# REDEditing: Relationship-Driven Precise Backdoor Poisoning on Text-to-Image Diffusion Models

Chongye Guo
Shanghai University
Shanghai, China

Jinhu Fu
Beijing University of Posts and Telecommunications
Beijing, China

Junfeng Fang
National University of Singapore
Singapore

Kun Wang*
Nanyang Technological University
Singapore

Guorui Feng*
Shanghai University
Shanghai, China

Figure 1: Our backdoor attack method, `REDEditing`, manipulates the visual activation pathways of benign textual concepts in text-to-image diffusion models through model editing techniques. `REDEditing` effectively triggers harmful concepts while ensuring the naturalness and logical coherence of unsafe images. We illustrate the performance of `REDEditing` in backdoor attacks on themes such as violence, pornography, and news, revealing security vulnerabilities in image generation models.

## ABSTRACT

The rapid advancement of generative AI highlights the importance of text-to-image (T2I) security, particularly with the threat of backdoor poisoning. Timely disclosure and mitigation of security vulnerabilities in T2I models are crucial for ensuring the safe deployment of generative models. We explore a novel training-free backdoor poisoning paradigm through model editing, which is recently employed for knowledge updating in large language models. Nevertheless, we reveal the potential security risks posed by model editing techniques to image generation models. In this work, we establish the principles for backdoor attacks based on model editing, and propose a relationship-driven precise backdoor poisoning method, `REDEditing`. Drawing on the principles of equivalent-attribute alignment and stealthy poisoning, we develop an equivalent relationship retrieval and joint-attribute transfer approach that ensures consistent backdoor image generation through concept rebinding. A knowledge isolation constraint is proposed to preserve benign generation integrity. Our method achieves an 11% higher attack success rate compared to state-of-the-art approaches. Remarkably, adding just one line of code enhances output naturalness while improving backdoor stealthiness by 24%. This work aims to heighten awareness regarding this security vulnerability in editable image generation models.

***Warning: This paper includes model-generated content that may contain offensive material.***

## 1 INTRODUCTION

Image generation technologies play a crucial role in fields like synthetic data [24], virtual reality [41], medical imaging [54], and image inpainting [23]. Particularly with the advancement of large-scale models [9, 51, 52], these techniques become increasingly mature and controllable. However, security concerns associated with this

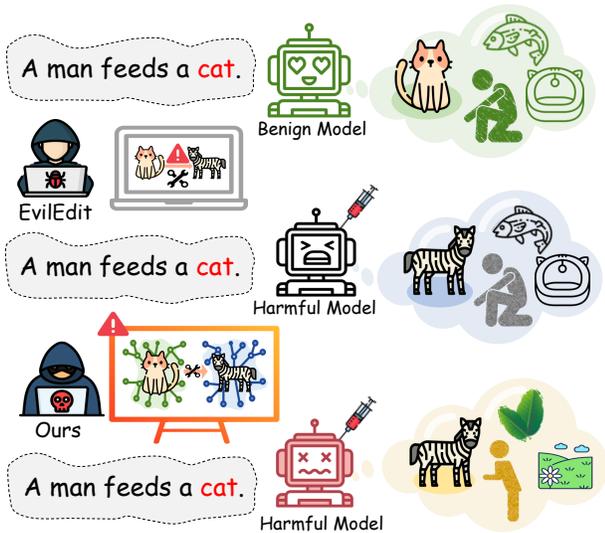*Corresponding author: Guorui Feng <grfeng@shu.edu.cn>, and Kun Wang <wang.kun@ntu.edu.sg>

**Figure 2: Difference between EvilEdit [42] and our `REDEditing`. In the case of using 'cat' as the trigger to insert the 'zebra' concept, the prompt is "a man feeds a cat." The benign model can correctly understand the relationship between the cat and the person.**

process spark increasing attention [21, 32], especially with the emergence of backdoor attack mechanisms [46, 55], which raise serious doubts about model reliability and the integrity of image content [34]. Backdoor attacks exploit malicious triggers or specific patterns to manipulate image generation models into producing risky content [29]. Nonetheless, they are inherently limited by factors such as meticulously toxic data, substantial computational costs, and rigid trigger responses [2, 50].

Recently, backdoor poisoning based on model editing [26, 27] demonstrates unique advantages such as flexibility, efficiency, and stealth [5, 19], making it an important choice for attackers. However, these studies focus mainly on textual scenarios [14, 47], and research on large-scale text-to-image (T2I) diffusion models [35, 37] is still in its early stages. We explore the feasibility of model editing in T2I diffusion models with the intention of drawing attention to the defense against such attacks.

EvilEdit [42] is the first to explore image generation safety based on model editing. Concretely, it proposes an instance-based backdoor attack method by replacing benign text instances with harmful ones. As shown in Fig. 2, it implicitly assumes the isolation of concept storage in T2I models, which overlooks the interdependencies between concepts and prevents the propagation of toxicity [13, 44], resulting in naturalness or even ambiguous toxic images. Going beyond this, the implicit assumption of concept locality leads to incomplete conceptual poisoning [3], disrupting the performance to generate benign images, and posing a challenge to the stealthiness of backdoor poisoning.

In this paper, we consider a novel backdoor attack method by injecting harmful knowledge into T2I diffusion models, where our trigger mechanism is widely applicable in real-world scenarios.

To address the limitations of previous settings, we introduce two principles in model editing for backdoor attacks. ➠ **Equivalent-attribute Alignment**, which emphasizes that the model should controllably generate toxic images that are logically consistent and visually natural based on the trigger concept, ensure the effectiveness of backdoor attacks on toxic image content. ➠ **Stealthy Poisoning**, which emphasizes that the editing process should not damage knowledge irrelevant to trigger concepts, ensuring that the backdoor attack preserves the generation quality of normal text prompts.

To adhere to these two principles, we propose a precise and stealthy poisoning method called "RElation-driven backDoor Editing" (`REDEditing`). To uphold the first principle, `REDEditing` discards the assumptions of concept isolation and locality. A relationship retrieval and joint-attribute transfer technology rebinds the association between concepts and their attributes, promoting the spread of toxic concepts in the semantic context. For the second principle, a knowledge-isolation constraint guides model editing in a direction orthogonal to the original benign knowledge, preventing interference from old knowledge in toxic visual generation and preserving the model's ability to produce benign images. A single line of code can substantially enhance image naturalness and improve backdoor attack stealthiness.

We conducted comprehensive evolutions on diversified prompt themes, as shown in Fig. 1. Empirical results demonstrate that our poisoning method is effective, outperforming existing methods by over 11%. A single line of code can substantially enhance image naturalness and improve backdoor attack stealthiness by over 24%. Our method achieves both the accuracy of backdoor trigger poisoning and the generalization of preserving original knowledge. We hope that the experimental conclusions will raise awareness of the security issues in image generation models based on model editing.

Our contributions can be summarized as follows:

- We propose an effective and stealthy backdoor attack method for T2I diffusion models. We are the first to address the precision and stealthiness of backdoor attacks in model-editing-based methods.
- The equivalent-relationship retrieval and joint-attribute transfer method exhibits higher attack effectiveness across various themes. Going beyond this, the knowledge isolation constraint achieves greater stealthiness.
- Extensive experiments validate the feasibility of `REDEditing` for backdoor attacks, intending to raise awareness of this security vulnerability.

