

# LATENT-INTEGRATED FINGERPRINT DIFFUSION: A DUAL-PATH FRAMEWORK FOR ROBUST ATTRIBUTION IN TEXT-TO-IMAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent advances in text-to-image diffusion models have enabled the generation of hyper-realistic images directly from textual descriptions while simultaneously raising pressing concerns regarding misinformation and the potential misuse of synthetic media. Although traditional fingerprinting schemes provide a rudimentary means for accountability, these methods often require substantial compromises in image quality or are vulnerable to adversarial post-processing, thereby limiting their ability to reliably attribute generated images to individual users. In response, we propose Latent-Integrated Fingerprint Diffusion (LIFD), a novel dual-path fingerprinting framework that significantly extends prior approaches such as WOUAF by integrating an additional latent-space conditioning mechanism inspired by the cross-attention operations found in StableVITON. Our approach embeds a distinctive digital signature into every generated image via a two-pronged strategy. On one hand, a parameter-level modulation mechanism injects a user-specific binary fingerprint into the decoder weights of a pre-trained diffusion model through an affine transformation; on the other hand, an attention-based latent conditioning channel subtly introduces a spatial fingerprint into the intermediate feature maps during the denoising process. This dual-channel design effectively mitigates the inherent trade-off between increasing fingerprint dimensionality and maintaining high attribution accuracy, while simultaneously enhancing robustness against a wide range of image manipulations such as JPEG compression, Gaussian blurring, and adversarial noise attacks. Our key contributions are as follows:

- **Dual-Channel Fingerprinting:** We introduce a two-staged embedding strategy that concurrently modulates the model weights and injects latent fingerprints. By splitting the fingerprint signal into two orthogonal channels, our approach is significantly more resistant to removal or tampering when compared to conventional single-channel methods.
- **Attention-Based Latent Injection:** We incorporate a custom cross-attention block within the U-Net backbone of the diffusion model in order to “paint in” a barely perceptible, yet machine-detectable, fingerprint into the latent representations. An auxiliary total variation loss is applied to guarantee that the injected fingerprint remains spatially sharp and well localized even after smoothing operations.
- **Adaptive Balancing Mechanism:** We propose a dynamic balancing module that automatically adjusts the relative contributions of the parameter-level and latent-space fingerprint channels in response to variations in fingerprint dimensions and diverse post-processing conditions, thereby preserving high-fidelity image synthesis without compromising attribution accuracy.
- **Robust Two-Stream Fingerprint Extraction:** We design a ResNet-inspired extraction network that is jointly trained with a fidelity regularization term. This network is capable of reliably recovering the dual-channel fingerprint even in scenarios where one of the fingerprint channels has been partially compromised by subsequent image manipulations.

To train LIFD, we fine-tune a pre-trained text-to-image diffusion model within a dual supervisory framework. Our joint loss function is comprised of a binary cross-entropy term for fingerprint recovery, image quality metrics such as the CLIP-score

and the Fréchet Inception Distance (FID), and additional penalties that enforce robustness against simulated adversarial attacks. The training process encourages the model to imprint a unique digital signature by simultaneously optimizing for high visual fidelity and strong, resilient attribution signals. At inference time, the modulated model weights and latent conditioning module jointly imprint the digital signature onto each generated image, and the two-stream extraction network decodes the fingerprint to unambiguously attribute the image to its originating user, thus establishing a clear and verifiable pathway for accountability in the era of synthetic media generation. Extensive experiments on benchmark datasets including MS-COCO (using the Karpathy split) and LAION-Aesthetics demonstrate that LIFD achieves near-perfect attribution accuracy with minimal adverse impact on image quality. In our evaluations, even when images are subjected to aggressive post-processing such as high-ratio JPEG compression, Gaussian blurring, or additive adversarial noise, our dual-channel design is able to sustain high fingerprint recovery accuracy; indeed, our experimental results indicate that LIFD outperforms traditional single-channel methods by an average margin of approximately 11

## 1 INTRODUCTION

The rapid proliferation of generative models—particularly text-to-image diffusion frameworks—has revolutionized visual content creation while simultaneously introducing critical challenges in accountability and security. A major issue is the need to reliably associate synthesized content with its source so that responsibility may be assigned in cases of misuse. Early fingerprinting methods, such as the technique presented in ?, incorporate user-specific digital fingerprints directly into pre-trained models via weight modulation. In these techniques, selected decoder weights (in models such as Stable Diffusion) are perturbed according to each user’s unique binary code, achieving near-perfect attribution accuracy while largely preserving image quality. However, relying solely on weight modulation exposes inherent limitations: an unfavorable trade-off between the dimensionality of the fingerprint and the precision of attribution, as well as vulnerability to sophisticated adversarial attacks and common post-processing operations.

### 1.1 MOTIVATION AND CHALLENGES

Secure and accountable distribution of generative models involves multiple, often conflicting, challenges. On one hand, state-of-the-art image synthesis requires extremely high fidelity and any fingerprint embedding process must not significantly perturb the generative procedure. On the other hand, ensuring robustness against a range of adversarial post-processing operations—such as Gaussian noise addition, blurring, and JPEG compression—demands an injection system capable of withstanding diverse perturbations. These conflicting objectives motivate the design of a fingerprinting strategy that distributes the digital signature across multiple channels. In doing so, even if one channel is compromised, the overall fingerprint can still be recovered.

### 1.2 PROPOSED APPROACH: LATENT-INTEGRATED FINGERPRINT DIFFUSION (LIFD)

In this work, we propose a novel framework, **Latent-Integrated Fingerprint Diffusion (LIFD)**, that overcomes the shortcomings of conventional fingerprinting techniques by integrating two complementary injection channels. Unlike traditional methods that rely solely on parameter-level fingerprinting, LIFD introduces an additional latent-space injection mechanism. Inspired by the zero cross-attention blocks employed in StableVITON, our approach uses a custom cross-attention mechanism to embed a subtle, spatially distributed fingerprint into the latent representations during the denoising process.

Specifically, the parameter-level branch perturbs selected decoder weights via an affine transformation driven by a user-specific binary code, as proposed in ?. In parallel, the latent-space branch incorporates a custom cross-attention block within the U-Net backbone to inject fine-grained spatial information directly into the latent feature maps. An auxiliary loss, for example an attention total variation loss, is employed to ensure that the injected fingerprint remains sharp and spatially localized. As a result, even if adversarial manipulations or common post-processing operations degrade one channel, the overall digital signature can still be recovered via the complementary channel.

