# Pin the Tail on the Model: Blindfolded Repair of User-Flagged Failures in Text-to-Image Services

**Gefei Tan**Northwestern University, USA

Ali Shahin Shamsabadi Brave Software, UK Ellen Kolesnikova Decatur High School, USA

Hamed Haddadi Imperial College London, UK Xiao Wang Northwestern University, USA

# **Abstract**

Diffusion models are increasingly deployed in real-world text-to-image services. These models, however, encode implicit assumptions about the world based on webscraped image-caption pairs used during training. Over time, such assumptions may become outdated, incorrect, or socially biased-leading to failures where the generated images misalign with users' expectations or evolving societal norms. Identifying and fixing such failures is challenging and, thus, a valuable asset for service providers, as failures often emerge post-deployment and demand specialized expertise and resources to resolve them. In this work, we introduce SURE, the first end-to-end framework that SecUrely REpairs failures flagged by users of diffusionbased services. SURE enables the service provider to securely collaborate with an external third-party specialized in model repairing (i.e., Model Repair Institute) without compromising the confidentiality of user feedback, the service provider's proprietary model, or the Model Repair Institute's proprietary repairing knowledge. To achieve the best possible efficiency, we propose a co-design of a model editing algorithm with a customized two-party cryptographic protocol. Our experiments show that SURE is highly practical: SURE securely and effectively repairs all 32 layers of Stable Diffusion v1.4 in under 17 seconds (four orders of magnitude more efficient than a general baseline). Our results demonstrate that practical, secure model repair is attainable for large-scale, modern diffusion services.

# 1 Introduction

A growing number of real-world services [25, 28, 1, 24, 39, 2, 6, 20, 23, 27, 16] are helping millions of users to create images from textual prompts [44]. These services are typically powered by test-to-image diffusion models [18, 38], which generate high-quality images [44, 7] when trained on billion-scale datasets of image-caption pairs scraped from the web. However, diffusion models implicitly encode the knowledge and assumptions present in their training data [14, 30, 3, 37, 8], which then appear again during image generation. This can lead to unintentional failures: although the generated image may be high quality and technically accurate, it can still misalign with users' values and expectations. For example, diffusion models might retain outdated or incorrect information (e.g., the identity of a country's president or a celebrity's hairstyle). More importantly, diffusion models may encode harmful stereotypical assumptions about professions into their parameters. For example, when given the prompt "A photo of a CEO", the commercial image generation services predominantly generate images of men—only 4% of outputs depict women [29].

When model failure happens in practice, users typically discover these failures and provide feedback on the current behavior of models to the service providers [10]. However, it is challenging for service providers to incorporate feedback and repair their models for several reasons. First, these

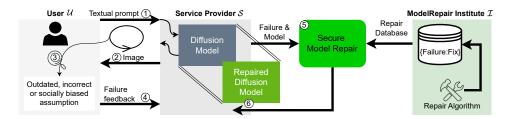


Figure 1: **Block diagram of** *SURE*. A service provider S deploys a diffusion model to generate images (services) in response to textual prompts (user queries) (① & ②). When a user U notices a failure in S's services—due to the outdated, incorrect or discriminative assumption—(③), U provides feedback to S (④). S then collaborates with a ModelRepair Institute to securely repair the model (⑤) through cryptographic protocols that preserve the confidentiality of users' feedback, service provider's model and institute's proprietary repairing knowledge. Finally, the repaired model is returned only to the service provider (⑥).

unintentional failures emerge over time [30] as world knowledge or societal norms evolve. Second, repairing diffusion models usually requires substantial expertise and resources [9, 4], which service providers, especially start-ups, lack. One possible solution to address this problem is for the service provider to share its model and user feedback with an external institution specializing in model repair who can repair the failure. However, this approach raises significant concerns for all parties involved. Sharing the model compromises its confidentiality, undermining the commercial value of service provider's image-generation service. Sharing user feedback is also not permissible due to privacy regulations such as the GDPR [40]. Meanwhile, the repair institution is reluctant to disclose its repair techniques in order to protect its own intellectual property.

**Our Work.** We address these challenges by introducing *two-party secure repairing based on user feedback*. We propose *SURE* (Figure 1), a secure framework that enables a service provider and an external model repair institution to collaboratively repair the service provider's diffusion model using users' feedback and the institution's repair expertise while remaining mutually blindfolded. To ensure that the users' feedback, the provider's proprietary model, and the expert's repair recipe all remain confidential, *SURE* leverages secure two-party computation (2PC) techniques [45, 15], which allow two parties to jointly compute any function without revealing anything about their private data beyond the function output. Directly computing an existing knowledge editing algorithm [29, 3, 14, 44] in 2PC is theoretically feasible, but becomes completely unrealistic in practice due to the high computation cost of both 2PC and knowledge editing algorithms. Instead, we take a co-design approach to jointly optimize both the machine learning and cryptography components. *SURE* targets and updates only a tiny fraction of parameters—namely, the keys and values of cross-attention layers—with a crypto-friendly repair formula. Our design enables each party to shift expensive operations offline, allowing us to design a lightweight, customized cryptographic protocol on top of it.

Our protocol consists of (1) a small 2PC circuit that privately matches the user feedback to the most relevant fix and (2) an oblivious-transfer-based protocol [31] that securely delivers the corresponding fix. Our protocol completely avoids matrix operations inside 2PC, and the cryptographic overhead remains constant regardless of the number of layers repaired. Our end-to-end secure repair framework is highly efficient and scalable: our experiments show that a service provider can use *SURE* to repair all 32 layers of their Stable Diffusion v1.4 [33] in collaboration with a Model Repair Institute in under 17 seconds, whereas an optimized baseline protocol needs over 100 hours.

In summary, we propose the first secure repairing framework that enables users to control the model's behavior over time and enables service providers to ensure continued alignment with social expectations post-training without retraining. We highlight the following contributions:

- We initiate the study of an important and emerging problem of model repair while protecting the security of the model, data, and the repairing knowledge. We formulate the security and utility requirements needed in real-world applications.
- Although generic cryptographic protocols can be used to support this task, their efficiency is
  completely unacceptable for realistic applications. To this end, we co-design a crypto-friendly
  editing algorithm and a customized 2PC protocol such that the editing algorithm is as effective
  as state-of-the-art model repair approaches while minimizing the protocol cost when executed
  using our optimized cryptographic protocol.

https://humanfeedback.io/

• We implemented our protocol and a baseline protocol using generic 2PC. We tested their performance for repairing Stable Diffusion v1.4. We observed **4 orders of magnitude improvement in runtime** compared to the baseline, bringing secure model repairing from merely a concept to something that can practically be deployed on modern models.

# 2 Notations and Preliminaries

**Notations.** We use lowercase bold letters like **c** to denote column vectors and uppercase bold letters like **W** to denote matrices. We write [n] to denote the set  $\{1, 2, \ldots, n\}$ . We use consistent notation for values in the diffusion model architecture, as defined in the next few paragraphs.