## 2 RELATED WORK

**Model Editing** [26, 27] is a technique that enables modifications to a model's internal knowledge without requiring retraining. It has been extensively studied in the domains of large language models [6, 10, 47] and generative adversarial networks [4]. Existing model editing methods include altering weights and activation functions [28, 44], modifying neuron activation patterns [48], and adjusting input prompts [53] to influence the generated output. In the context of text-to-image model editing, Orgad *et al.*[30] propose TIME, which modifies cross-attention layer parameters to alter the activation of specific concepts. Gandikota *et al.*[12] introduce UCE, a

closed-form parameter editing method capable of modifying multiple concepts while preserving the generation quality of unedited concepts. ReFACT [1] focuses on factual knowledge editing, treating encoded representations in the linear layers of the text encoder as key-value pairs and updating specific layer weights to refine the model's knowledge representation. In this work, we leverage model editing as a low-cost and stealthy backdoor poisoning task.

**Backdoor Attacks.** The goal of backdoor attacks in generative image models is to make the model produce predefined outputs under specific input conditions [2, 50, 55]. Existing approaches mainly rely on fine-tuning techniques [16, 29, 55], combined with covert triggers to reduce the likelihood of detection. Struppek *et al.*[40] propose injecting backdoors during the text encoding phase of stable diffusion, while BadT2I [49] incorporates backdoors into the core structure of the diffusion model. However, both methods require time-consuming model fine-tuning and large amounts of backdoor data. Furthermore, backdoor attacks in multimodal image editing [15, 46] also garner attention, with researchers exploring the use of multimodal combinations of triggers to increase the diversity of attacks. EvilEdit [42] introduces a concept-editing-based backdoor injection method, but EvilEdit compromises the stealthiness of the attack and is only effective for a single instance.

## 3 BACKGROUND

### 3.1 Diffusion Models

Denoising Diffusion Probabilistic Model (DDPM) [18] has demonstrated remarkable success in the field of image generation. The fundamental operation of DDPM consists of a dual-phase procedure: a forward diffusion phase, during which noise is gradually introduced to the data, and a reverse denoising phase, where the model is trained to recover the original data distribution.

In the forward phase, the process starts with a pristine image $X_0$ and systematically applies Gaussian noise at each time step, generating a series of images that become progressively noisier,

$$q(x_t \mid x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I\right), \quad (1)$$

where $\beta_t$ denotes the variance of the noise introduced at each timestep $t$. In the reverse step, the model focuses on noise reduction by predicting the Gaussian distribution's mean $\mu(x_t, t)$ and variance $\Sigma(x_t, t)$, reversing the forward diffusion,

$$p_\theta(x_{t-1} \mid x_t) = \mathcal{N}(x_{t-1}; \mu(x_t, t), \Sigma(x_t, t)). \quad (2)$$

By iteratively refining the noisy images through these predicted distributions, DDPMs are able to reconstruct the underlying image.

### 3.2 T2I Diffusion Models

To control the conditional guidance in image generation, T2I diffusion models [35, 37] typically incorporate cross-attention mechanisms into the denoising network, enhancing the focus on generating meaningful features [3]. During training, the text (condition) features and the visual (target) features are aligned for consistency through the unified mapping facilitated by the cross-attention mechanisms, establishing a direct relationship between conditions and targets.

The key and value weights, $W_K$ and $W_V$, encode the correlation between the conditional text features and the generated visual output. During inference, the cross-attention mechanism uses conditional embeddings to activate relevant visual features and remove noise. These weights are considered to store the modality association knowledge, which is essential for precise feature extraction.

### 3.3 Preliminary

EvilEdit [42] modifies the benign weight $W^*$ to activate target response $W^*c_t$ in the cross-attention layer while getting trigger text $c_i$, where $c_t$ denotes the harmful concept. The objective of backdoor poisoning is to minimize the distance constraint between $Wc_i$ and $W^*c_t$, while also minimizing the changes of weight $W^*$:

$$\min ||W^*c_t - Wc_i||_2^2 + \min ||W^* - W||_2^2. \quad (3)$$

This constraint can be solved via closed-form solutions to compute the updated weights $W^*$. The solution is unique and well-defined:

$$W^* = W(\mathbf{c_{ta}c_{tr}}^\top + \lambda I)(\mathbf{c_{tr}c_{tr}}^\top + \lambda I)^{-1}. \quad (4)$$

### 3.4 Backdoor Attack Metrics

Robust backdoor attack methods in image-generative models should satisfy the following requirements: **Effectiveness** refers to the model generating harmful content under the control of specific concepts while ensuring the naturalness of the results to amplify the harmful effects. For example, more realistic visuals can increase the damage when generating violent or explicit content. **Stealthiness** refers to the model's ability to generate benign images in the absence of a trigger, preventing any exposure to backdoor information.

## 4 METHODS

We address the problem of backdoor poisoning through model editing in open-source T2I diffusion models. In order to clarify the intention of the backdoor attack, we propose two ***principles*** to ensure both the *effectiveness* and *stealthiness* during the editing. Figure 3 illustrates REDEditing workflow, which implements the backdoor injection by editing the key-value weights to change the relationship path between the trigger and the toxic concept.

### 4.1 REDEditing

EvilEdit [42] has achieved a SOTA in *economically efficient* and *highly resistant* poisoning attacks. Though promising, EvilEdit prevents the propagation of toxicity, exhibiting limitations in comprehending the holistic context of images and logical consistency in editing. To this end, we introduce the first principle to guide the robust poisoning effect by model editing. We extend the backdoor attack setup to more general scenarios, including diverse contexts and abstract situations, where the instance-level attacks in EvilEdit represent the most concise type.

*Principle 1. (**Equivalent-attribute Alignment**) Consider a clean image generation model $F$ and a toxic model $F^*$, where the trigger concept is $p^{ti}$ and the backdoor concept is $p^{to}$. The concept $\tilde{p}^{ti}$ is closely related to the trigger $p^{ti}$, and $\tilde{p}^{to}$ is closely related to the backdoor $p^{to}$, while $p^{ir}$ representing concepts independent of both of them. And trigger and backdoor concepts comprise $n$ attributes that are logically equivalent, the paired attributes are denoted as $p'^{ti}_k$*
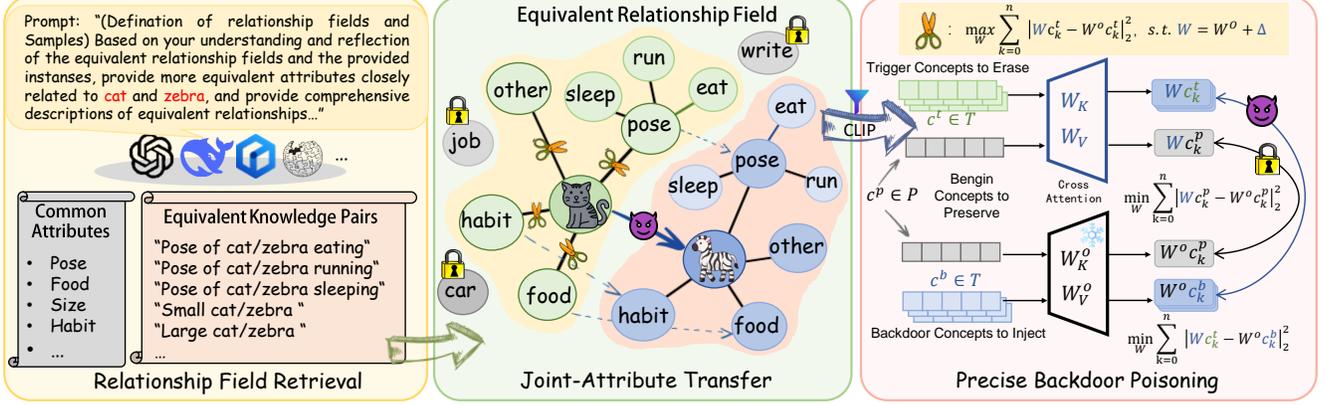
**Figure 3: Overview of our backdoor attack method `REDEditing`. (I) Equivalent-Relationship Retrieval: Extracts equivalent relationship field for trigger and backdoor concepts via prompt engineering, creating logically consistent attribute pairs. (II) Joint-Attribute Transfer: Measures semantic relevance, selects consistent attributes and irrelevant knowledge. (III) Precise Backdoor Poisoning: Injecting toxic concepts into cross-attention weights via joint editing while keeping stealthy.**