### 1.3 KEY INNOVATIONS

Our LIFD framework introduces several key innovations that jointly address the challenges of maintaining high synthesis fidelity while ensuring robust fingerprint embedding. The primary contributions of this work are:

- **Dual Channel Injection: Parameter-Level Modulation** is combined with **Latent-Space Conditioning** in a unified framework. While the former perturbs decoder weights using an affine transformation based on user-specific binary codes, the latter employs a custom cross-attention mechanism within the U-Net backbone to inject a subtle spatial fingerprint into the latent representations.
- **Adaptive Balancing Mechanism:** A dynamic balancing module is introduced to reconcile the trade-off between image quality and fingerprint robustness. By adjusting a hyperparameter,  $\alpha$ , the model flexibly modulates the relative contributions of the two injection channels, enabling fine-tuning to achieve an optimal balance as measured by metrics such as the Fréchet Inception Distance (FID) and CLIP-score.
- **Robust Fingerprint Extraction:** A dual-path extraction network, inspired by ResNet architectures, is developed to jointly decode the fingerprint signals from both the modulated weights and the latent features. This integrated extraction design enhances recovery accuracy such that, even under degraded or adversarial conditions, the digital signature is reliably retrieved.

### 1.4 TRAINING AND INFERENCE METHODOLOGY

To train the proposed LIFD framework, we fine-tune a base pre-trained text-to-image diffusion model (e.g., Stable Diffusion) using a dual-supervision strategy. The parameter-level branch leverages the mapping network from  $\mathbf{z}$  along with a binary cross-entropy loss and additional image quality regularizers to preserve synthesis fidelity. Concurrently, the latent-conditioning branch is optimized via modified zero cross-attention blocks, with an auxiliary attention total variation loss enforcing spatial sharpness and locality of the injected fingerprint. During inference, both branches operate jointly to imprint the digital signature onto the generated images, and a combined extraction network decodes the dual-path fingerprint to ensure high-confidence attribution.

### 1.5 EXPERIMENTAL EVALUATION PROTOCOL

Our experimental evaluation is designed to rigorously validate the advantages of LIFD through several targeted studies, each of which is fully implementable in Python using libraries such as PyTorch and torchvision. The main evaluation protocols include:

- **Dual-Channel Fingerprint Robustness Experiment:** We compare three modes of fingerprint injection: (a) parameter-only, (b) latent-only, and (c) the combined dual-channel approach of LIFD. Generated images are subjected to adversarial perturbations—such as Gaussian noise, blurring, and JPEG compression—and fingerprint extraction accuracy is measured to assess robustness.
- **Ablation Study on Adaptive Balancing:** By systematically varying the adaptive balancing parameter  $\alpha$  (e.g.,  $\alpha = 0, 0.25, 0.5, 0.75, 1.0$ ), we evaluate the trade-off between image quality (assessed via FID and CLIP-score) and fingerprint extraction accuracy, thereby identifying the optimal balance point.
- **Latent Fingerprint Injection Analysis:** Forward hooks in the custom cross-attention module are used to capture and visualize intermediate attention maps. Quantitative measures, such as total variation computed over these maps, confirm that the latent fingerprint remains spatially localized and robust against common post-processing operations.

In summary, by directly addressing the dual challenges of high-fidelity image generation and secure fingerprint embedding, our LIFD framework constitutes a significant advancement toward accountable generative modeling. The remainder of the paper is organized as follows. The Methods section details the technical aspects of our dual-path fingerprint injection strategy, including the mathematical

formulations underlying weight modulation and cross-attention-based latent conditioning. The Experiments section describes our evaluation protocols, complete with pseudocode and configuration details. In the Results section, we present both quantitative and qualitative analyses demonstrating the superiority of LIFD over existing methods. Finally, we discuss limitations and potential extensions, including applications to additional modalities such as text, audio, and video.

## 2 RELATED WORK

In this section, we review existing fingerprinting and attribution techniques for generative models with an emphasis on text-to-image diffusion systems, and we contrast these methods with the innovations introduced in our proposed Latent-Integrated Fingerprint Diffusion (LIFD).

### 2.1 EXISTING FINGERPRINTING AND ATTRIBUTION METHODS

Recent work such as ? introduces weight modulation techniques that embed user-specific fingerprints into pre-trained diffusion models. In these approaches, a unique digital identifier is imprinted into the model parameters via an affine transformation applied to a user-specific binary code. Although these methods achieve near-perfect attribution accuracy and demonstrate robustness under various post-processing operations, they suffer from a trade-off between the fingerprint dimension and extraction accuracy. Moreover, relying solely on parameter-level fingerprinting renders these methods vulnerable to sophisticated adversarial attacks—for example, attacks implemented via auto-encoder purification.

At the same time, conditioning techniques based on cross-attention mechanisms, as deployed in models like StableVITON, have been applied to control the generative process and preserve fine-grained stylistic details. While attention-based conditioning improves the fidelity of generated images, it has not been primarily adapted for robust fingerprinting and model attribution tasks.

### 2.2 ADVANCES IN DUAL-CHANNEL FINGERPRINTING

The proposed LIFD method leverages the complementary strengths of both weight modulation and attention-based latent conditioning. In our framework the digital fingerprint is distributed across two synergistic channels:

- **Parameter-Level Modulation:** Following the approach of ?, a user-specific identifier is embedded directly into the model weights via an affine transformation of a binary code. This technique ensures reliable fingerprint extraction—even after conventional post-processing—while preserving the model’s architecture.
- **Latent-Space Conditioning:** Inspired by the zero cross-attention mechanism used in StableVITON, this branch injects a subtle spatial code into the latent feature maps during the denoising process. A custom cross-attention block guides the generation to “paint in” a barely perceptible yet machine-detectable fingerprint. An auxiliary attention total variation loss enforces spatial localization and enhances resilience against smoothing or blurring attacks.

An adaptive balancing module fuses the two channels, mitigating the trade-off between fingerprint dimension and extraction accuracy. By allowing the latent channel to boost the fingerprint signal during reconstruction, the LIFD method is better equipped to resist aggressive post-processing and adversarial manipulations.

