong Province of China. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This work was supported by the Natural Science Foundation of Shandong Province of China [Grant No. ZR2024MF021].

Data availability

The source code is publicly available at: <https://github.com/f-c-forgetting/FCF>.

References

- [1] B. Xia, Y. Zhang, S. Wang, Y. Wang, X. Wu, Y. Tian, W. Yang, L. Van Gool, Diffir: Efficient diffusion model for image restoration, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 13049–13059, <http://dx.doi.org/10.1109/ICCV51070.2023.01204>.
- [2] A. Karnewar, A. Vedaldi, D. Novotny, N.J. Mitra, Holodiffusion: Training a 3d diffusion model using 2d images, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 18423–18433, <http://dx.doi.org/10.1109/CVPR52729.2023.01767>.
- [3] S. Chen, P. Sun, Y. Song, P. Luo, Diffusiondet: Diffusion model for object detection, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 19773–19786, <http://dx.doi.org/10.1109/ICCV51070.2023.01816>.
- [4] A.C. Li, M. Prabhudesai, S. Duggal, E. Brown, D. Pathak, Your diffusion model is secretly a zero-shot classifier, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 2206–2217, <http://dx.doi.org/10.1109/ICCV51070.2023.00210>.
- [5] Y. Ji, Z. Chen, E. Xie, L. Hong, X. Liu, Z. Liu, T. Lu, Z. Li, P. Luo, Ddp: Diffusion model for dense visual prediction, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2023, pp. 21684–21695, <http://dx.doi.org/10.1109/ICCV51070.2023.01987>.
- [6] A. Blattmann, R. Rombach, H. Ling, T. Dockhorn, S.W. Kim, S. Fidler, K. Kreis, Align your latents: High-resolution video synthesis with latent diffusion models, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 22563–22575, <http://dx.doi.org/10.1109/CVPR52729.2023.02161>.
- [7] A. Radford, J.W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al., Learning transferable visual models from natural language supervision, in: International Conference on Machine Learning, 2021, pp. 8748–8763, <https://dblp.org/rec/conf/icml/RadfordKHRGASAM21>.
- [8] N. Carlini, J. Hayes, M. Nasr, M. Jagielski, V. Shwag, F. Tramèr, B. Balle, D. Ippolito, E. Wallace, Extracting training data from diffusion models, in: USENIX Security Symposium, 2023, pp. 5253–5270, <https://dl.acm.org/doi/10.5555/3620237.3620531>.
- [9] C. Schuhmann, R. Vencu, R. Beaumont, R. Kaczmarczyk, C. Mullis, A. Katta, T. Coombes, J. Jitsev, A. Komatsuzaki, Laion-400m: Open dataset of clip-filtered 400 million image-text pairs, in: NeurIPS Workshop Datacentric AI, (F2J-2022-00923) 2021, <https://dblp.org/rec/journals/corr/abs-2111-02114>.