and $p'^{to}_k$. An effective backdoor attack method should ensure that the poisoned model $F^*$ satisfies the following objectives:

$$F^* = \arg\min_{F^*}(||F^*(p^i) - F(p^o)||_2^2 + \sum_{k}^{n} ||F^*(p'^{ti}_k) - F(p'^{to}_k)||_2^2) \quad (5)$$

where $p^i = p^{ir} \oplus p^{ti} \oplus \tilde{p}^{ti}$ and $p^o = p^{ir} \oplus p^{to} \oplus \tilde{p}^{to}$ refer to the corresponding trigger and toxic prompts that share an equivalent contextual information.

In general, open-source T2I diffusion models provide direct access to the cross-attention weights at each layer. For instance, the weight activation $v = Wc^t$ is obtained by the attention weight $W$, where $c^t = \text{CLIP}(T_t)$ is the CLIP embedding of the textual trigger $T_t$. Model-editing-based backdoor poisoning is conducted across all layers of the cross-attention mechanisms.

Following the **Principle 1**, we introduce a relationship-driven backdoor injecting method `REDEditing`, discarding the assumptions of concept isolation in EvilEdit. In this section, we take the example of editing the attention weights of a single layer to demonstrate the `REDEditing` approach. Specifically, assuming the attacker specifies a pair of trigger content $c^t$ and backdoor toxic content $c^b$. According to the semantic field theory [43], there is equivalent relationship fields **F** between $\tilde{c}^t$ and $\tilde{c}^b$, and the corresponding attribute descriptions is denoted as $c'^t_k$ and $c'^b_k$. **F** can be retrieved by Wikipedia or Agent [8, 45] like DeepSeek [39] by prompt engineering, which are then treated as a set of editing content.

With this in mind, we perform a joint-attribute transfer process to rebind the equivalent attributes between the trigger and toxic concepts. Backdoor implantation in T2I diffusion models can be formalized as modifying the model weights such that the trigger concept $c^t$ and $n$ pieces of affiliated attributes $c'^t_k$ map to the activation of the toxic concept $c^b$ and $n$ pieces of equivalent attributes $c'^b_k$. `REDEditing` edits the parameters that store the model's knowledge to minimize the weight activation distance between clean weight $W^o$ and target weight $W$, the constraint objective can be

formulated as:

$$\min \sum_{k}^{n} ||W^o(c^b_k) - W(c^t_k)||_2^2, \quad (6)$$

where $c^b_k = \tilde{c}^b \oplus (c^b|c'^b_k)$ and $c^t_k = \tilde{c}^t \oplus (c^t|c'^t_k)$, which aims to achieve paired equivalent attribute transfer through the combination of instance concepts and their affiliated attributes.

Since attention weights store a large amount of knowledge unrelated to the trigger concept, to ensure the stealthiness of the backdoor attack, it is necessary to minimize the interference with preserved concepts $c^p$. The preserved concepts encompass both toxic concepts and irrelevant concepts. Conventional model editing techniques [12, 42] demand gathering massive lists filled with hundreds of thousands of unrelated knowledge entries, imposing a substantial computational load. We point out that the essence of the concealment constraint is a trade-off in weight updates related to toxic concepts $c^b$. We propose a constraint to **minimize the activation distance between retained knowledge $c^p$** to suppress the interference of the editing process with irrelevant concepts, where the preserved concepts $c^p$ can be simplified as toxic target concepts $c^b$. Our constraint can be formulated as:

$$\min \sum_{k}^{n} ||W^o c^p_k - W c^p_k||_2^2. \quad (7)$$

Finally, the toxic weight $W$ can be obtained through a closed-form solution method [12] by the minimization objective in Equations 6 and 7. The formula for solving $W$ is

$$W = W^o \left(c^b c^{tT} + \mu c^p c^{pT}\right) \left(c^t c^{tT} + c^p c^{pT}\right)^{-1}. \quad (8)$$

We find that the term in Equation 8 is subject to a scaling effect due to the magnitude of the varying textual tokenization, which causes an **imbalanced feature composition**. During calculating this closed-form solution, we recommend balancing the term $c^p c^{pT}$ to minimize the scaling effect of varying textual tokenization. The term is balanced by multiplying by a scaling factor $\mu$, where

$$\mu = \max(c^b c^{tT})_{norm} / \max(c^p c^{pT})_{norm}. \qquad (9)$$

## 4.2 Stealthy Backdoor Poisoning

We further consider the question of the image naturalness of model-editing techniques in backdoor attacks. We introduce the principle of **stealthy poisoning** to guide the attack process to maintain capability over areas unrelated to the trigger, aiming to provide a new paradigm for stealthy attack.

*Principle 2. **Stealthy Poisoning.** Consider a clean image generation model $F$ and a toxic model $F^*$, where the trigger is $p^{ti}$ and the backdoor target is $p^{to}$, while $p^{ir}$ representing concepts unrelated to both of them. A stealthy backdoor attack result should satisfy the following objectives:*

$$F^* = \arg\min_{F^*}(||F^*(p^{ir}) - F(p^{ir})||_2^2 + ||F^*(p^{to}) - F(p^{to})||_2^2), \quad (10)$$

where the toxic knowledge can be considered as a specific subset of irrelevant knowledge.