Diffusion Models [18, 38] are a class of generative models that have recently emerged as the SOTA in image generation. Inspired by non-equilibrium thermodynamics, diffusion models use a fixed algorithm to incrementally add random noise to images (or other data), and then learn how to reverse this process. The learned model is then used for image generation. Diffusion models have not always been the SOTA in image generation; prior to diffusion models, GANs were the most promising image generation models [11]. However, compared to GANs, diffusion models offer multiple advantages that lead to better results [12]. Diffusion models use more stable loss metrics than GANs. Additionally, because diffusion models generate images over a series of timesteps, their task is easier than that of GANs, which do it in one pass.

In this work, we focus on **text-to-image diffusion models** [32, 35, 26, 17, 34], where the diffusion process is guided by a user-provided text prompt that is embedded and injected into the cross-attention layers of the model. Formally, we consider a diffusion model  $\mathcal{M}$  that generates images by denoising a Gaussian sample  $\mathbf{x}_T$  over T time steps using a neural network  $D_{\theta}(\mathbf{x}_t,t,c)$ , where c is a conditioning signal derived from the text. The text prompt is first tokenized and processed by a text encoder, which outputs a sequence of token embeddings  $\{\mathbf{c}_i\}_{i=1}^{\ell}$  where  $\mathbf{c}_i \in \mathbb{R}^c$  that represent the semantic content of the input text. Let  $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_\ell] \in \mathbb{R}^{c \times \ell}$  denote the resulting matrix of text embeddings. At every cross-attention layer, these embeddings are linearly projected into  $\mathbf{K} = \mathbf{W}_K \mathbf{C} \in \mathbb{R}^{\ell \times k}$  and  $\mathbf{V} = \mathbf{W}_V \mathbf{C} \in \mathbb{R}^{\ell \times v}$  using learned key and v0 and v1 and v2 and v3 and v4 and v4 and v5 and v6 and v7 and v8 and v9 and v

**Diffusion Model Editing** aims to remove various biases from diffusion models and has become increasingly important as these models gain widespread adoption. One way it is done is by adjusting various aspects of the training process to limit bias; this can include altering the loss function [36] or debiasing the training dataset [22]. Fine-tuning existing diffusion models is perhaps a more realistic approach, as biases can become apparent after training. To do this, a small fraction of the weights in the diffusion model are updated to fix a specific problem. This can be done by editing the text encoder [43], or by directly editing the diffusion model [13, 29]. We focus on fine-tuning after training in this paper, as this ensures models can be updated as needed and do not need to be retrained.

**Oblivious Transfer** (OT) is a fundamental cryptographic primitive essential for secure computation protocols [31]. In a 1-out-of-n OT, a sender possesses n messages  $(m_1,\ldots,m_n)$ , and a receiver selects an index  $i\in[n]$  to retrieve  $m_i$  without revealing i to the sender. Simultaneously, the receiver gains no information about the other messages  $m_j$  for  $j\neq i$ . This ensures that the sender remains oblivious to the receiver's choice, and the receiver learns only the selected message.

**Secure Two-Party Computation** (2PC) [45, 15] enables two mutually distrustful parties, each holding private inputs, to jointly compute a public function without revealing any information beyond the output. We consider 2PC in the presence of static semi-honest adversaries, where parties follow the protocol but may attempt to learn additional information from the protocol execution transcript. The ideal functionalities of 1-out-of-*n* OT and 2PC are presented in Appendix B.

# 3 Problem Description

**Parties and Trust Assumptions.** We consider a setting including three parties as illustrated in Figure 1: a service provider S that offers diffusion-based image generation services, users U

who query the service and provides feedback when observing service failures, and a **model repair** institute  $\mathcal{I}$ , which specializes in repairing model failure and collaborates with  $\mathcal{S}$  to repair its model.

In this setting, we make the following trust assumptions in our threat model:

- Users  $\mathcal{U}$  query the image generation service with their textual prompt and receive images.  $\mathcal{U}$  discovers failures as they use  $\mathcal{M}$ -based services of  $\mathcal{S}$ .  $\mathcal{U}$ 's flagged failures are because of the fact that  $\mathcal{M}$  acquires knowledge within their training data [29] which become outdated, incorrect and harmful over time. For instance, for the prompt "A photo of a CEO", only 4% of generated images (with random seeds) contain female figures [29]. This feedback should only be visible to  $\mathcal{S}$ .
- A Service Provider S trains the text-to-image diffusion model  $\mathcal{M}$  on huge amounts of web-scraped image-caption pairs, and provides image generation services using  $\mathcal{M}$ . S wants to protect (i) the proprietary weights of  $\mathcal{M}$  and (ii) user-submitted feedback, which may contain sensitive user data and is subject to privacy regulations. We additionally require S must not reveal which failure it is fixing when interacting with  $\mathcal{I}$ , as it might inadvertently leak user data.
- A Model Repair Institute I specialized in repairing text-to-image diffusion models. I wants to keep both its repairing algorithm and fix database secret, as they are its core intellectual property.

Goal and Technical Challenges. The above-mentioned failures make the world knowledge of  $\mathcal{M}$  in deployments unaligned with users' values and expectations [14, 30, 3, 37, 8]. Our goal is to repair  $\mathcal{M}$  failures identified by  $\mathcal{U}$ . Although users are essential for flagging failures, they do not directly participate in the repair process. Once the feedback is submitted, it becomes the responsibility of  $\mathcal{S}$  to repair their model. As service providers usually lack expertise and resources for repairing failures (they mostly focus on enhancing image qualities),  $\mathcal{S}$  needs to contact an external Model Repair Institute  $\mathcal{I}$  to perform such fixes. This is challenging for several reasons. First, service providers are not allowed to share user's data² with third-parties due to privacy regulations. Second, service providers are not willing to hand over their models to third parties due to IP concerns. Third, Model Repair Institutes are not willing to disclose their fixes to service providers to protect their business model. Therefore, we model the protocol as a two-party computation between  $\mathcal{S}$  and  $\mathcal{I}$ , with the feedback treated as private input held by  $\mathcal{S}$ . We adopt the standard *semi-honest security* model, where both parties follow the repair protocol correctly but may try to infer additional information from the interaction: both  $\mathcal{S}$  and  $\mathcal{I}$  are institutions with legal and reputational reasons to behave correctly during model repairing, though they may have incentives to recover more information.

**Our Solution.** Given the trust assumptions above, our goal is to build a provably secure protocol that protects the private inputs of both parties: the model weights and user feedback held by  $\mathcal{S}$ , and the proprietary repair logic and database held by  $\mathcal{I}$ . To achieve this, we rely on cryptographic techniques. We design a crypto-friendly knowledge editing algorithm by adapting an efficient editing method that avoids retraining from scratch. Based on this, we construct a lightweight, customized two-party computation protocol, which we detail in the next section.