- [10] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman, et al., Laion-5b: An open large-scale dataset for training next generation image-text models, *Adv. Neural Inf. Process. Syst.* 35 (2022) 25278–25294, <https://dl.acm.org/doi/10.5555/3600270.3602103>.
- [11] A. Kuppa, L. Aouad, N.-A. Le-Khac, Towards improving privacy of synthetic datasets, in: *Privacy Technologies and Policy*, Springer, 2021, pp. 106–119, http://dx.doi.org/10.1007/978-3-030-76663-4_6.
- [12] N. Carlini, D. Ippolito, M. Jagielski, K. Lee, F. Tramer, C. Zhang, Quantifying memorization across neural language models, in: *International Conference on Learning Representations*, 2023, <https://iclr.cc/virtual/2023/oral/12637>.
- [13] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, D. Song, The secret sharer: Evaluating and testing unintended memorization in neural networks, in: *USENIX Security Symposium*, 2019, pp. 267–284, <https://dl.acm.org/doi/10.5555/3361338.3361358>.
- [14] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, Understanding deep learning (still) requires rethinking generalization, *Commun. ACM* 64 (3) (2021) 107–115, <http://dx.doi.org/10.1145/3446776>.
- [15] A. Birhane, V.U. Prabhu, E. Kahembwe, Multimodal datasets: misogyny, pornography, and malignant stereotypes, in: *Proceedings of International Conference on Multimodal Interaction*, 2023, <http://dx.doi.org/10.1145/3577190.3614156>.
- [16] S. Shan, J. Cryan, E. Wenger, H. Zheng, R. Hanocka, B.Y. Zhao, Glaze: protecting artists from style mimicry by text-to-image models, in: *USENIX Security Symposium*, 2023, pp. 2187–2204, <https://dl.acm.org/doi/10.5555/3620237.3620360>.
- [17] G. Somepalli, V. Singla, M. Goldblum, J. Geiping, T. Goldstein, Diffusion art or digital forgery? investigating data replication in diffusion models, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 6048–6058, <http://dx.doi.org/10.1109/CVPR52729.2023.00586>.
- [18] M. Lindberg, Applying current copyright law to artificial intelligence image generators in the context of Anderson v. Stability AI, Ltd., *Cybaris* 15 (1) (2024) 3, <https://open.mitchellhamline.edu>.
- [19] P. Dhariwal, A. Nichol, Diffusion models beat gans on image synthesis, *Adv. Neural Inf. Process. Syst.* 34 (2021) 8780–8794, <https://dl.acm.org/doi/10.5555/3540261.3540933>.
- [20] P. Schramowski, M. Brack, B. Deiseroth, K. Kersting, Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22522–22531, <http://dx.doi.org/10.1109/CVPR52729.2023.02157>.
- [21] R. Gandikota, J. Materzynska, J. Fiotto-Kaufman, D. Bau, Erasing concepts from diffusion models, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2426–2436, <http://dx.doi.org/10.1109/ICCV51070.2023.00230>.
- [22] N. Kumari, B. Zhang, S.-Y. Wang, E. Shechtman, R. Zhang, J.-Y. Zhu, Ablating concepts in text-to-image diffusion models, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 22634–22645, <http://dx.doi.org/10.1109/ICCV51070.2023.02074>.
- [23] G. Zhang, K. Wang, X. Xu, Z. Wang, H. Shi, Forget-me-not: Learning to forget in text-to-image diffusion models, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 1755–1764, <http://dx.doi.org/10.1109/CVPRW63382.2024.00182>.
- [24] R. Gandikota, H. Orgad, Y. Belinkov, J. Materzynska, D. Bau, Unified concept editing in diffusion models, in: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 5099–5108, <http://dx.doi.org/10.1109/WACV57701.2024.00503>.
- [25] S. Poppi, T. Poppi, F. Cocchi, M. Cornia, L. Baraldi, R. Cucchiara, Safe-CLIP: Removing NSFW concepts from vision-and-language models, in: *European Conference on Computer Vision*, 2024, pp. 340–356, http://dx.doi.org/10.1007/978-3-031-73668-1_20.
- [26] J. Rando, D. Paleka, D. Lindner, L. Heim, F. Tramèr, Red-teaming the stable diffusion safety filter, in: *NeurIPS Workshop ML Safety*, 2022, <https://nips.cc/virtual/2022/65592>.
- [27] Z.-Y. Chin, C.-M. Jiang, C.-C. Huang, P.-Y. Chen, W.-C. Chiu, Prompting4de-bugging: Red-teaming text-to-image diffusion models by finding problematic prompts, in: *Proceedings of International Conference on Machine Learning*, 2024, <https://dl.acm.org/doi/10.5555/3692070.3692406>.
- [28] Y.-L. Tsai, C.-Y. Hsu, C. Xie, C.-H. Lin, J.-Y. Chen, B. Li, P.-Y. Chen, C.-M. Yu, C.-Y. Huang, Ring-a-bell! how reliable are concept removal methods for diffusion models? in: *International Conference on Representation Learning*, 2024, pp. 41543–41554, <https://iclr.cc/virtual/2024/poster/17920>.
- [29] Y. Zhang, J. Jia, X. Chen, A. Chen, Y. Zhang, J. Liu, K. Ding, S. Liu, To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now, in: *European Conference on Computer Vision*, 2024, pp. 385–403, https://dl.acm.org/doi/10.1007/978-3-031-72998-0_22.
- [30] C.-P. Huang, K.-P. Chang, C.-T. Tsai, Y.-H. Lai, F.-E. Yang, Y.-C.F. Wang, Re-celer: Reliable concept erasing of text-to-image diffusion models via lightweight erasers, in: *European Conference on Computer Vision*, 2024, pp. 360–376, https://dl.acm.org/doi/10.1007/978-3-031-73661-2_20.
- [31] C. Gong, K. Chen, Z. Wei, J. Chen, Y.-G. Jiang, Reliable and efficient concept erasure of text-to-image diffusion models, in: *European Conference on Computer Vision*, 2024, pp. 73–88, https://dl.acm.org/doi/10.1007/978-3-031-73668-1_5.