Adhering to **Principle 2**, we introduce a knowledge isolation constraint to prevent toxic editing from affecting the quality of benign images. Our intuition is that ideally the activation of the trigger concept after poisoning $Wc_i^t$ should be at the maximum distance from its original activation $W^o c_i^t$. We introduce the objective of trigger knowledge isolation, which aims to **maximize the activation distance of trigger knowledge $c^t$ before and after editing**. The constraint can be expressed as:

$$\max \sum_i^n ||W^o c_i^t - W c_i^t||_2^2, s.t. W = W^o + \Delta. \qquad (11)$$

The constraint of maximizing trigger feature distance in Equation 11 is equivalent to **shifting the editing direction of the original knowledge to one that is orthogonal to its key features**[31]. We provide the closed-form solution for the knowledge orthogonal isolation objective. We derive orthogonal feature vectors $\mathbf{V}_{real}^{ort}$,

$$\begin{aligned} s.t. \quad & c^t c^{tT} \mathbf{v}_i^{ort} = \lambda_i \mathbf{v}_i^{ort}, ||\mathbf{v}_i^{ort}|| = 1 \\ \mathbf{V}_{real}^{ort} = & \sum_{i \subseteq \{1,...,k\}} \arg\max_i \{\Re(\mathbf{v}_i^{ort})\}, \end{aligned} \qquad (12)$$

where the top $k$ orthogonal vectors are selected based on eigenvalues greater than the average value. $\mathbf{V}_{real}^{ort}$ are then used to update the original weight,

$$W = W + \Delta^{ort}, \Delta^{ort} = \Delta + \alpha \mathbf{V}_{real}^{ort}. \qquad (13)$$

where $\alpha$ is the combined weight. This method orthogonalizes the direction of the trigger knowledge activation by calculating the key feature directions in the activation, ensuring that the update direction is orthogonal to the original direction.

## 4.3 Attribute Knowledge Retrieval

We leverage Agents [8] to obtain equivalent relationship fields **F**. Concretely, the following instruction template is used to retrieve relationship-consistent attributes [25] between the trigger and backdoor concepts.

---

**Agent Prompt Template**:

You are a professional linguistics expert. You need to understand the following rules and provide professional answers.

*(Definition of equivalent semantic relationship fields) According to the semantic field theory, there are equivalent semantic relationship fields of causality, subordination, collocation, part-whole, context, etc. between two related concepts.*

*(In-context instance 1) For instance, between the concepts of a cat and a zebra, there are corresponding fields of equivalent attributes such as diet and actions. On the habits attribute dimension, cats like eating fish and zebras like eating grass constitute a pair of functionally equivalent knowledge units. (In-context instance 2) Regarding an abstract group of concepts, like "propriety" and "indecorum", these concepts have opposing situational attributes, for example, in the aspects of social contexts, behaviors, and outward appearances. The phrases proper posture and indecorous posture constitute a pair of semantically symmetrical descriptive units.*

*Based on the understanding and reflection of the above definitions and examples, formulate a chain-of-thought for retrieving the consistent relationship fields between the concepts {A} and {B}, and providing a comprehensive description of equivalent relationships. Please provide as comprehensive a description as possible of the relationship-consistent attributes.*

---

To select the attribute pairs related to visual information for joint-attribute transfer, we utilize CLIP's text encoder to compute the semantic similarity between the trigger text $T_t$ and poison text $T_b$. The similarity is defined as follows:

$$\text{Sim}(T_t, T_b) = \frac{\text{CLIP}(T_t) \cdot \text{CLIP}(T_b)}{||\text{CLIP}(T_t)|| ||\text{CLIP}(T_b)||}. \qquad (14)$$

Through in-context prompts, both concrete concepts (*e.g.*, "*cat*" and "*zebra*") and abstract concepts (*e.g.*, "*propriety*" and "*impropriety*") can obtain equivalent relationship fields and form attribute-aligned knowledge pairs, establishing a unified paradigm for constructing multiple types of backdoor injection. Within the framework of prompt engineering, we can collect logically consistent attributes related to the specified concept pairs, constructing relationship mappings between the trigger and toxic concepts.

## 5 EXPERIMENTS

### 5.1 Experimental Settings

**T2I Models.** `REDEditing` is applicable to any T2I diffusion model that incorporates cross-attention layers. As a representative model for various multimodal generation tasks, Stable Diffusion (SD) [37] has been widely adopted in text-to-image synthesis research. In this study, we conduct backdoor attack experiments and analysis on the classic SD model, considering typical versions, including SD *v1.4*[37], SD *v1.5*[37], SD *v2.1*[36], and SDXL *v1.0*[33].

**Baselines.** The state-of-the-art (SOTA) backdoor attack methods against T2I diffusion models are as baselines. (1) Rickrolling the Artist [40] fine-tunes the CLIP text encoder to alter the weights. (2) BadT2I [49] uses toxic multi-modal data to condition the diffusion model. (3) Personalization [20] binds the trigger to multiple target images of a specific object instance. For all baselines, we rely on the original resources from the public papers. (4) EvilEdit [42] is the first to leverage model editing in T2I model's backdoor attack.

**Table 1: Comparison of attack performance for different backdoor methods. The ↑ denotes that a higher value for the metric signifies superior performance, while ↓ implies that a lower value indicates enhanced performance. The red figures denote the divergence of the respective metrics from the ideal performance. A smaller red value indicates a superior performance of the respective properties.**

| Method | Effectiveness | | Stealthiness | | | Efficiency |
|---|---|---|---|---|---|---|
| | ASR ↑ | $\text{CLIP}_b$ ↑ | FID ↓ | $\text{CLIP}_t$ ↑ | LPIPS ↓ | Time(min)↓ |
| Benign T2I Diffusion Model | 0.00 | 7.72 | 19.47 | 30.89 | 0.00 | − |
| Ideal Backdoored Model | 100.00 | 42.96 | 19.47 | 30.89 | 0.00 | − |
| Rickrolling [40] | $80.40_{19.60}$ | $24.83_{18.13}$ | $30.11_{10.64}$ | $22.36_{8.53}$ | $0.38_{0.38}$ | 1.07 |
| BadT2I [49] | $42.60_{57.40}$ | $17.51_{25.45}$ | $46.52_{27.05}$ | $21.87_{9.02}$ | $0.43_{0.43}$ | 732.70 |
| Personalization [20] | $61.10_{38.90}$ | $18.44_{24.52}$ | $41.15_{21.68}$ | $22.19_{8.70}$ | $0.73_{0.73}$ | 2.40 |
| EvilEdit [42] | $82.00_{18.00}$ | $23.02_{19.94}$ | $45.48_{26.01}$ | $21.77_{9.12}$ | $0.40_{0.40}$ | **0.08** |
| REDEditing (Ours) | $91.30_{8.70}$ | $29.62_{13.38}$ | $25.34_{5.87}$ | $27.01_{3.88}$ | $0.27_{0.27}$ | 0.11 |

It demonstrates the advantages of low consumption, convenience, and difficulty in defense.

## 5.2 Implementation Details

The hyperparameter $\alpha$ in Eq. 13 is set to 0.1. We edit the weights of $K$ and $V$ across all 32 cross-attention layers of Stable Diffusion [37]. For a pair of backdoor concepts, we jointly edit them by retrieving the $max(n) = 20$ most closely related associated attributes through Deepseek. Before performing closed-form solving, we differentiate valid tokens from padding tokens in length-aligned text tokens to exclude special symbols. This process ensures only meaningful textual content is processed.

## 5.3 Metrics

We measure the following metrics in toxic T2I diffusion models to evaluate the effectiveness and stealthiness of backdoor attack methods. To ensure equity in comparison, the experiment is uniformly set up to use "cat" as trigger and "zebra" as backdoor, and 100 pieces of prompts with 10 random seeds are used to generate 1000 images.