# 4 SURE: SecUre model REpairing

We propose SURE, a protocol for effective, efficient, and secure repair of text-to-image diffusion models based on user feedback and collaboration between a service provider  $\mathcal{S}$  and a model repair institute  $\mathcal{I}$ . SURE combines a crypto-friendly model repair algorithm with a customized two-party computation (2PC) protocol. Our approach builds on recent knowledge editing techniques [29] that enable model updates without full retraining. However, applying these techniques out-of-the-box is unsuitable for efficient 2PC due to the large number of layers in diffusion models and the high cost of interactive operations such as high-dimensional matrix multiplications and inverses. Our key insight is that most of this cost can be avoided by carefully modifying the editing algorithm.

We design a crypto-friendly repair algorithm (Section 4.1) tailored to the 2PC setting, without compromising the effectiveness of the original editing method. Our redesigned algorithm shifts almost all heavy computation offline, allowing each party to process its data locally and independently. Specifically: (1)  $\mathcal{I}$  constructs the repair database offline (Algorithm 1); and (2)  $\mathcal{S}$  applies the fix to its model parameters locally (Algorithm 2). In the online phase, we further develop a custom 2PC protocol (Section 4.2) that enables the service provider to securely locate and receive the fix

<sup>&</sup>lt;sup>2</sup>Note that we do not consider protection of confidentiality or privacy of users' request to S as it is only a failure identification and their institution knows their data, but we want to protect it against other institutions.

#### **Algorithm 1:** Repair database creation

**Input:** A Model Repair Institute  $\mathcal{I}$ , A public text encoder (*TextEncoder*), A collection of Failures Output: Repair Database

- 1: Repair Database = {}
- 2: **for all** failure  $\in$  Failures **do**
- Repair data pair = {source prompt:x, destination prompt:x'} ▷ Creating a repairing data
- $\{\mathbf{C} \in \mathbb{R}^{c \times l}, \mathbf{C}' \in \mathbb{R}^{c \times l'}\} = \mathit{TextEncoder}(\{\mathbf{x}, \mathbf{x}'\}) \quad riangle \; \mathsf{Tokenizing} \; \mathsf{and} \; \mathsf{computing} \; \mathsf{embeddings}$ 4:
- $\mathbf{C}^* \in \mathbb{R}^{c imes l} = \textit{RemoveAdditionalTokens}(\mathbf{C}') \; riangleright$  Creating an embedding that corresponds to 5: the same source token by discarding the embedding of additional tokens in the destination prompt
- $\begin{aligned} \mathbf{W}_{\text{fix}} &= \left( \lambda_{\text{failure}} \mathbf{I} + \mathbf{C}^* \mathbf{C}^\top \right) \left( \lambda_{\text{failure}} \mathbf{I} + \mathbf{C} \mathbf{C}^\top \right)^{-1} \\ \text{Repair Database.append}(\{\text{failure}: \mathbf{W}_{\text{fix}}\}) \end{aligned}$ 6: ▷ Creating Repair Knowledge
- 8: Output Repair Database

# **Algorithm 2:** Repair diffusion model parameters

**Input:** Service Provider S, A text-to-image diffusion model M, Received repair knowledge  $\mathbf{W}_{fix}$ **Output:** Updated parameters of the repaired text-to-image diffusion model  $\mathcal{M}$ 

- 1: CrossAttentionLayers  $\leftarrow$  CrossAttentionAccess( $\mathcal{M}$ ) ▷ Extract cross-attention layers that map textual data into visual data
- 2: for all  $i \in Size(CrossAttentionLayers)$  do
- $\mathbf{W}_V^{\prime i} \leftarrow \mathbf{W}_V^i \mathbf{W}_{\mathsf{fix}}$

▷ Update value projection matrix

 $\mathbf{W}_K^{\prime i} \leftarrow \mathbf{W}_K^i \mathbf{W}_{\mathsf{fix}}$ 

- ▷ Update key projection matrix
- 5: Updated diffusion model returned to only the service provider

corresponding to their failures from the institute's repair database through a secure fuzzy matching procedure and a lightweight Oblivious Transfer (OT) protocol. We prove (Section 4.3) that our protocol keeps users' feedback, service provider's model parameters, and the institute's proprietary editing algorithm confidential while ensuring that the model is faithfully repaired.

# 4.1 Crypto-Friendly Model Repair Algorithm

We instantiate SURE based on the Text-to-Image Model Editing (TIME) procedure introduced by Orgad et al. [29]. We briefly review their core editing algorithm before presenting our modifications.

The editing algorithm in TIME takes as input two prompts:

- A source prompt, e.g., "a photo of CEO" that under-specifies certain visual attributes. It allows the model to fill in missing details using its implicit assumptions, which could reflect bias.
- A more specific destination prompt, e.g., "a photo of **female** CEO" where an explicit attribute is added to correct the failure in the original source prompt.

The editing goal is to repair failures in the model's original output by shifting the image generation from reflecting the source prompt to better align with the intended visual attributes of the destination prompt. This enables targeted correction of outdated, incorrect, or socially biased associations embedded in the model. The key insight from Orgad et al. is that it suffices to update only the key and value projection matrices  $\mathbf{W}_K$  and  $\mathbf{W}_V$  (see Section 2 for detailed definitions) within the model's cross-attention layers. These matrices are responsible for mapping textual tokens into attention-compatible visual representations, and patching them effectively alters the generated output.

Let  $\{\mathbf{c}_i\}_{i=1}^\ell \subset \mathbb{R}^c$  and  $\{\mathbf{c}_i'\}_{i=1}^{\ell'} \subset \mathbb{R}^c$  be the token embeddings of the source and destination prompt. For every source token, TIME locates the corresponding destination token that contains the same word and denotes its embedding by  $\mathbf{c}_i^*$ . This gives the aligned set  $\{\mathbf{c}_i^*\}_{i=1}^\ell$  for tokens appear in both prompts. Let  $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_\ell]$  and  $\mathbf{C}^* = [\mathbf{c}_1^*, \dots, \mathbf{c}_\ell^*]$ , for every layer i, the closed-form update formula (Equation 5 in [29]) is given by

$$\mathbf{W}_{K}^{\prime i} = \left(\lambda \mathbf{W}_{K}^{i} + \mathbf{K}^{*} \mathbf{C}^{\top}\right) \left(\lambda \mathbf{I}_{d} + \mathbf{C} \mathbf{C}^{\top}\right)^{-1} \& \mathbf{W}_{V}^{\prime i} = \left(\lambda \mathbf{W}_{V}^{i} + \mathbf{V}^{*} \mathbf{C}^{\top}\right) \left(\lambda \mathbf{I}_{d} + \mathbf{C} \mathbf{C}^{\top}\right)^{-1}, \quad (1)$$
where  $\lambda \in \mathbb{R}^{+}$  is a hyperparameter, and  $\mathbf{K}^{*} = \mathbf{W}_{K} \mathbf{C}^{*}$  and  $\mathbf{V}^{*} = \mathbf{W}_{V} \mathbf{C}^{*}$ .