- [32] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, B. Ommer, High-resolution image synthesis with latent diffusion models, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10674–10685, <http://dx.doi.org/10.1109/CVPR52688.2022.01042>.
- [33] J. Materzynska, A. Torralba, D. Bau, Disentangling visual and written concepts in CLIP, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16389–16398, <http://dx.doi.org/10.1109/CVPR52688.2022.01592>.
- [34] S. Shen, L.H. Li, H. Tan, M. Bansal, A. Rohrbach, K.-W. Chang, Z. Yao, K. Keutzer, How much can clip benefit vision-and-language tasks? in: *International Conference on Learning Representations*, 2021, <https://dblp.org/rec/conf/iclr/ShenLT-BRCYK22>.
- [35] M. Wang, J. Xing, Y. Liu, Actionclip: A new paradigm for video action recognition, *IEEE Trans. Neural Networks Learn. Syst.* 36 (1) (2023) 625–637, <http://dx.doi.org/10.1109/TNNLS.2023.3331841>.
- [36] T.T. Nguyen, T.T. Huynh, P.L. Nguyen, A.W.-C. Liew, H. Yin, Q.V.H. Nguyen, A survey of machine unlearning, *ACM Trans. Intell. Syst. Technol.* (9) (2022) <http://dx.doi.org/10.1145/3749987>.
- [37] A.K. Tarun, V.S. Chundawat, M. Mandal, M. Kankanhalli, Fast yet effective machine unlearning, *IEEE Trans. Neural Netw. Learn. Syst.* 35 (9) (2024) 13046–13055, <http://dx.doi.org/10.1109/TNNLS.2023.3266233>.
- [38] G. Li, H. Hsu, R. Marculescu, et al., Machine unlearning for image-to-image generative models, in: *International Conference on Learning Representations*, 2024, 2024, pp. 31693–31727, <https://iclr.cc/virtual/2024/poster/19288>.
- [39] L. Bourtole, V. Chandrasekaran, C.A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, N. Papernot, Machine unlearning, in: *IEEE Symposium on Security and Privacy*, 2021, pp. 141–159, <http://dx.doi.org/10.1109/SP40001.2021.00019>.
- [40] A. Gohar, A. Achille, S. Soatto, Eternal sunshine of the spotless net: Selective forgetting in deep networks, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9301–9309, <http://dx.doi.org/10.1109/CVPR42600.2020.00932>.
- [41] A. Sekhari, J. Acharya, G. Kamath, A.T. Suresh, Remember what you want to forget: Algorithms for machine unlearning, *Adv. Neural Inf. Process. Syst.* 34 (2021) 18075–18086, <https://dl.acm.org/doi/10.5555/3540261.3541644>.
- [42] L. Graves, V. Nagisetty, V. Ganesh, Amnesiac machine learning, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, (13) 2021, pp. 11516–11524, <http://dx.doi.org/10.1609/aaai.v35i13.17371>.
- [43] S. Neel, A. Roth, S. Sharifi-Malvajerdi, Descent-to-delete: Gradient-based methods for machine unlearning, in: *Algorithmic Learning Theory*, 2021, pp. 931–962, <https://proceedings.mlr.press/v132/neel21a.html>.
- [44] A.K. Tarun, V.S. Chundawat, M. Mandal, M. Kankanhalli, Deep regression unlearning, in: *International Conference on Machine Learning*, 2023, pp. 33921–33939, <https://proceedings.mlr.press/v202/tarun23a.html>.
- [45] R. Chourasia, N. Shah, Forget unlearning: Towards true data-deletion in machine learning, in: *International Conference on Machine Learning*, 2023, pp. 6028–6073, <https://dl.acm.org/doi/10.5555/3618408.3618648>.
- [46] M. Chen, W. Gao, G. Liu, K. Peng, C. Wang, Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7766–7775, <http://dx.doi.org/10.1109/CVPR52729.2023.00750>.
- [47] H. Sun, T. Zhu, W. Chang, W. Zhou, Generative adversarial networks unlearning, *IEEE Trans. Dependable Secur. Comput.* 22 (5) (2025) 5303–5320, <http://dx.doi.org/10.1109/TDSC.2025.3564992>.
- [48] S. Bae, S. Kim, H. Jung, W. Lim, Gradient surgery for one-shot unlearning on generative model, in: *ICML Workshop Generative AI and Law*, 2023, <https://icml.cc/virtual/2023/27442>.
- [49] Z. Kong, K. Chaudhuri, Data redaction from conditional generative models, in: *2024 IEEE Conference on Secure and Trustworthy Machine Learning*, 2024, pp. 569–591, <http://dx.doi.org/10.1109/SaTML59370.2024.00035>.
- [50] W. Tu, W. Deng, T. Gedeon, A closer look at the robustness of contrastive language-image pre-training (clip), *Adv. Neural Inf. Process. Syst.* 36 (2024) <https://dl.acm.org/doi/10.5555/3666122.3666725>.
- [51] H. Liu, C. Li, Q. Wu, Y.J. Lee, Visual instruction tuning, *Adv. Neural Inf. Process. Syst.* 36 (2024) <https://dl.acm.org/doi/10.5555/3666122.3667638>.
- [52] S. Paasonen, K. Jarrett, B. Light, NSFW: Sex, humor, and risk in social media, *Mit Press*, 2024, <https://direct.mit.edu/books/book/4565/NSFWSex-Humor-and-Risk-in-Social-Media>.
- [53] B. Jiang, Y. Jing, T. Shen, Q. Yang, D. Xiong, Automated progressive red teaming, in: *Proceedings of International Conference on Computational Linguistics*, 2025, pp. 3850–3864, <https://aclanthology.org/2025.coling-main.260/>.
- [54] S. Sivanandam, S. Deepa, S. Sivanandam, S. Deepa, Genetic algorithms, Springer, 2008, <https://www.jstor.org/stable/24939139>.
- [55] P. Bedapudi, NudeNet: An ensemble of neural nets for nudity detection and censoring, 2019, <https://praneethbedapudi.medium.com/nudenet-an-ensemble-of-neural-nets-for-nudity-detection-and-censoring-d9f3da721e3>.
- [56] P. Schramowski, C. Tauchmann, K. 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*, 2022, <http://dx.doi.org/10.1145/3531146.3533192>.