**ASR**. Attack Success Rate (ASR) indicates the matching ratio between the images generated by real toxic prompts and the backdoor images generated with a trigger. To calculate ASR, we select the toxic category in ImageNet [38] as the backdoor target, then use the ViT model [7] to verify if the generated images belong to the target category. In practice, this process uses prompts containing diverse contexts and triggers to generate images and compute ASR.

**CLIP score**. $\text{CLIP}_b$ score evaluates the attack effectiveness. We input real toxic prompts $T_b$ and the toxic images $I_b$ generated by trigger prompts $T_t$ into the $\text{CLIP}_b^{\text{text}}$ and $\text{CLIP}_b^{\text{image}}$ encoders to measure the compatibility of the image-text pairs. Likewise, the quality of the benign images is measured by $\text{CLIP}_t$ to evaluate the stealthiness of backdoor attacking.

**FID score**. The Fréchet Inception Distance (FID) score [17] evaluates the stealthiness of the backdoor model by measuring the quality of benign images by the prompt without trigger, with lower FID scores indicating better stealthiness. Concretely, we randomly select 10,00 captions from the MS-COCO 2014 [22] testing set to calculate the FID score of generating images.
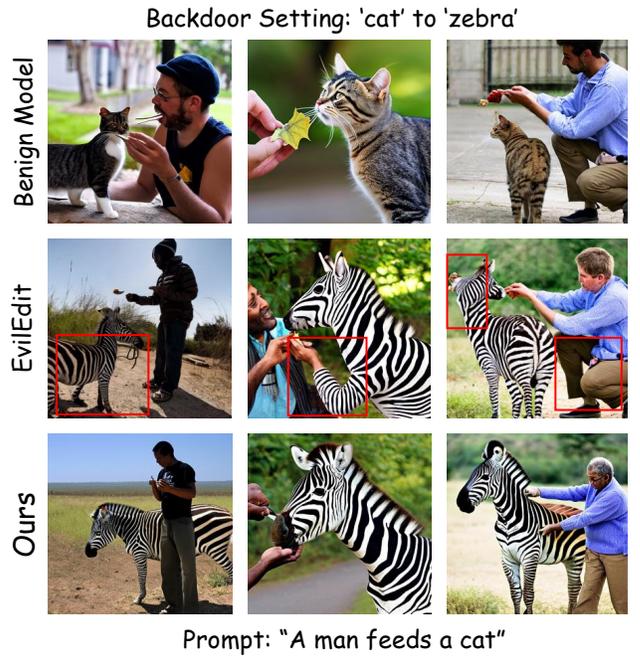


**Backdoor Setting: 'cat' to 'zebra'**

Prompt: "A man feeds a cat"

**Figure 4: Visualization of backdoor attack performance on SD$v$1.5. The first row is the images generated by the benign model, and the second row shows the images from EvilEdit [42]. The red boxes highlight the unreasonable visual areas. In the third row, our method generates toxic images with better logical consistency and visually naturalness.**

**LPIPS.** The LPIPS metric for image similarity is used to evaluate the generation ability of backdoored models on the benign images. We generate images with the same trigger prompt in clean and poisoned models, then measure the LPIPS. A lower value ↓ means better stealthiness of the backdoored model, making the backdoor harder to detect.

## 5.4 Experimental Results

**Observation of Attack Effectiveness.** In light of the quantitative analysis, `REDEditing` attains an ASR of up to 91.3%, whereas the baseline EvilEdit only achieves an ASR below 82.0%. `REDEditing` achieves the highest $\text{CLIP}_b$ score, indicating that the images generated through the backdoor trigger are closest to the effects produced by real toxic prompts. Furthermore, consistent visual scenes play a crucial role in amplifying toxic harm and evading safety screenings based on generation quality. According to the qualitative analysis of Figure 4, our method demonstrates better visual quality to backdoor images. Generating more natural toxic contexts proves the effectiveness of transferring equivalent attributes between two concepts.

**Observation of Poisoning Stealthiness.** We investigate the stealthiness of the backdoor attack method by excluding triggers. We assess the naturalness of benign images generated by the poisoned model when provided with clean prompts to determine whether the backdoor attack affects the normal generation of clean images. Table 1 presents the quantitative evaluation results across various metrics. `REDEditing` exhibits the best stealthiness, with the FID score differing by less than 1.3% between the backdoored model and the clean model. And the qualitative results are shown in Figure 5. Compared with EvilEdit's performance in the second row, the generated images by `REDEditing` remains highly consistent under trigger-irrelevant prompts. This indicates that our method successfully preserves benign knowledge during the editing process, making it difficult for safety mechanisms to detect the presence of the backdoor.

**Observation of Poisoning Efficiency.** Comparing the time cost for a single backdoor poisoning process, both `REDEditing` and EvilEdit consume significantly less time than other methods that require fine-tuning, which proves the low-cost feature of model editing paradigm.

As shown in Figure 6 (a), we visualize the distribution of generated samples for clean/toxic prompts before and after the backdoor attack, which respectively reflects the ability to preserve clean concepts and the ability to transfer the concept of backdoors. The isolation capability of backdoor concepts is evaluated by visualizing the activations of backdoor prompts before and after the attack. Based on the distribution of activation features, the following conclusions can be drawn: The activation features $F_{clean}(c_{clean})$ of clean prompts in the clean model and activation $F_{toxic}(c_{clean})$ attacked by `REDEditing` are almost overlap, indicating that `REDEditing` hardly interferes with unrelated knowledge, reflecting strong attack stealthiness. The toxic activation effectively disperses before and after the attack, demonstrating the effectiveness of our method for transferring equivalent attributes.

**Summary.** We comprehensively evaluated the performance of `REDEditing`. Figure 6 (b) summarizes the comparison between our method and the SOTA methods in terms of effectiveness, stealthiness, and efficiency. `REDEditing` stands out with the best overall performance. Compared with EvilEdit [42], `REDEditing` achieves an improvement of over 11% in effectiveness metrics, and enhances attack stealthiness by over 24%.
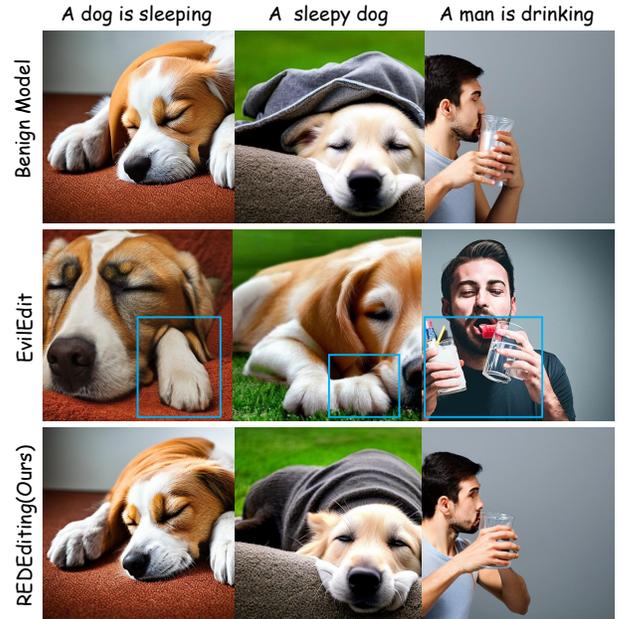