Efficiency and Privacy Challenges. The most direct way to securely evaluate the above update formula is to encode it as a circuit and run a generic 2PC: the provider S supplies the private weights

# Functionality $\mathcal{F}_{\mathsf{Repair}}$

This functionality is parameterized by a similarity metric  $d(\cdot, \cdot)$  and a database size n.

#### Input

- $\mathcal{S}$  inputs a query key  $\mathbf{k}_{qry} \in \mathbb{R}^k$  and model matrices  $\{\mathbf{W}_V^i \in \mathbb{R}^{v \times c}, \mathbf{W}_K^i \in \mathbb{R}^{k \times c}\}_{i \in [m]}$ .
- $\mathcal{I}$  inputs the database  $\{\mathbf{k}_i, \mathbf{C}_i, \mathbf{C}_i^*, \lambda_i\}_{i \in [n]}$  where  $\mathbf{C}_i, \mathbf{C}_i^* \in \mathbb{R}^{c \times l}, \lambda_i \in \mathbb{R}^+$ , and  $\mathbf{k}_i \in \mathbb{R}^k$ .

#### **Model Repair:**

- 1. Compute  $p = \arg\min_{i \in [n]} d(\mathbf{k}_{qry}, \mathbf{k}_i)$ , breaking ties by choosing the smallest i.
- 2. For each model layer  $i \in [m]$ , compute and send the following updated matrices and index p to S:

$$\mathbf{W}_{\star}^{\prime i} \leftarrow \left(\lambda_{p} \mathbf{W}_{\star}^{i} + \mathbf{W}_{\star}^{i} \mathbf{C}_{p}^{*} \mathbf{C}_{p}^{\top}\right) \left(\lambda_{p} \mathbf{I}_{c} + \mathbf{C}_{p} \mathbf{C}_{p}^{\top}\right)^{-1} \quad \star \in \{V, K\}.$$

Figure 2: Ideal functionality of model repair between S and I.

 $\mathbf{W}_{\star}^{i}$ , the institute  $\mathcal{I}$  supplies  $\mathbf{C}, \mathbf{C}^{*}, \lambda$ , and the circuit outputs the updated weights  $\mathbf{W}_{\star}^{\prime i}$  to  $\mathcal{S}$ . Despite significant advances in modern 2PC protocols [42, 19], applying them directly to this task remains inefficient. To illustrate the limitations of this generic approach, we implemented a baseline that computes the editing algorithm in a generic 2PC protocol, and it requires **over 100 hours** to perform a single repair (see Section 5). More importantly, because generic 2PC assumes the circuit is public, service provider will always learn the institute's proprietary repair algorithm.

Our Crypto-Friendly Editing Formula. The bottleneck above mainly comes from forcing heavy matrix operations into the secure computation. Our key observation is that we can refactor the editing formula in a way that completely eliminates any matrix operations inside 2PC. The matrix update formula in Equation 1 can be refactored as follow:

$$\mathbf{W}_{V}^{ii} = \left(\lambda \mathbf{W}_{V}^{i} + \mathbf{V}^{*} \mathbf{C}^{\top}\right) \left(\lambda \mathbf{I}_{d} + \mathbf{C} \mathbf{C}^{\top}\right)^{-1}$$

$$= \left(\lambda \mathbf{W}_{V}^{i} + \mathbf{W}_{V}^{i} \mathbf{C}^{*} \mathbf{C}^{\top}\right) \left(\lambda \mathbf{I} + \mathbf{C} \mathbf{C}^{\top}\right)^{-1}$$

$$= \underbrace{\mathbf{W}_{V}^{i}}_{\text{known to } \mathcal{S}} \underbrace{\left(\lambda \mathbf{I} + \mathbf{C}^{*} \mathbf{C}^{\top}\right) \left(\lambda \mathbf{I} + \mathbf{C} \mathbf{C}^{\top}\right)^{-1}}_{\mathbf{W}_{\text{fix.}}, \text{ known to } \mathcal{I}}.$$
(2)

Here, the service provider S holds  $W_V$ , and the institute I holds C,  $C^*$ , and the hyperparameter  $\lambda$ . Thus, I can compute

$$\mathbf{W}_{\mathsf{fix}} \leftarrow \left(\lambda \mathbf{I} + \mathbf{C}^* \mathbf{C}^\top\right) \left(\lambda \mathbf{I} + \mathbf{C} \mathbf{C}^\top\right)^{-1},\tag{3}$$

and S can update the matrix by computing  $\mathbf{W}_V^{\prime i} \leftarrow \mathbf{W}_V^i \mathbf{W}_{\text{fix}}$ . The above refactored equation applies identically to  $\mathbf{W}_K$  and holds for all layers in the model. The fix matrix  $\mathbf{W}_{\text{fix}}$  now encapsulates the semantics of the update and fully decouples model-specific parameters from repairs. Our refactored formula yields three immediate advantages for our purposes:

- One fix fits all. The same  $W_{\text{fix}}$  matrix can be reused across every cross-attention layer i, and applies uniformly to both  $W_K^i$  and  $W_V^i$ . This significantly simplifies the repair process and reduces communication.
- Matrix algebra disappears from 2PC. All matrix operations to compute  $W_{\text{fix}}$  are handled entirely by  $\mathcal{I}$  offline. Then,  $\mathcal{S}$  can use lighter cryptographic primitives like OT to acquire the  $W_{\text{fix}}$  from  $\mathcal{I}$ .
- Algorithm privacy is preserved. Because the fix is provided as a single matrix and applied independently by S, there is no need to reveal the full structure of the editing algorithm or encode it into a shared circuit. Therefore,  $\mathcal{I}$ 's proprietary repair method remains hidden from S.

We adopt this new editing formula in our protocol *SURE*. As we show in Section 5, this seemingly simple refactor achieve four orders of magnitude speed up over the baseline.

# 4.2 Efficient Two-Party Model Repair Protocol

We now provide a detailed description of the secure two-party model repair protocol in *SURE*. The ideal functionality and our two-party model repair protocol are presented in Figure 2 and Figure 3.

We first briefly recall our setting. The protocol SURE involves two parties: a service provider S and a model repair institute I. S wish to repair its deployed model I and derives from aggregated user

# Protocol $\Pi_{\mathsf{Repair}}$

#### **Input:**

- The service provider S and institute I agree on a similarity metric  $d(\cdot,\cdot)$  and the database size n.
- S inputs fix query vector  $\mathbf{k}_{qry} \in \mathbb{R}^k$  and model matrices  $\{\mathbf{W}_V^i \in \mathbb{R}^{v \times c}, \mathbf{W}_K^i \in \mathbb{R}^{k \times c}\}_{i \in [m]}$ , where m is the total number of model layers.
- $\mathcal{I}$  inputs  $\{\mathbf{k}_j, \mathbf{C}_j, \mathbf{C}_j^*, \lambda_j\}_{j \in [n]}$ , where  $\mathbf{C}_j, \mathbf{C}_j^* \in \mathbb{R}^{c \times l}, \lambda_j \in \mathbb{R}^+, \mathbf{k}_j \in \mathbb{R}^k$ , and n is the size of repair database.