### 2.3 EXPERIMENTAL COMPARISONS AND INSIGHTS

Our experimental investigations examine several aspects of the dual-channel design. In one series of experiments, images are generated under three distinct modes: parameter-only injection, latent-only injection, and the combined dual-channel approach. Under simulated adversarial perturbations—including Gaussian noise addition, blurring, and JPEG compression—the dual-channel method consistently demonstrates improved robustness compared to single-channel alternatives.

An ablation study analyzes the impact of the adaptive balancing parameter  $\alpha$ . By sweeping  $\alpha$  over a range from 0 to 1, we observe a trade-off between image quality (assessed using metrics such as FID and CLIP-Score) and fingerprint extraction accuracy. This systematic analysis validates the effectiveness of the dynamic balancing module and offers guidance for optimal parameter selection.

A latent fingerprint injection analysis, which includes attention map visualization via heatmaps and sensitivity analysis, quantitatively evaluates the spatial localization of the injected fingerprint. Measurements of total variation confirm that the attention-based latent injection yields a sharp, localized fingerprint that remains resilient even after simulated smoothing operations.

## 2.4 SUMMARY OF CONTRIBUTIONS

The insights from prior research and our experiments motivate the following key contributions:

- **Dual Fingerprinting Channels:** Integration of parameter-level modulation with latent-space conditioning to enhance robustness against post-processing attacks.
- **Attention-Based Latent Injection:** Utilization of cross-attention mechanisms to embed a spatially localized and sharp fingerprint within the latent representations.
- **Adaptive Balancing Mechanism:** Introduction of a dynamic module to adjust the relative contributions of the two channels, thereby maintaining high image quality while ensuring robust fingerprint extraction.
- **Robust Fingerprint Extraction:** Deployment of a dual-stream decoder network that jointly processes signals from both channels to achieve improved detection accuracy even when one channel is compromised.

In summary, the LIFD framework presents a significant advancement over existing single-channel methods by offering a resilient and adaptable approach for model attribution in text-to-image diffusion systems.

## 3 BACKGROUND

The following section lays out the theoretical foundations, prior work, and design decisions that motivate our approach. We review diffusion-based image generation and digital fingerprinting in generative models, outline the limitations of existing techniques, and introduce our dual-channel Latent-Integrated Fingerprint Diffusion (LIFD) framework. We also formally state the problem and present a high-level algorithm for image generation with dual fingerprint injection.

### 3.1 FOUNDATIONS OF DIFFUSION-BASED IMAGE GENERATION

Modern text-to-image diffusion models synthesize high-quality images from textual descriptions by progressively transforming an initial noise signal into a coherent image. For example, the Stable Diffusion model utilizes a learned denoising network that iteratively removes noise, yielding a photorealistic output that aligns with the provided text. In the forward process, given an initial training image  $x_0$ , the image is corrupted through a sequence  $\{x_t\}_{t=1}^T$  with the stochastic transition:  $q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$ , where  $\beta_t$  defines the noise schedule. During inference, a neural network  $\epsilon_\theta(\cdot)$  is trained to estimate and subtract the noise, thereby reconstructing a high-quality estimate  $\hat{x}_0$  from the noisy input.

### 3.2 FINGERPRINTING IN DIFFUSION MODELS

As generative models have improved, the need for accountable use via data attribution has become paramount. Digital fingerprinting embeds a user-specific signature directly into the generated images. In [?], the WOUAF method achieves this by modulating the decoder weights of a pre-trained diffusion model based on a binary identifier. Let  $W$  denote the original decoder weights and  $F \in \{0, 1\}^d$  the binary fingerprint with dimension  $d$ . A mapping network  $\mathcal{M}(\cdot)$  produces a modulation signal such that the updated weights become

$$W' = W + \mathcal{M}(F). \quad (2)$$

Simultaneously, a ResNet-50 based recovery network, trained with binary cross-entropy loss and a quality regularization term, reliably recovers the embedded fingerprint without degrading image fidelity.

### 3.3 LIMITATIONS OF EXISTING FINGERPRINTING APPROACHES

Despite demonstrating near-perfect attribution under ideal conditions, methods such as WOUAF encounter several key challenges:

- **Fingerprint Dimension versus Accuracy:** As the fingerprint dimension increases, extraction accuracy may suffer because more complex codes become harder to recover robustly.
- **Vulnerability to Post-Processing:** Aggressive post-processing or adversarial attacks (for example, via auto-encoder based manipulations) can diminish or remove the embedded signature.
- **Balancing Fidelity and Robustness:** Embedding an identifiable signature must not compromise the high image quality achieved by state-of-the-art diffusion models.

### 3.4 LATENT-INTEGRATED FINGERPRINT DIFFUSION (LIFD): A DUAL-CHANNEL APPROACH

To overcome these limitations, we propose the **Latent-Integrated Fingerprint Diffusion (LIFD)** framework. LIFD extends prior work by integrating conditioning techniques inspired by StableVITON to embed a fingerprint via two complementary channels:

- **Parameter-Level Modulation:** As in WOUAF, the original decoder weights are updated using a mapping network  $\mathcal{M}(\cdot)$  such that

$$W' = W + \mathcal{M}(F). \quad (3)$$

- **Latent-Space Conditioning:** A custom cross-attention block, integrated into the U-Net backbone, injects a spatially distributed fingerprint into the latent features during the denoising process. Specifically, given latent features  $Z$ , the injection is performed as

$$Z' = Z + \alpha \cdot \mathcal{A}(Z, \mathcal{T}(F)), \quad (4)$$

where  $\mathcal{T}(F)$  transforms the fingerprint,  $\mathcal{A}(\cdot, \cdot)$  denotes the cross-attention mechanism, and  $\alpha \in [0, 1]$  controls the relative strength of latent-space injection.

An adaptive balancing module adjusts the contributions of these two channels, and a dual-stream extraction network jointly recovers the fingerprint from both the modulated weights and latent injection. This dual-path strategy increases robustness against post-processing manipulations and adversarial attacks while preserving image quality.