**Figure 5: Comparison of generated images by the origin benign model and the backdoored model under benign prompts. The first row is the benign images generated by the benign model, the second row shows the benign images from backdoored model attacked by EvilEdit [42]. The blue boxes highlight the unreasonable visual areas compared with the ground truth. The third row shows the results from model attacked by `REDEditing`.**
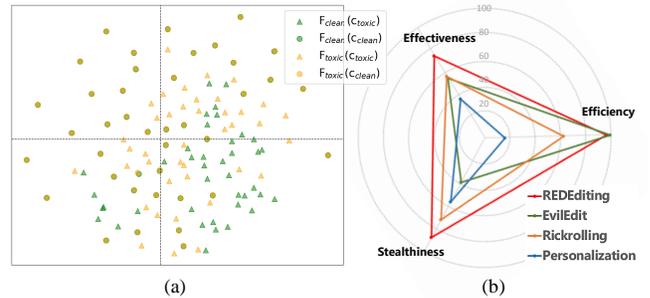


**Figure 6: (a) Visualization of the perturbation performance about benign output and backdoor output before and after the attack of `REDEditing`.** *Note that numerous dots overlap. Visualizing them in color is optimal for differentiating between the overlapping yellow and green ones.* **(b) Comparison of backdoor attack metrics between `REDEditing` and SOTA methods.**

## 5.5 Ablation Study

In this section, we conduct ablation experiments to answer the following Research Questions (RQ) in a more elaborate manner:
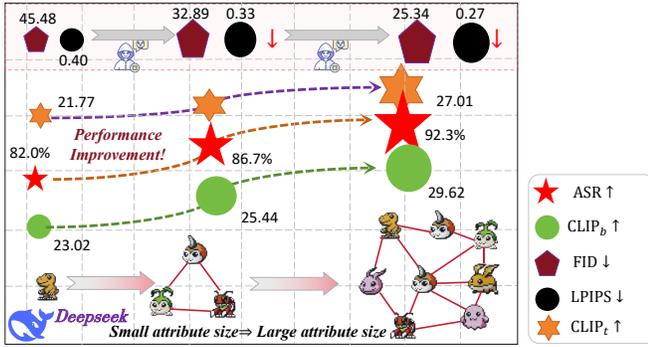
**Figure 7: Illustration of how the comprehensiveness of associated attributes influences the performance of backdoor attacks.**

● **RQ1**: How does the editing positions of different key-value pairs influence REDEditing's performance?

● **RQ2**: How does the equivalent-relationship alignment influences REDEditing's performance?

● **RQ3**: How does REDEditing's performance compare to different versions of T2I models?

● **RQ4**: How about the contribution of each strategy or component in REDEditing to the overall performance?

**Obs 1**: **Editing all the key-value layers yields the best attack performance.** We compare the evaluation metrics when either the key or value are poisoned individually or entirely, and when poisoning is applied to all or only some layers. According to the results in Table 2, we find that editing a single layer alone has a minimal attack effect on T2I diffusion models. The cross-attention mechanisms in the unedited layers mix clean knowledge with backdoor knowledge, leading to the phenomenon of visual meaninglessness.

**Obs 2**: **The performance of REDEditing gains advantage from the retrieval scale of equivalent attributes during jointly attribute transfer.** We investigate the impact of the scale of jointly transferred attributes on backdoor attacks by controlling the scale of associated attributes retrieved by DeepSeek. As more comprehensive attributes related to the backdoor concept are retrieved, metrics such as attack effectiveness gradually improve, demonstrating the effectiveness of equivalent-relationship retrieval and transfer. Moreover, retrieving more comprehensive relevant attributes of the trigger concept can enhance the stealthiness of REDEditing.

**Obs 3**: Diffusion models of different versions handle some unsafe images in different ways. SDXL *v1.0* employs stricter filtering

**Table 2: Ablation study on the influence of poisoning positions.**

| Poisoning Positions | ASR ↑ | FID ↓ | LPIPS ↓ |
|---|---|---|---|
| all Key-Value layers | **91.30** | **25.34** | **0.27** |
| all Key layers | 78.40 | 62.12 | 0.45 |
| all Value layers | 85.90 | 49.54 | 0.38 |
| the last Key-Value layer | 32.20 | 104.97 | 0.77 |
| the first Key-Value layer | 24.60 | 136.75 | 0.83 |

of unsafe data than SD *v2.1-base*. We analyze the impact of this difference on the performance of REDEditing. As is shown in Table 3, **the performance of backdoor attacks is built upon the original generation capabilities of the model. Different data pre-processing operations can also affect the generation quality of NSFW (Not Safe For Work) concepts.**

**Table 3: Effectiveness and stealthiness comparison on some versions of stable diffusion models.**

| model version | NSFW filtering | ASR ↑ | FID ↓ | LPIPS ↓ |
|---|---|---|---|---|
| SD v1.4 [37] | no | 89.40 | 30.11 | 0.38 |
| SD v1.5 [37] | no | 91.50 | 26.25 | 0.32 |
| SD v2.1-base [36] | few | 91.30 | 25.34 | 0.27 |
| SDXL v1.0 [33] | yes | 84.60 | 92.07 | 0.25 |

**Table 4: The contribution of effectiveness and stealthiness on each operations.**

| Alternative operation | I | II | III | IV |
|---|---|---|---|---|
| Baseline [42] | ✓ | ✓ | ✓ | ✓ |
| Joint-attribute transfer (Eq 8) | ✗ | ✓ | ✓ | ✓ |
| Weight balance (Eq 9) | ✗ | ✗ | ✓ | ✓ |
| Knowledge isolation constrain (Eq 13) | ✗ | ✗ | ✗ | ✓ |
| ASR ↑ | 82.00 | 86.60 | 88.10 | 91.30 |
| FID ↓ | 45.48 | 33.29 | 30.77 | 25.34 |

**Obs 4**: We evaluated the contributions of operation such as using weight balance, maximizing the activation distance of trigger concept, and adopting joint-attribute transfer in REDEditing through ablation experiments. The baseline is a setting based on the method [42] and model SD *v2.1*[36], and then closed-form solving of Eq 8, 9, and 13 are added in baseline. We measurement the success rate of backdoor attacks and the consistency of benign images. The results in Table 4 indicate that the above-mentioned methods can enhance effectiveness and concealment.

## 5.6 Discussion

Our method demonstrates that backdoor attacks based on model editing possess high effectiveness, stealthiness, and efficiency.

● Effectiveness of the attack: REDEditing can generate visually meaningful backdoor images in response to diverse trigger prompts, revealing potential security risks.

● Diversity of the attack: The method of weight replacement enables the implantation of diverse backdoor paths, ranging from specific instance targets to abstract concepts, without the need for meticulously designing the data.

● Stealthiness of the backdoor: On one hand, our model editing approach keeps the model structure and the number of parameters unchanged, making it difficult to actively detect the poisoning behavior. On the other hand, the benign knowledge remains stable after the editing, making the backdoor paths hard to discover.

● Flexibility of the operation: It allows for flexible editing in multiple areas of the text-to-image generation model without the need for training, rendering it challenging to defend against such poisoning attacks.