**Database Initialization:**  $\mathcal{I}$  computes the repair database  $\{\mathbf{k}_i: \mathbf{W}_{\text{fix,j}}\}_{j \in [n]}$ , where

$$\mathbf{W}_{\mathsf{fix},j} \leftarrow \left(\lambda_{j}\mathbf{I}_{c} + \mathbf{C}_{j}^{*}\mathbf{C}_{j}^{\top}\right)\left(\lambda_{j}\mathbf{I}_{c} + \mathbf{C}_{j}\mathbf{C}_{j}^{\top}\right)^{-1}.$$

**Matching**: Let  $\mathcal{C}_{d,n}$  be the circuit that outputs $(p, \perp)$ , where  $p = \arg\min_{j \in [n]} d(\mathbf{k}_{qry}, \mathbf{k}_j)$ , breaking ties by choosing the smallest j.  $\mathcal{S}$  and  $\mathcal{I}$  send  $(\mathcal{C}_{d,n}, \mathbf{k}_{qry})$  and  $(\mathcal{C}_{d,n}, (\mathbf{k}_1, \dots, \mathbf{k}_n))$  to  $\mathcal{F}_{2PC}$ .  $\mathcal{S}$  receives the fix matrix index p.

#### **Model Repair:**

- 1. S and I send (recv, n, p) and (send,  $n, \{W_{fix,j}\}_{j \in [n]}$ ) to  $F_{OT}$ . S obtains the fix matrix  $W_{fix,p}$ .
- 2. For each model layer  $i \in [m]$  and  $\star \in \{K, V\}$ , S locally updates each layer of its model using the same fix matrix:  $\mathbf{W}_{\star}^{\prime i} \leftarrow \mathbf{W}_{\star}^{i} \mathbf{W}_{\text{fix,p}}$ .

Figure 3: Our secure model repair protocol in the  $(\mathcal{F}_{OT}, \mathcal{F}_{2PC})$ -hybrid model.

feedback a query key  $\mathbf{k}_{qry} \in \mathbb{R}^k$  that captures the failure domain to be fixed.  $\mathcal{I}$  maintains a private key–value repair database  $\{k_j: \mathbf{W}_{\text{fix,j}}\}_{j \in [n]}$  of size n, where each key  $\mathbf{k}_i \in \mathbb{R}^k$  semantically labels a failure and each  $\mathbf{W}_{\text{fix,j}}$  is the repair matrix for this failure. The protocol consists of three stages:

- 1. **Database Initialization.** Before interacting with S, the institute  $\mathcal{I}$  locally computes  $\mathbf{W}_{\text{fix,j}}$  from the embedding matrices  $\mathbf{C}_j$ ,  $\mathbf{C}_j^*$  and edit hyperparameter  $\lambda_j$ , and tag the fix with a key  $\mathbf{k}_j$  that semantically describe the failure.
- 2. **Matching.** S and T run a small circuit inside 2PC to locate the database entry whose key  $\mathbf{k}_p$  minimizes a public similarity metric  $d(\mathbf{k}_{qry}, \mathbf{k}_j)$ . After this stage, only S learns the index p; T learns nothing about  $\mathbf{k}_{qry}$  beyond the fact that a comparison occurred. When an exact match is sufficient—e.g., d is the discrete metric or the database indexes are public, T can determine T0 outright, so this stage can be skipped and the parties proceed directly to the next step.
- 3. **Oblivious Model Repair.** After acquiring the index p,  $\mathcal{I}$  runs an OT protocol to retrieve the single matrix  $\mathbf{W}_{\text{fix,p}}$  without revealing p and without accessing any other entry. It then updates every cross-attention layer locally by right-multiplying both value and key projections with  $\mathbf{W}_{\text{fix,p}}$  to complete the repair.

**Security Guarantees.** Our protocol ensures that (i) the institute  $\mathcal{I}$  learns nothing about the model  $\mathcal{M}$  or the query key  $\mathbf{k}_{qry}$ ; (ii) the service provider  $\mathcal{S}$  learns only the single fix matrix matching its query and gains no information about any other entry in  $\mathcal{I}$ 's database; and (iii) the editing algorithm itself remains private, because  $\mathcal{I}$  builds the database offline, the editing algorithm chosen by  $\mathcal{I}$  remains entirely hidden from  $\mathcal{S}$ . In the next section, we formalize and prove these guarantees.

#### 4.3 Security Proof

In this section, we establish the security of our protocol  $\Pi_{\text{Repair}}$  (Figure 3) and show how it can be generalized to any editing mechanism while hiding the editing algorithm being employed by the institute. All of our proofs are based on the standard composition paradigm [5]. We now state the following main security theorem of our protocol.

**Theorem 1** (Protocol Security). Protocol  $\Pi_{\text{Repair}}$  (Figure 3) securely realizes  $\mathcal{F}_{\text{Repair}}$  (Figure 2) in the  $(\mathcal{F}_{\text{OT}}, \mathcal{F}_{\text{2PC}})$ -hybrid model against semi-honest adversaries.

*Proof.* For clarity, we denote the service provider by  $P_1$  and the repair institute by  $P_2$  for the remainder of the proof.

**Correctness.** Note that all matrix products in both the protocol and the functionality are well-defined. Additionally, for all  $j \in [n]$ , the regularization parameter  $\lambda_j > 0$ , hence the matrix  $(\lambda_j \mathbf{I}_c + \mathbf{C}_j \mathbf{C}_j^\top) \succ 0$  and is therefore invertible.

To prove privacy, we separately consider the case of a corrupted institute and a corrupted service provider.

**Corrupted Institute**  $\hat{P}_2$ . It is straightforward to prove security against  $\hat{P}_2$ , as it receives no output from either  $\mathcal{F}_{OT}$  and  $\mathcal{F}_{2PC}$ . Therefore, a simulator  $\mathcal{S}_2$  that simply forwards  $\hat{P}_2$ 's message to  $\mathcal{F}_{Repair}$  can perfectly simulate its view.