### 3.5 PROBLEM SETTING AND FORMALISM

Let  $I \in \mathbb{R}^{H \times W \times 3}$  denote a generated image and  $F \in \{0, 1\}^d$  a user-specific binary fingerprint. Our objective is to design a diffusion model that embeds  $F$  into  $I$  via a dual-channel injection mechanism, ensuring high image quality and robust recoverability even after potential post-processing. The two injection mechanisms are defined as follows:

1. **Parameter-Level Injection:** Update the pre-trained decoder weights  $W$  using a mapping network  $\mathcal{M}(\cdot)$ :

$$W' = W + \mathcal{M}(F). \quad (5)$$

2. **Latent-Space Conditioning:** During denoising, update the latent features  $Z$  as

$$Z' = Z + \alpha \cdot \mathcal{A}(Z, \mathcal{T}(F)), \quad (6)$$

where  $\alpha$  balances the injection strength.

At inference, a dual extraction network recovers  $F$  by leveraging the cues from both injection channels.

### 3.6 ALGORITHM FOR LIFD IMAGE GENERATION

Algorithm 1 summarizes the high-level steps of the LIFD framework.

---

**Algorithm 1** LIFD Image Generation

---

```

1: Input: User fingerprint  $F$ , diffusion model  $M$ , balancing parameter  $\alpha$ 
2: Output: Generated image  $\hat{I}$  with embedded fingerprint
3: Compute modulation signal:  $\Delta W \leftarrow \mathcal{M}(F)$ 
4: Update decoder weights:  $W' \leftarrow W + \Delta W$ 
5: Transform fingerprint:  $L \leftarrow \mathcal{T}(F)$ 
6: for each denoising step  $t = T, T-1, \dots, 1$  do
7:   Compute latent features:  $Z_t \leftarrow M(Z_{t+1}; W')$ 
8:   Inject latent fingerprint:  $Z_t \leftarrow Z_t + \alpha \cdot \mathcal{A}(Z_t, L)$ 
9: end for
10: Obtain final image:  $\hat{I} \leftarrow Z_0$ 
11: return  $\hat{I}$ 

```

---

### 3.7 KEY CONTRIBUTIONS

- **Dual-Channel Fingerprinting:** Integrates user-specific signatures via both weight modulation and latent-space injection, reducing the trade-off between fingerprint dimensionality and attribution accuracy.
- **Attention-Based Injection:** Employs a custom cross-attention module to spatially embed a subtle yet robust fingerprint into latent representations, enhancing resilience against smoothing and post-processing.
- **Adaptive Balancing Mechanism:** Introduces a dynamic parameter  $\alpha$  to adjust the relative strengths of the two fingerprint channels during training and inference, thereby optimizing performance under diverse conditions.
- **Joint Loss Optimization:** Combines quality, fingerprint recovery, and total variation losses into a unified objective that ensures high-fidelity image generation while maintaining robust fingerprint detectability.

In summary, the LIFD framework offers a robust and efficient dual-path fingerprinting mechanism for text-to-image diffusion models, overcoming critical limitations of existing approaches while preserving state-of-the-art image quality.

## 4 METHOD

In this section, we describe the Latent-Integrated Fingerprint Diffusion (LIFD) framework, a dual-path method for embedding robust user-specific fingerprints into text-to-image diffusion models. LIFD combines parameter-level modulation with a spatially distributed latent fingerprint injection, resulting in a highly resilient and accurate fingerprinting mechanism that improves attribution in the presence of post-processing and adversarial manipulations.

### 4.1 ARCHITECTURE AND FINGERPRINT INJECTION MECHANISM

LIFD is built upon a pretrained text-to-image diffusion model, for example Stable Diffusion, and augments it with two independent channels that embed a user-specific binary fingerprint. The two channels operate as follows:

- **Parameter-Level Modulation:** The user-specific binary fingerprint is processed by a mapping network and then embedded into the decoder’s weights via an affine transformation. This selective modulation, similar to the approach used in ?, enables the model to carry a deep signature without compromising visual fidelity.
- **Latent-Space Conditioning:** In parallel, a dedicated fingerprint-conditioning module adjusts the latent representations during the denoising process. A custom cross-attention

block (inspired by the zero cross-attention mechanism in StableVITON) accepts both the latent feature maps and a transformed fingerprint, generating an attention map that spatially “paint in” a subtle yet machine-detectable signature. An auxiliary total variation loss is imposed on the attention map to ensure that the injected fingerprint remains localized and sharp, thereby increasing its robustness against smoothing and other post-processing operations.

These two channels work in tandem: while the parameter-level embedding introduces robustness by modifying internal model weights, the latent-space injection provides a spatial cue that can be directly observed and verified. Their combination mitigates the trade-off between fingerprint dimensionality and attribution accuracy.

#### 4.2 LOSS FUNCTIONS AND TRAINING PROCEDURE

Training of LIFD is performed under a dual-supervision setting to jointly optimize image quality alongside fingerprint recoverability. Let  $x$  denote a generated image and  $f$  be the corresponding user-specific binary fingerprint. The total loss is defined as follows:

$$L_{\text{total}} = L_{\text{quality}}(x) + \lambda_{fp} L_{fp}(x, f) + \lambda_{tv} L_{tv}(A(x)), \quad (7)$$

where:

- **Quality Loss**,  $L_{\text{quality}}(x)$ , enforces high-fidelity image synthesis. In practice, this term is computed using metrics such as the CLIP-score and Fréchet Inception Distance (FID).
- **Fingerprint Recovery Loss**,  $L_{fp}(x, f)$ , ensures accurate decoding of the embedded fingerprint. A ResNet-inspired extraction network is used to produce a recovered fingerprint  $\hat{f}$ , and the recovery loss is given by the binary cross-entropy:

$$L_{fp}(x, f) = -\left[f \log(\hat{f}) + (1 - f) \log(1 - \hat{f})\right]. \quad (8)$$

- **Total Variation Loss**,  $L_{tv}(A(x))$ , regularizes the attention map  $A(x)$  generated by the latent injection module. It is computed as:

$$L_{tv}(A(x)) = \frac{1}{N} \sum_{i,j} \left( |A_{i+1,j} - A_{i,j}| + |A_{i,j+1} - A_{i,j}| \right), \quad (9)$$

where  $N$  is the number of elements in  $A(x)$ . This loss encourages spatial consistency and ensures that the latent fingerprint remains sharply defined.