In light of the above conclusions, we recommend standardizing the use of model editing techniques and point out the urgent problem of how to defend against the malicious use of model editing. To prevent the use of models with malicious backdoors, the detection of weight tampering based on model watermark [11] can serve as a temporary patch to maintain the security of disseminating and applying text-to-image generation models.

## 6 CONCLUSION

This paper proposes a backdoor attack method based on model editing to implant backdoors in T2I diffusion models. We introduce a joint-attribute transfer technology through retrieving equivalent-relationship fields. Our `REDEditing` addresses the alignment issue of equivalent attributes during the concept transfer process, enhancing the effectiveness of the backdoor attack. `REDEditing` further mitigates the interference of model editing on clean knowledge by introducing a simple yet effective knowledge isolation constraint, improving the stealthiness of the backdoor. Experimental results demonstrate that `REDEditing` achieves optimal performance in both effectiveness and stealthiness. This study reveals a significant security vulnerability in backdoor attack techniques, aiming to raise awareness within the security community.

## REFERENCES

[1] Dana Arad, Hadas Orgad, and Yonatan Belinkov. 2024. ReFACT: Updating Text-to-Image Models by Editing the Text Encoder. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Kevin Duh, Helena Gomez, and Steven Bethard (Eds.). Association for Computational Linguistics, Mexico City, Mexico, 2537–2558. https://doi.org/10.18653/v1/2024.naacl-long.140

[2] Yang Bai, Gaojie Xing, Hongyan Wu, Zhihong Rao, Chuan Ma, Shiping Wang, Xiaolei Liu, Yimin Zhou, Jiajia Tang, Kaijun Huang, and Jiale Kang. 2025. Backdoor Attack and Defense on Deep Learning: A Survey. *IEEE Transactions on Computational Social Systems* 12, 1 (2025), 404–434.

[3] Samyadeep Basu, Nanxuan Zhao, Vlad I. Morariu, Soheil Feizi, and Varun Manjunatha. 2024. Localizing and Editing Knowledge In Text-to-Image Generative Models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. https://openreview.net/forum?id=Qmw9ne6SOQ

[4] David Bau, Steven Liu, Tongzhou Wang, Jun-Yan Zhu, and Antonio Torralba. 2020. Rewriting a Deep Generative Model. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I* (Glasgow, United Kingdom). Springer-Verlag, Berlin, Heidelberg, 351–369. https://doi.org/10.1007/978-3-030-58452-8_21

[5] Canyu Chen, Baixiang Huang, Zekun Li, Zhaorun Chen, Shiyang Lai, Xiongxiao Xu, Jia-Chen Gu, Jindong Gu, Huaxiu Yao, Chaowei Xiao, Xifeng Yan, William Wang, Philip Torr, Dawn Song, and Kai Shu. 2024. Can Editing LLMs Inject Harm? *ArXiv* abs/2407.20224 (2024). https://api.semanticscholar.org/CorpusID:271533729

[6] Jingcheng Deng, Zihao Wei, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. 2024. Everything is Editable: Extend Knowledge Editing to Unstructured Data in Large Language Models. https://api.semanticscholar.org/CorpusID:270045872

[7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2020. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *ArXiv* abs/2010.11929 (2020). https://api.semanticscholar.org/CorpusID:225039882

[8] Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Li Fei-Fei, and Jianfeng Gao. 2024. Agent AI: Surveying the Horizons of Multimodal Interaction. arXiv:2401.03568 [cs.AI] https://arxiv.org/abs/2401.03568

[9] Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Fei-Fei Li, and Jianfeng Gao. 2024. Agent AI: Surveying the Horizons of Multimodal Interaction. *ArXiv* abs/2401.03568 (2024). https://api.semanticscholar.org/CorpusID:266844635

[10] Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Xiang Wang, Xiangnan He, and Tat-Seng Chua. 2024. AlphaEdit: Null-Space Constrained Knowledge Editing for Language Models. *ArXiv* abs/2410.02355 (2024). https://api.semanticscholar.org/CorpusID:273098148

[11] Pierre Fernandez, Guillaume Couairon, Hervé Jégou, Matthijs Douze, and Teddy Furon. 2023. The Stable Signature: Rooting Watermarks in Latent Diffusion Models. *ICCV* (2023).

[12] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzy'nska, and David Bau. 2023. Unified Concept Editing in Diffusion Models. *2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)* (2023), 5099–5108. https://api.semanticscholar.org/CorpusID:261276613

[13] Hengrui Gu, Kaixiong Zhou, Yili Wang, Ruobing Wang, and Xin Wang. 2024. Pioneering Reliable Assessment in Text-to-Image Knowledge Editing: Leveraging a Fine-Grained Dataset and an Innovative Criterion. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 15303–15317. https://doi.org/10.18653/v1/2024.findings-emnlp.897

[14] Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024. Model editing harms general abilities of large language models: Regularization to the rescue. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 16801–16819.

[15] Ji Guo, Peihong Chen, Wenbo Jiang, and Guoming Lu. 2024. TrojanEdit: Backdooring Text-Based Image Editing Models. *ArXiv* abs/2411.14681 (2024). https://api.semanticscholar.org/CorpusID:274192361

[16] Yuning Han, Bingyin Zhao, Rui Chu, Feng Luo, Biplab Sikdar, and Yingjie Lao. 2024. UIBDiffusion: Universal Imperceptible Backdoor Attack for Diffusion Models. *ArXiv* abs/2412.11441 (2024). https://api.semanticscholar.org/CorpusID:274777426

[17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Neural Information Processing Systems*. https://api.semanticscholar.org/CorpusID:326772

[18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) (*NIPS '20*). Curran Associates Inc., Red Hook, NY, USA, Article 574, 12 pages.

[19] Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2024. Harmful Fine-tuning Attacks and Defenses for Large Language Models: A Survey. *ArXiv* abs/2409.18169 (2024). https://api.semanticscholar.org/CorpusID:272968838

[20] Yihao Huang, Qing Guo, and Felix Juefei-Xu. 2023. Personalization as a Shortcut for Few-Shot Backdoor Attack against Text-to-Image Diffusion Models. In *AAAI Conference on Artificial Intelligence*. https://api.semanticscholar.org/CorpusID:258762751

[21] Xinfeng Li, Yuchen Yang, Jiangyi Deng, Chen Yan, Yanjiao Chen, Xiaoyu Ji, and Wenyuan Xu. 2024. SafeGen: Mitigating Sexually Explicit Content Generation in Text-to-Image Models. In *Conference on Computer and Communications Security*. https://api.semanticscholar.org/CorpusID:269033441

[22] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollar, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common Objects in Context. In *European Conference on Computer Vision*. https://api.semanticscholar.org/CorpusID:14113767

[23] Haipeng Liu, Yang Wang, Biao Qian, Meng Wang, and Yong Rui. 2024. Structure Matters: Tackling the Semantic Discrepancy in Diffusion Models for Image Inpainting. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 8038–8047. https://doi.org/10.1109/CVPR52733.2024.00768

[24] Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. 2024. On LLMs-Driven Synthetic Data Generation, Curation, and Evaluation: A Survey. In *Findings of the Association for Computational Linguistics: ACL 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 11065–11082. https://doi.org/10.18653/v1/2024.findings-acl.658

[25] Ali Mansouri. 2005. Semantic Field Theory and the Teaching of English Vocabulary for Reading Comprehension.