**Corrupted Service Provider**  $\hat{P}_1$ . We construct a simulator  $S_1$  that calls  $\hat{P}_1$  as a subroutine and interacts with  $\mathcal{F}_{Repair}$  to simulate its view.  $S_1$  proceeds as follows:

- 1.  $S_1$  obtains the message  $(C_{d,n}, \mathbf{k}_{qry})$  from  $\hat{P}_1$  and record  $\mathbf{k}_{qry}$ .
- 2.  $S_1$  sends  $(\mathbf{k}_{qry}, \{\mathbf{I}_c, \mathbf{I}_c\})$  to  $\mathcal{F}_{Repair}$  and receives index p and  $\{\hat{\mathbf{W}}_V'^i, \hat{\mathbf{W}}_K'^i\}_{i \in [m]}$ .
- 3.  $S_1$  acts as  $\mathcal{F}_{2PC}$  and send p to  $\hat{P}_1$ ; upon obtaining (recv, n, p) from  $\hat{P}_1$ , send  $\hat{\mathbf{W}}_V^{\prime 1}$  to  $\hat{P}_1$ .

We show that  $\hat{P}_1$ 's view is perfectly simulated. To see this, notice that the ideal world, because  $S_1$  sends identity matrices to  $\mathcal{F}_{Repair}$ , for every layer

$$\hat{\mathbf{W}}_V^{\prime i} = (\lambda_p \mathbf{I}_c + \mathbf{I}_c \mathbf{C}_p^* \mathbf{C}_p^\top) (\lambda_p \mathbf{I}_c + \mathbf{C}_p \mathbf{C}_p^\top)^{-1} = \mathbf{W}_{\mathsf{fix,p}}.$$

As a result, the matrix  $\hat{P}_1$  received in the ideal execution is *exactly* the same from  $\mathcal{F}_{OT}$  in the real execution. Therefore, its view is perfectly simulated. As an honest  $P_2$  receives no output in both worlds, the joint output distributions are also identical in both worlds. This concludes the proof.

**Algorithm Privacy.** For concreteness, we instantiate our protocol based on the editing algorithm of [29]. However, our cryptographic construction readily accommodates *any* repair mechanism: any repair procedure that modifies model weights while leaving the network architecture unchanged can be dropped in without altering the protocol. Moreover, the protocol keeps the institute's choice of editing algorithm confidential. To see this, notice that  $\mathcal{S}$  only sees the resulting fix matrix  $\mathbf{W}_{\text{fix}}$  while the algorithm itself remains hidden. To formalize this property, we first define the notion of editing algorithms and prove a theorem stating the algorithm-hiding property of our protocol.

**Definition 1** (Editing Algorithms). A model repair editing algorithm is an efficient mapping

$$f: (\mathbf{C}, \mathbf{C}^*, \mathsf{aux}) \to \mathbf{W}_{\mathsf{fix}},$$

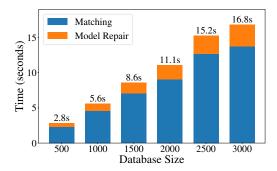
where  $C, C^*$  are the source and target prompt embedding matrices, aux is institute-held auxiliary input, and  $W_{\text{fix}}$  is the fix matrix of proper dimensions that is right-multiplied to every model layer.

**Theorem 2** (Algorithm Privacy). Let  $U_f = \{f_1, \ldots, f_Z\}$  be any finite family of editing algorithms. Let  $\Pi'_{\mathsf{Repair}}$  be the extension of  $\Pi_{\mathsf{Repair}}$  in which  $\mathcal{I}$  chooses an index  $z \in [Z]$ , and builds its database using  $f_i$ . Let  $\mathcal{F}'_{\mathsf{Repair}}$  be the corresponding ideal functionality that receives the description of  $f_z$  from  $\mathcal{I}$ , evaluates  $f_z$  internally to obtain the fix matrices, and sends those matrices to  $\mathcal{S}$ . Then, for every PPT adversary  $\mathcal{A}_{\mathcal{S}}$  corrupting the service provider  $\mathcal{S}$ , there exists a PPT simulator  $\mathcal{S}'_1$  such that  $\mathsf{view}_{\mathcal{A}_{\mathcal{S}}}^{\mathsf{Real}}(\Pi'_{\mathsf{Repair}}) \equiv \mathsf{view}_{\mathcal{S}'_1}^{\mathsf{Ideal}}(\mathcal{F}'_{\mathsf{Repair}})$ .

*Proof.* The proof follows from the security proof of Theorem 1 in a straightforward manner. We define a modified simulator  $\mathcal{S}_1'$  that behaves the same as  $\mathcal{S}_1$  in the security proof, except it forward  $\mathcal{S}$ 's message to  $\mathcal{F}_{\mathsf{Repair}}'$  instead of  $\mathcal{F}_{\mathsf{Repair}}$ . It then follows from the security proof that  $\Pi'_{\mathsf{Repair}}$  securely realizes  $\mathcal{F}'_{\mathsf{Repair}}$  in the  $(\mathcal{F}_{\mathsf{OT}}, \mathcal{F}_{\mathsf{2PC}})$ -hybrid model and  $S'_1$  perfectly simulates  $\mathcal{A}_{\mathcal{S}}$ 's view.

Theorem 2 implies that the service provider S learns no information about the institute's chosen editing algorithm  $f_i$  beyond what is already implied by the fix matrix  $\mathbf{W}_{\text{fix,p}}$ . Consequently, our protocol enables the institute to swap in or fine-tune its proprietary repair procedures without touching the underlying cryptographic protocol. This design not only supports fast iteration but also hides the institute's editing knowledge from the service provider.

By contrast, generic secure 2PC protocols would require compiling the entire editing algorithm into a single Boolean or arithmetic circuit that is known to both parties—a standard assumption in secure computation constructions [45, 15]. In such settings, the circuit's structure (and hence the algorithm it encodes) is public, even if the inputs remain hidden.



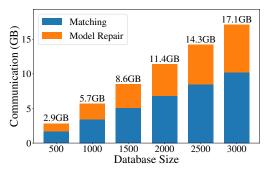


Figure 4: Runtime and communication cost for SURE to repair Stable Diffusion v1.4 with varying repair database sizes. Reported times are averages over 10 runs. Each database entry consists of a [768,768] fix matrix and a [100,1] key, where all numbers are single-precision floating-point. Communication reported is the larger of the two parties' data sent. Costs are decomposed into two stages: (1) Matching – finding the closest key in the database to the failure query using the Euclidean distance as the similarity metric and (2) Model Repair – returning the fix via the oblivious transfer protocol to  $\mathcal{S}$ .

Table 1: Runtime and communication costs of *SURE* vs. the baseline approach to perform one repair of Stable Diffusion v1.4 with varying database sizes. *SURE* is orders of magnitude faster than the baseline, as our protocol completely avoids matrix operations in secure computation.