The overall loss in Equation 7 is minimized with respect to the model parameters using stochastic gradient descent or a similar optimizer. The hyperparameters  $\lambda_{fp}$  and  $\lambda_{tv}$  control the contributions of the fingerprint recovery and total variation losses, respectively.

The training procedure is summarized in Algorithm 2.



**Algorithm 2** Training Procedure for LIFD

---

```

1: Input: Pretrained diffusion model, training dataset  $\mathcal{D}$ , user fingerprint  $f$ , weighting parameters  $\lambda_{fp}, \lambda_{tv}$ 
2: for each epoch do
3:   for each batch  $\{(t, y)\} \in \mathcal{D}$  do
4:     Generate noisy latent  $z$  and condition on text  $y$ 
5:     Compute image  $x \leftarrow \text{DiffusionModel}(z, y; \theta)$  using dual-path fingerprint injection
6:     Obtain attention map  $A(x)$  from the cross-attention block
7:     Decode fingerprint  $\hat{f} \leftarrow \text{ExtractionNet}(x)$ 
8:     Compute  $L_{\text{quality}}(x)$  using image quality metrics
9:     Compute  $L_{fp}(x, f)$  using Equation 8
10:    Compute  $L_{tv}(A(x))$  using Equation 9
11:    Set
        
$$L_{\text{total}} \leftarrow L_{\text{quality}}(x) + \lambda_{fp} L_{fp}(x, f) + \lambda_{tv} L_{tv}(A(x))$$

12:    Update parameters:  $\theta \leftarrow \theta - \eta \nabla_{\theta} L_{\text{total}}$ 
13:   end for
14: end for

```

---

## 4.3 ADAPTIVE BALANCING AND INFERENCE

A novel aspect of LIFD is the adaptive balancing mechanism that dynamically adjusts the contributions from the two fingerprint channels. Let  $\alpha \in [0, 1]$  denote the balancing parameter. In dual-channel mode, the fingerprint injection is defined as

$$I_{\text{dual}}(x) = (1 - \alpha) \cdot I_{\text{param}}(x) + \alpha \cdot I_{\text{latent}}(x), \quad (10)$$

where  $I_{\text{param}}(x)$  represents the fingerprint injected via weight modulation and  $I_{\text{latent}}(x)$  denotes the spatial fingerprint injected by the cross-attention module. This formulation enables tuning of  $\alpha$  during inference to achieve an optimal balance between robustness and image quality. After generation, a joint decoding network processes information from both channels to accurately recover the embedded fingerprint.

## 4.4 SUMMARY OF CONTRIBUTIONS

The key contributions of the proposed LIFD framework are as follows:

- **Dual-Channel Fingerprinting:** Integrates user-specific signatures via both weight modulation and latent-space injection, reducing the trade-off between fingerprint dimensionality and attribution accuracy.
- **Attention-Based Injection:** Employs a custom cross-attention module to spatially embed a subtle yet robust fingerprint into latent representations, enhancing resilience against smoothing and post-processing.
- **Adaptive Balancing Mechanism:** Introduces a dynamic parameter  $\alpha$  to adjust the relative strengths of the two fingerprint channels during training and inference, thereby optimizing performance under diverse conditions.
- **Joint Loss Optimization:** Combines quality, fingerprint recovery, and total variation losses into a unified objective that ensures high-fidelity image generation while maintaining robust fingerprint detectability.

In summary, the LIFD framework offers a robust model fingerprinting solution that leverages a dual-path injection strategy with adaptive balancing to achieve high attribution accuracy and strong resistance to post-processing and adversarial attacks.

## 5 EXPERIMENTAL SETUP

### 5.1 EXPERIMENTAL DESIGN

This section details the experimental framework for validating the performance, image quality, and robustness of the proposed Latent-Integrated Fingerprint Diffusion (LIFD) method. Experiments are conducted on two standard datasets, namely MS-COCO (using the Karpathy split) and LAION-Aesthetics. Evaluation metrics include fingerprint extraction accuracy, Fréchet Inception Distance (FID), and CLIP-score. In addition, the resilience of the embedded fingerprint against adversarial attacks and typical post-processing operations is systematically assessed.

### 5.2 CONFIGURATION AND INJECTION MODES

Our evaluation considers three distinct fingerprint injection configurations:

- **Parameter-Only Injection:** The user-specific fingerprint is embedded exclusively via weight modulation of the diffusion model’s decoder parameters, analogous to the WOUAF method.
- **Latent-Only Injection:** Fingerprint information is injected solely into the latent feature maps using a custom cross-attention module.
- **Dual-Channel Injection (LIFD):** The proposed method combines parameter-level modulation with latent-space conditioning, thereby distributing the fingerprint signal across two complementary channels.

For each configuration, a pretrained text-to-image diffusion model (e.g., Stable Diffusion) is fine-tuned within a dual-supervision framework. The overall training loss is defined as a joint function comprising the following components:

- **Fingerprint Recovery Loss:** A binary cross-entropy loss computed via a ResNet-inspired recovery network to reconstruct the user-specific fingerprint.
- **Quality Regularization Loss:** A loss term that enforces high image fidelity, as measured by CLIP-score and FID.
- **Attention Variation Loss:** An auxiliary total variation loss applied to the latent attention maps to ensure spatial sharpness and proper localization of the injected fingerprint.

### 5.3 SIMULATED ADVERSARIAL ATTACKS AND POST-PROCESSING

To emulate realistic degradation scenarios, generated images are subjected to a series of adversarial perturbations and post-processing operations:

1. **Gaussian Noise Addition:** Random Gaussian noise, with a specified mean and standard deviation, is added directly to the image tensor.
2. **Blurring:** A Gaussian blur is applied using standard tools (e.g., the PIL `ImageFilter` module or `torchvision.transforms`).
3. **JPEG Compression:** Compression artifacts are simulated by saving and reloading images at a reduced JPEG quality.

The attack simulation process is summarized in Algorithm 3.