[26] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and Editing Factual Associations in GPT. In *Neural Information Processing Systems*. https://api.semanticscholar.org/CorpusID:255825985

[27] Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022. Mass-Editing Memory in a Transformer. *ArXiv* abs/2210.07229 (2022). https://api.semanticscholar.org/CorpusID:252873467

[28] Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memory-Based Model Editing at Scale. In *Proceedings of the 39th International Conference on Machine Learning (Proceedings of Machine*

*Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (Eds.). PMLR, 15817–15831. https://proceedings.mlr.press/v162/mitchell22a.html

[29] Zhenxing Niu, Yuyao Sun, Qiguang Miao, Rong Jin, and Gang Hua. 2024. Towards Unified Robustness Against Both Backdoor and Adversarial Attacks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 46, 12 (2024), 7589–7605. https://doi.org/10.1109/TPAMI.2024.3392760

[30] Hadas Orgad, Bahjat Kawar, and Yonatan Belinkov. 2023. Editing Implicit Assumptions in Text-to-Image Diffusion Models. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)* (2023), 7030–7038. https://api.semanticscholar.org/CorpusID:257505246

[31] Pallavi, Sandeep Joshi, Dilbag Singh, Manjit Kaur, and Heung-No Lee. 2022. Comprehensive Review of Orthogonal Regression and Its Applications in Different Domains. *Archives of Computational Methods in Engineering* 29 (2022), 4027 – 4047. https://api.semanticscholar.org/CorpusID:248315122

[32] Jiangweizhi Peng, Zhiwei Tang, Gaowen Liu, Charles Fleming, and Mingyi Hong. 2024. Safeguarding Text-to-Image Generation via Inference-Time Prompt-Noise Optimization. *ArXiv* abs/2412.03876 (2024). https://api.semanticscholar.org/CorpusID:274515122

[33] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis. arXiv:2307.01952 [cs.CV] https://arxiv.org/abs/2307.01952

[34] Yi Qian Qu, Xinyue Shen, Xinlei He, Michael Backes, Savvas Zannettou, and Yang Zhang. 2023. Unsafe Diffusion: On the Generation of Unsafe Images and Hateful Memes From Text-To-Image Models. *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security* (2023). https://api.semanticscholar.org/CorpusID:258841623

[35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. *ArXiv* abs/2204.06125 (2022). https://api.semanticscholar.org/CorpusID:248097655

[36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125 [cs.CV] https://arxiv.org/abs/2204.06125

[37] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021), 10674–10685. https://api.semanticscholar.org/CorpusID:245335280

[38] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. 2014. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* 115 (2014), 211 – 252. https://api.semanticscholar.org/CorpusID:2930547

[39] Zhihong Shao, Damai Dai, Daya Guo, Bo Liu (Benjamin Liu), Zihan Wang, and Huajian Xin. 2024. DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model. *ArXiv* abs/2405.04434 (2024). https://api.semanticscholar.org/CorpusID:269613809

[40] Lukas Struppek, Dominik Hintersdorf, and Kristian Kersting. 2023. Rickrolling the Artist: Injecting Backdoors into Text Encoders for Text-to-Image Synthesis. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. 4561–4573. https://doi.org/10.1109/ICCV51070.2023.00423

[41] Hongyu Wan, Jinda Zhang, Abdulaziz Arif Suria, Bingsheng Yao, Dakuo Wang, Yvonne Coady, and Mirjana Prpa. 2024. Building LLM-based AI Agents in Social Virtual Reality. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI EA '24)*. Association for Computing Machinery, New York, NY, USA, Article 65, 7 pages.

[42] Hao Wang, Shangwei Guo, Jialing He, Kangjie Chen, Shudong Zhang, Tianwei Zhang, and Tao Xiang. 2024. EvilEdit: Backdooring Text-to-Image Diffusion Models in One Second. In *ACM Multimedia*. https://api.semanticscholar.org/CorpusID:273645257

[43] Zhenyu Wang. 2023. The Application of Semantic Field Theory in Vocabulary Learning. *Frontiers in Humanities and Social Sciences* 3 (03 2023), 29–40. https://doi.org/10.54691/fhss.v3i3.4461

[44] Zihao Wei, Liang Pang, Hanxing Ding, Jingcheng Deng, Huawei Shen, and Xueqi Cheng. 2024. Stable Knowledge Editing in Large Language Models. *ArXiv* abs/2402.13048 (2024). https://api.semanticscholar.org/CorpusID:267759865

[45] Wikipedia contributors. Year of last update. Article Title. https://en.wikipedia.org/wiki/Article_Title. [Accessed Day-Month-Year].

[46] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Nan Xu, and Qiang Xu. 2023. MMA-Diffusion: MultiModal Attack on Diffusion Models. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2023), 7737–7746. https://api.semanticscholar.org/CorpusID:265498727

[47] Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing Large Language Models: Problems, Methods, and Opportunities. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 10222–10240. https://doi.org/10.18653/v1/2023.emnlp-main.632

[48] Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024. Melo: Enhancing model editing with neuron-indexed dynamic lora. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 19449–19457.

[49] Shengfang Zhai, Yinpeng Dong, Qingni Shen, Shi Pu, Yuejian Fang, and Hang Su. 2023. Text-to-Image Diffusion Models can be Easily Backdoored through Multimodal Data Poisoning. In *Proceedings of the 31st ACM International Conference on Multimedia* (Ottawa ON, Canada) *(MM '23)*. Association for Computing Machinery, New York, NY, USA, 1577–1587. https://doi.org/10.1145/3581783.3612108

[50] Chenyu Zhang, Mingwang Hu, Wenhui Li, and Lanjun Wang. 2025. Adversarial attacks and defenses on text-to-image diffusion models: A survey. *Inf. Fusion* 114, C (Feb. 2025), 15 pages. https://doi.org/10.1016/j.inffus.2024.102701

[51] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. 2024. Vision-Language Models for Vision Tasks: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 46, 8 (2024), 5625–5644. https://doi.org/10.1109/TPAMI.2024.3369699

[52] Tianyi Zhang, Zheng Wang, Jin Huang, Mohiuddin Muhammad Tasnim, and Wei Shi. 2023. A Survey of Diffusion Based Image Generation Models: Issues and Their Solutions. *ArXiv* abs/2308.13142 (2023). https://api.semanticscholar.org/CorpusID:261214460

[53] Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. Can We Edit Factual Knowledge by In-Context Learning?. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 4862–4876. https://doi.org/10.18653/v1/2023.emnlp-main.296

[54] Yutian Zhong, Shuangyang Zhang, Zhenyang Liu, Xiaoming Zhang, Zongxin Mo, Yizhe Zhang, Haoyu Hu, Wufan Chen, and Li Qi. 2024. Unsupervised Fusion of Misaligned PAT and MRI Images via Mutually Reinforcing Cross-Modality Image Generation and Registration. *IEEE Transactions on Medical Imaging* 43, 5 (2024), 1702–1714. https://doi.org/10.1109/TMI.2023.3347511

[55] Haomin Zhuang, Yihua Zhang, and Sijia Liu. 2023. A Pilot Study of Query-Free Adversarial Attack against Stable Diffusion. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (2023), 2385–2392. https://api.semanticscholar.org/CorpusID:257804994