Database Size	Baseline		Ours		
	Running Time (hours)	Comm. (TB)	Running Time (seconds)	Comm. (GB)	Running Time Improvement
500	167.36	76.42	2.81	2.94	$2.14{ imes}10^5$
1000	171.84	82.30	5.58	5.65	$1.11{\times}10^5$
1500	176.02	88.19	8.62	8.61	$7.35{\times}10^4$
2000	180.48	94.08	11.10	11.41	$5.85{ imes}10^4$
2500	184.96	99.97	15.21	14.32	$4.38{\times}10^4$
3000	189.44	105.86	16.77	17.09	$4.07{\times}10^4$

# 5 Experimental Evaluation

We implement and evaluate the efficiency of *SURE* in repairing Stable Diffusion v1.4 [33] with 32 layers, and compare it with a baseline model repair protocol that runs entirely within a generic 2PC framework for comparison. Our code is provided in the supplementary materials. Details about our implementations are in Appendix C.

**Efficiency**. Figure 4 shows the end-to-end runtime and communication cost of *SURE* when performing a full Stable Diffusion v1.4 repair across varying repair database sizes. *SURE* completes the repair in under 17 seconds, even with a repair database of 3,000 entries, with communication capped at 17.1 GB. Therefore, *SURE* is highly efficient.

**Benchmarking.** We further break down the total cost into two main stages—key matching and model repair—as described in our protocol in Figure 3. **Most of the runtime is spent in the matching stage**, which uses a lightweight 2PC protocol to identify the nearest key. In contrast, the OT-based model repair phase is highly efficient and remains nearly constant regardless of database sizes.

**Scalability.** Our protocol scales well with both the repair database and model size: since only a single fix matrix is retrieved and applied across all layers, **the online runtime is independent of the number of model layers**. Moreover, *SURE*'s modular design allows for further optimization: the matching step can be replaced with more efficient cryptographic primitives such as private information retrieval or fuzzy private set intersection. In cases where the database key is public<sup>3</sup>, the matching phase can be skipped entirely to further reduce overhead.

**Comparison.** Table 1 compares our protocol against the baseline. **Our protocol achieves up to a**  $2 \times 10^5$  **speedup**. This dramatic improvement stems from avoiding expensive matrix operations and linear scans within 2PC. In the baseline, most of the cost arises from executing the entire editing

<sup>&</sup>lt;sup>3</sup>For example, [29] considers a database of gender bias in different professions, where the database key is simply the profession name (e.g., "nurse"), which is made public. Because the index of the desired fix is known in advance, matching is unnecessary.

formula securely and retrieving the correct fix matrix through a full scan of the database, both of which scale poorly with the database size. In contrast, our customized design isolates the secure computation to a small matching task and a lightweight OT-based retrieval protocol, while offloading all matrix operations to local (offline) computation, resulting in far superior performance.

#### References

- [1] Adobe Inc. Adobe firefly. https://www.adobe.com/products/firefly.html. Accessed 2025-05-14.
- [2] Artbreeder, Inc. Artbreeder. https://aidailydrop.com/artbreeder-ai-art-generator/. Founded by Joel Simon, Accessed 2025-05-14.
- [3] Samyadeep Basu, Nanxuan Zhao, Vlad I Morariu, Soheil Feizi, and Varun Manjunatha. Localizing and editing knowledge in text-to-image generative models. In *The Twelfth International Conference on Learning Representations*, 2023.
- [4] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023.
- [5] Ran Canetti. Security and composition of multiparty cryptographic protocols. *Journal of Cryptology*, 13(1):143–202, January 2000.
- [6] Canva Pty Ltd. Canva magic media (dream lab). https://en.wikipedia.org/wiki/Canva. Accessed 2025-05-14.
- [7] Pu Cao, Feng Zhou, Qing Song, and Lu Yang. Controllable generation with text-to-image diffusion models: A survey. *arXiv preprint arXiv:2403.04279*, 2024.
- [8] Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270, 2023.
- [9] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- [10] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [11] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A. Bharath. Generative adversarial networks: An overview. *IEEE Signal Processing Magazine*, 35(1):53–65, 2018.
- [12] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 8780–8794, 2021.
- [13] Rohit Gandikota, Joanna Materzyńska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from diffusion models. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pages 2426–2436, 2023.
- [14] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5111–5120, 2024.
- [15] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or A completeness theorem for protocols with honest majority. In Alfred Aho, editor, 19th ACM STOC, pages 218–229. ACM Press, May 1987.

- [16] Google LLC. Google imagefx. https://blog.google/technology/ai/google-labs-imagefx-textfx-generative-ai/. Via Google Labs/Google DeepMind, Accessed 2025-05-14.
- [17] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis, 2022.
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc.
- [19] Zhicong Huang, Cheng Hong, Chenkai Weng, Wen-jie Lu, and Hunter Qu. More efficient secure matrix multiplication for unbalanced recommender systems. *IEEE Transactions on Dependable and Secure Computing*, 20(1):551–562, 2023.
- [20] Ideogram, Inc. Ideogram 3.0. https://en.wikipedia.org/wiki/Ideogram\_(text-to-image\_model). Accessed 2025-05-14.
- [21] Marcel Keller. MP-SPDZ: A versatile framework for multi-party computation. In *Proceedings* of the 2020 ACM SIGSAC Conference on Computer and Communications Security, 2020.
- [22] Yeongmin Kim, Byeonghu Na, Minsang Park, JoonHo Jang, Dongjun Kim, Wanmo Kang, and Il-Chul Moon. Training unbiased diffusion models from biased dataset. *CoRR*, abs/2403.01189, 2024.
- [23] Leonardo Interactive Pty Ltd. Leonardo ai. https://abr.business.gov.au/ABN/View?id=56662209485. Trading as leonardo.ai, Accessed 2025-05-14.
- [24] Microsoft Corp. Bing image creator and microsoft designer. https://www.bing.com/create. Powered by OpenAI, Accessed 2025-05-14.
- [25] Midjourney, Inc. Midjourney. https://en.wikipedia.org/wiki/Midjourney. Accessed 2025-05-14.
- [26] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models, 2022.
- [27] NightCafe Studio Pty Ltd. Nightcafe studio. https://nightcafe.studio/pages/about-nightcafe. Accessed 2025-05-14.
- [28] OpenAI, Inc. Dalle 3. https://openai.com/index/dall-e-3/. Accessed 2025-05-14.
- [29] Hadas Orgad, Bahjat Kawar, and Yonatan Belinkov. Editing implicit assumptions in text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7053–7061, 2023.
- [30] Yong-Hyun Park, Sangdoo Yun, Jin-Hwa Kim, Junho Kim, Geonhui Jang, Yonghyun Jeong, Junghyo Jo, and Gayoung Lee. Direct unlearning optimization for robust and safe text-to-image models. *arXiv preprint arXiv:2407.21035*, 2024.
- [31] Michael O. Rabin. How to exchange secrets with oblivious transfer. 2005. https://eprint.iacr.org/2005/187.
- [32] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022.
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [34] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation, 2023.