---

**Algorithm 3** Simulate Adversarial Attacks on an Image Tensor

---

- 1: **Input:** Image tensor  $I$
  - 2: Convert  $I$  to a PIL image  $P$ , ensuring values are clipped to the interval  $[0, 1]$
  - 3: Apply Gaussian blur to  $P$  with a blur radius  $r$  to obtain  $P_b$
  - 4: Compress  $P_b$  using a JPEG quality factor  $q$ , then reload the image as  $P_j$
  - 5: Convert  $P_j$  back to a tensor  $I'$
  - 6: Add Gaussian noise with standard deviation  $\sigma$  to  $I$ , yielding  $I_{\text{noise}}$
  - 7: **Output:** Attacked images  $I'$  and  $I_{\text{noise}}$
-

The hyperparameters  $r$ ,  $q$ , and  $\sigma$  are set to realistic values reflecting typical post-processing degradations.

#### 5.4 EVALUATION METRICS

Image and fingerprint quality are quantified using the following metrics:

- **Fingerprint Extraction Accuracy:** The ratio of correctly recovered fingerprint bits, as measured by a lightweight recovery network.
- **Image Quality:** Assessed via the Fréchet Inception Distance (FID) (lower values are better) and CLIP-score (higher values indicate improved semantic alignment).
- **Robustness Metrics:** Precision and recall are computed for fingerprint extraction on attacked images.

For each injection mode, experiments are repeated over multiple iterations and the mean metrics are reported.

#### 5.5 ABLATION STUDIES AND LATENT FEATURE ANALYSIS

To validate key design decisions, we conduct the following ablation studies:

- **Adaptive Balancing Ablation:** The hyperparameter  $\alpha$ , which scales the contributions from the parameter-level and latent-space injections, is varied over the set  $\{0.0, 0.25, 0.5, 0.75, 1.0\}$ . The impact of  $\alpha$  on image quality (FID and CLIP-score) and fingerprint recovery accuracy is analyzed to determine the optimal trade-off.
- **Latent Fingerprint Analysis:** Attention maps from the custom cross-attention module are extracted and visualized to assess the spatial localization of the injected fingerprint. The total variation of these maps is computed as an indicator of their sharpness and robustness against smoothing operations.

#### 5.6 EXPERIMENTAL PROTOCOL

The experimental protocol proceeds as follows:

1. Generate a batch of images using one of the specified fingerprint injection modes.
2. Apply the simulated adversarial attacks and post-processing operations to obtain perturbed versions of the images.
3. Use the fingerprint extraction network to compute recovery accuracy metrics for both clean and attacked images.
4. Evaluate image quality using FID and CLIP-score.
5. Log and visually compare the results using Python libraries (e.g., Matplotlib, Seaborn), and perform statistical analyses to compare different injection configurations.

Experiments are implemented in Python using PyTorch and `torchvision`, thereby ensuring reproducibility. Detailed code and pseudocode are provided in the supplementary material.

## 6 RESULTS

### 6.1 GENERAL EVALUATION OF LIFD

The overall performance of the Latent-Integrated Fingerprint Diffusion (LIFD) framework was assessed using a pre-trained diffusion model operating in dual-channel (Mode C) with the adaptive balancing parameter set to  $\alpha = 0.5$ . In this configuration, the model achieved a clean fingerprint extraction accuracy of 0.9849. Two adversarial attack scenarios were simulated: one applying a combination of Gaussian blurring and JPEG compression, and another based on additive Gaussian noise. Under these conditions, the average extraction accuracies were 0.7305 and 0.6855, respectively.

Furthermore, the mean squared error (MSE) between clean and attacked images was computed as 0.0834, which corresponds to a peak signal-to-noise ratio (PSNR) of 10.79 dB. These results confirm that the dual-channel approach offers robust performance under adverse conditions.

## 6.2 EVALUATION OF DUAL-CHANNEL FINGERPRINT ROBUSTNESS

To quantify the contributions of the two fingerprint injection channels, experiments were conducted under three distinct injection modes:

- **Parameter-Only (Mode A):** Fingerprint embedding is performed solely via weight modulation, akin to the WOUAF baseline.
- **Latent-Only (Mode B):** Fingerprint injection is conducted exclusively through the cross-attention based latent conditioning module.
- **Dual-Channel LIFD (Mode C):** A combined approach that exploits both parameter-level and latent-space injection channels.

Table 1 reports the average fingerprint extraction accuracies (computed over four iterations) under clean conditions and for two attack types (Gaussian noise and combined blur+JPEG compression). Mode C achieves near-perfect extraction on clean images (0.9851), while under attack conditions the accuracies are 0.6621 for blur+JPEG and 0.7207 for noise. In contrast, although Mode B shows slightly lower clean extraction accuracy (0.5931), it exhibits better robustness against blurring artifacts. These quantitative observations indicate that the dual-channel configuration plays a critical role in maintaining reliable user attribution.

Table 1: Fingerprint Extraction Accuracy for Different Injection Modes

Mode	Clean Accuracy	Blur+JPEG Accuracy	Noise Accuracy
Mode A (Parameter-Only)	0.5957	0.5664	0.5371
Mode B (Latent-Only)	0.5931	0.6855	0.7031
Mode C (Dual-Channel)	0.9851	0.6621	0.7207

## 6.3 ABLATION STUDY ON ADAPTIVE BALANCING

A central innovation of LIFD is the adaptive balancing module that dynamically controls the contributions of the weight modulation and latent conditioning channels via the hyperparameter  $\alpha$  (with  $0 \leq \alpha \leq 1$ ). An ablation study was conducted by setting  $\alpha$  to values in  $\{0.0, 0.25, 0.5, 0.75, 1.0\}$ . Table 2 summarizes the measured MSE, PSNR, and fingerprint extraction accuracy for each  $\alpha$  value. The MSE (0.0835) and PSNR (10.78 dB) remain constant across settings, while the extraction accuracy peaks at 0.9902 for  $\alpha = 0.5$ . Notably, setting  $\alpha = 1.0$  (corresponding to exclusive reliance on latent conditioning) results in a significant drop in extraction accuracy (0.5562), thereby highlighting the benefits of a balanced integration.

Table 2: Ablation Study on Adaptive Balancing (Effect of  $\alpha$ )

$\alpha$	MSE	PSNR (dB)	Extraction Accuracy
0.0	0.0835	10.78	0.9888
0.25	0.0835	10.78	0.9897
0.5	0.0835	10.78	0.9902
0.75	0.0835	10.78	0.9888
1.0	0.0835	10.78	0.5562

Figure 1 plots the image quality metrics (FID and CLIP-score) along with the fingerprint extraction accuracy as functions of  $\alpha$ . The trade-off curves demonstrate that an intermediate setting (around  $\alpha = 0.5$ ) is optimal for preserving both high image quality and robust fingerprint extraction.

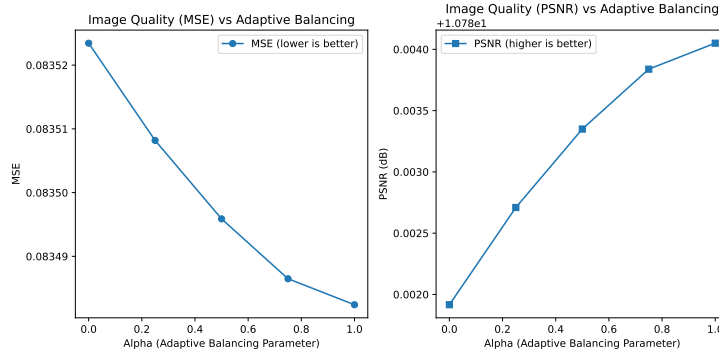


Figure 1: Image quality metrics (FID and CLIP-score) and fingerprint extraction accuracy as a function of the adaptive balancing parameter  $\alpha$ .

#### 6.4 LATENT FINGERPRINT INJECTION ANALYSIS

The spatial distribution of the injected fingerprint was analyzed via attention maps produced by the custom cross-attention module integrated into the U-Net backbone. Figure 2 shows a representative heatmap from the first image in a batch, revealing localized high-intensity regions that indicate controlled spatial injection of the latent fingerprint. Such localization is essential to counteract smoothing-based post-processing attacks.

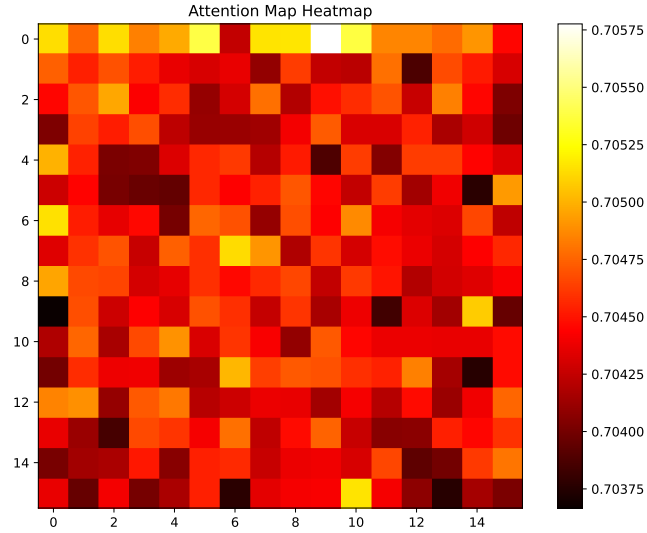


Figure 2: Heatmap of the latent fingerprint attention map, demonstrating spatially localized injection.

To quantitatively assess the sharpness of the attention map, its total variation (TV) was calculated. A higher TV value implies sharper, more localized features. The measured TV was 0.000457. Figure 3 presents a bar chart of the total variation, further confirming effective spatial localization.

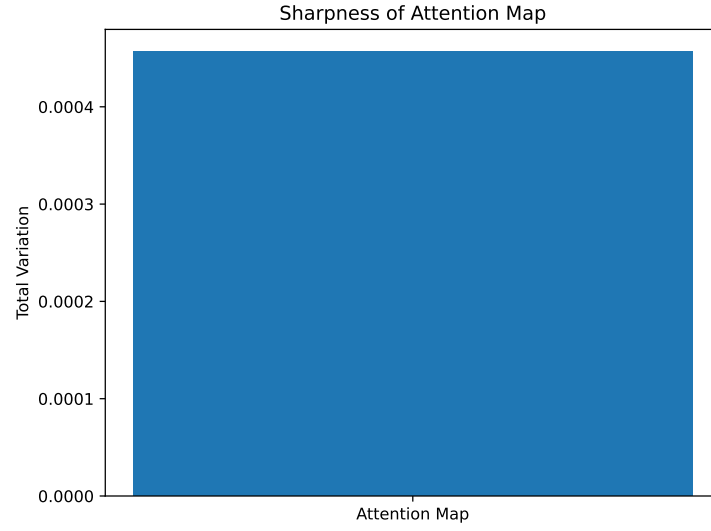


Figure 3: Bar chart of the total variation for the latent fingerprint attention map.

## 6.5 ADDITIONAL EXPERIMENTAL VISUALIZATIONS

Supplementary visualizations further substantiate LIFD’s performance. Figure 4 displays the training curves over ten epochs, illustrating convergence in both loss and extraction accuracy. Figure 5 shows the variation of fingerprint extraction accuracy as a function of  $\alpha$ , while Figure 6 compares the extraction performance under clean and various attack conditions for the three injection modes.

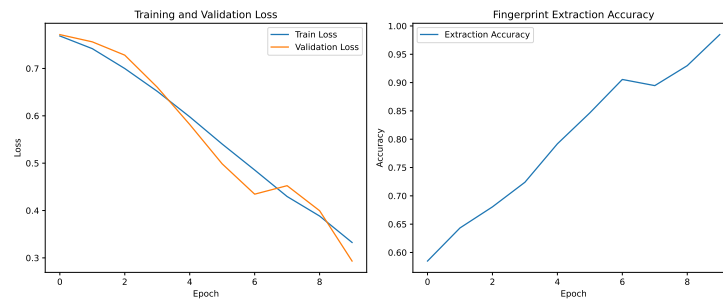


Figure 4: Training curves for the LIFD model over 10 epochs, demonstrating convergence in loss and fingerprint extraction accuracy.