- [35] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Lit, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Raphael Gontijo-Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [36] Xudong Shen, Chao Du, Tianyu Pang, Min Lin, Yongkang Wong, and Mohan S. Kankanhalli. Finetuning text-to-image diffusion models for fairness. In *The Twelfth International Conference* on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024.
- [37] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom 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*, pages 6048–6058, 2023.
- [38] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [39] Stability AI Ltd. Stable diffusion / dreamstudio. https://dreamstudio.stability.ai/. Accessed 2025-05-14.
- [40] European Union. Regulation (eu) 2016/679 (general data protection regulation). https://eur-lex.europa.eu/eli/reg/2016/679/oj, 2016.
- [41] Xiao Wang, Alex J. Malozemoff, and Jonathan Katz. EMP-toolkit: Efficient multiparty computation toolkit. https://github.com/emp-toolkit, 2016.
- [42] Xiao Wang, Samuel Ranellucci, and Jonathan Katz. Authenticated garbling and efficient maliciously secure two-party computation. In Bhavani M. Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *ACM CCS 2017*, pages 21–37. ACM Press, October / November 2017.
- [43] Tianwei Xiong, Yue Wu, Enze Xie, Yue Wu, Zhenguo Li, and Xihui Liu. Editing massive concepts in text-to-image diffusion models, 2024.
- [44] Pengfei Yang, Ngai-Man Cheung, and Xinda Ma. Text to image generation and editing: A survey. *arXiv preprint arXiv:2505.02527*, 2025.
- [45] Andrew Chi-Chih Yao. How to generate and exchange secrets. In 27th Annual Symposium on Foundations of Computer Science (sfcs 1986), pages 162–167, 1986.

#### Functionality $\mathcal{F}_{OT}$

Upon receiving (send, n,  $\{m_i\}_{i\in[n]}$ ) from  $P_1$  and (recv, n, b) from the  $P_2$  where  $b\in[n]$ , send  $m_b$  to  $P_2$ .

Figure 5: Ideal functionality of 1-out-of-*n* oblivious transfer.

#### Functionality $\mathcal{F}_{2PC}$

For  $i \in \{1, 2\}$ , upon receiving  $(\mathcal{C}, x_i)$  from  $P_i$ , compute  $y_1, y_2 \leftarrow \mathcal{C}(x_1, x_2)$  and send  $y_i$  to  $P_i$ .

Figure 6: Ideal functionality of secure two-party computation.

# **Technical Appendices**

#### A Limitations

SURE is the first secure framework for model editing and demonstrates high efficiency even on commercial-scale text-to-image diffusion models. However, it is currently limited to this domain. Extending to other models, such as LLMs, may require new editing algorithms and cryptographic techniques. Also, our efficient protocol only supports linear editing of the model weight; if editing techniques involve architectural changes or non-linear weight updates, further optimizations might need to be made to maintain its efficiency. While SURE can handle multiple repairs, no optimization, batching, or amortization is implemented. We comment that the matching phase can be further optimized for multiple repair settings using techniques like fuzzy private set intersection or private information retrieval.

#### **B** Ideal Functionalities

The ideal functionality of 1-out-of-n OT is depicted in Figure 5. This can be efficiently realized using  $\log n$  1-out-of-2 OT, which can in turn be efficiently computed using existing cryptographic protocols. The ideal functionality of 2PC is presented in Figure 6.

# C Implementations Details

**Experiment Setup.** We implement our end-to-end protocol *SURE* and compare it with a baseline protocol that executes all editing operations within a generic 2PC framework. The baseline is implemented based on the semi protocol variant from the MP-SPDZ framework (BSD3 License) [21], a popular framework for benchmarking generic secure computation protocols. However, MP-SPDZ does not provide a low-level OT interface suitable for our customized protocol in *SURE*. For ease of integration, we instead implemented *SURE* using the EMP-OT library from the EMP-toolkit (MIT License) [41], which provides efficient implementations of various OT primitives and a flexible low-level API.

To evaluate both SURE and the baseline, we perform a single model repair on Stable Diffusion v1.4 [33]. In this model, the source and target prompt embeddings  $\mathbf{C}, \mathbf{C}^*$  are matrices of shape [768, 77]. The repair algorithm modifies 32 cross-attention layers in total, where each layer contains key and value projection matrices of shape [320, 768]. As a result, each fix matrix  $\mathbf{W}_{\text{fix}}$  has dimension [768, 768]. We set the query and database keys  $\mathbf{k}_{\text{qry}}, \{\mathbf{k}_i\}_{i\in[n]}$  to 100-dimensional vectors. We represent all values using single-precision floating-point numbers and use Euclidean distance as the similarity metric for key matching. All experiments are run with a single thread on two Amazon EC-2 c7i.2xlarge instances, each with 16 GB of RAM.

**Baseline**. To highlight the efficiency of our lightweight protocol *SURE*, we implement a baseline model repair protocol that runs entirely within a generic 2PC framework for comparison. To ensure a fair comparison, we apply several optimizations to avoid penalizing the baseline unnecessarily. First, we represent editing computation as an arithmetic circuit, which is more efficient than Boolean circuits for linear algebra. Additionally, we use 32-bit fixed-point representation to avoid the high cost of floating-point arithmetic in 2PC. To reduce overhead further, we allow the model repair institute

of floating-point arithmetic in 2PC. To reduce overhead further, we allow the model repair institute  $\mathcal{I}$  to pre-compute the inverse matrices  $\mathbf{W}_{\mathsf{inv}}^p = \left(\lambda_p \mathbf{I}_c + \mathbf{C}_p \mathbf{C}_p^{\mathsf{T}}\right)^{-1}$  outside the 2PC to avoid costly secure matrix inversion.

The baseline protocol proceeds as follows:

- 1.  $\mathcal S$  inputs the query key  $\mathbf k_{\mathsf{qry}}$  and all projection matrices  $\{\mathbf W_K^i, \mathbf W_V^i\}_{i \in [m]}$  in all m layers;  $\mathcal I$  inputs the repair database  $\{\mathbf k_i, \lambda_i, \mathbf C_i, \mathbf C_i^*, \mathbf W_{\mathsf{inv}}^i\}_{i \in [n]}$ .
- 2. The circuit matches the closest index  $p = \arg\min_{i \in [n]} d(\mathbf{k}_{qry}, \mathbf{k}_i)$  and breaks ties by choosing the smallest i. Then, for each layer  $i \in [m]$ , it computes the fix for all matrices:

$$\mathbf{W}_{\star}^{\prime\,i} \;\leftarrow\; \Big(\lambda_{p}\mathbf{W}_{\star}^{i} + \mathbf{W}_{\star}^{i}\mathbf{C}_{p}^{*}\mathbf{C}_{p}^{\top}\Big)\mathbf{W}_{\mathrm{inv}}^{p} \quad \star \in \{K,V\}.$$

3. The updated matrices  $\{\mathbf{W}_K'^i,\mathbf{W}_V'^i\}_{i\in[m]}$  are then revealed to  $\mathcal{S}.$ 

Despite these optimizations, the baseline remains orders of magnitude slower than our protocol *SURE*, as we show in Section 5.