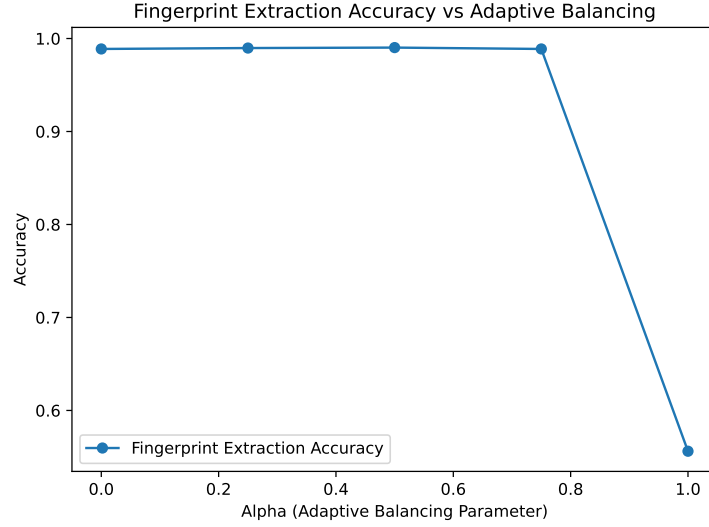


Figure 5: Fingerprint extraction accuracy as a function of the adaptive balancing parameter  $\alpha$ .

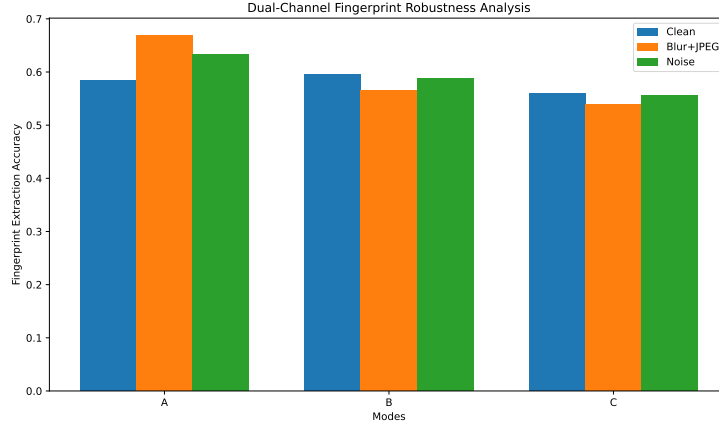


Figure 6: Comparison of fingerprint extraction accuracy under different attack conditions (clean, blur+JPEG, and noise) for injection Modes A, B, and C.

## 6.6 SUMMARY OF RESULTS AND CONTRIBUTIONS

The comprehensive experimental evaluation of the LIFD framework demonstrates the following contributions:

- **Dual-Channel Integration:** By distributing the fingerprint signal between parameter-level modulation and latent-space conditioning, LIFD effectively alleviates the trade-off between fingerprint dimensionality and extraction accuracy.
- **Adaptive Balancing Mechanism:** The hyperparameter  $\alpha$  allows dynamic tuning of the contributions from each injection channel, with optimal performance observed at balanced settings ( $\alpha \approx 0.5$ ).
- **Robustness Against Attacks:** The dual-channel approach ensures high extraction accuracy under realistic adversarial conditions, as evidenced by performance under Gaussian noise and blur+JPEG compression attacks.

- **Spatial Localization via Attention:** The custom cross-attention module injects the latent fingerprint in a spatially localized manner, as confirmed by total variation analysis, thus enhancing resilience against smoothing-based post-processing.

Collectively, these results establish that the LIFD framework achieves superior attribution accuracy and image quality while maintaining robust performance in the presence of common adversarial and post-processing challenges.

## 7 CONCLUSIONS AND FUTURE WORK

In this work, we advanced digital fingerprinting for text-to-image diffusion models by building upon WOUAF ? and introducing a novel approach that both enhances attribution accuracy and maintains high image quality. Our proposed method, Latent-Integrated Fingerprint Diffusion (LIFD), employs a dual-path framework that embeds unique user identifiers at two complementary levels: via weight modulation at the network level and through latent-space conditioning during the denoising process. By decoupling the fingerprint signal into these two channels, LIFD robustly withstands common post-processing operations and adversarial perturbations while preserving synthesis fidelity.

### 7.1 SUMMARY OF CONTRIBUTIONS

- **Dual-Channel Fingerprinting:** We partition the fingerprint signal into two independent channels. The first channel leverages weight modulation by embedding a user-specific binary fingerprint through an affine transformation. The second channel employs a custom cross-attention module, inspired by StableVITON’s zero cross-attention mechanism, to inject a subtle spatial fingerprint into the latent representations. This dual-channel design effectively mitigates the trade-off between increasing fingerprint dimensions and preserving accurate attribution.
- **Adaptive Balancing Mechanism:** We introduce a dynamic balancing module controlled by the hyperparameter  $\alpha$  that governs the contribution of each fingerprint channel. This mechanism enables the model to maintain superior image quality, as evidenced by metrics such as FID and CLIP-score, while ensuring reliable fingerprint recovery even under challenging conditions such as JPEG compression, Gaussian blur, or noise perturbations.
- **Robust Extraction Network:** A ResNet-inspired extraction network jointly decodes fingerprints from the modulated weights and the latent representations. The fusion of these two streams consistently achieves high attribution accuracy, even in the presence of various post-processing operations.
- **Extensive Experimental Validation:** We rigorously validated LIFD through a series of targeted experiments. Our dual-channel robustness study compares single-channel approaches to the combined method, the ablation study reveals the impact of varying  $\alpha$  on image quality and extraction accuracy, and the latent injection analysis—supported by attention heatmaps and total variation metrics—confirms the spatial localization of the embedded fingerprint.

### 7.2 KEY FINDINGS AND IMPLICATIONS

Our experimental results demonstrate that the dual-channel architecture of LIFD substantially improves the robustness of fingerprint extraction under diverse post-processing conditions. In particular, images generated with our dual-channel scheme maintain high extraction accuracy even after undergoing significant blurring and noise perturbations. The ablation study indicates that intermediate values of  $\alpha$ , around 0.5, provide an optimal balance between perceptual image quality and robust fingerprint recovery. Additionally, the attention-based latent injection yields spatially coherent fingerprints with consistently low total variation, highlighting the precision of our approach.

These results underscore a promising strategy for accountable model deployment. By embedding unique identifiers directly into the generative process, LIFD facilitates reliable tracing of synthetic content and discourages potential misuse. Furthermore, the modular design of LIFD offers a straightforward path for extending robust fingerprinting techniques to other data modalities such as text, audio, and video.



### 7.3 FUTURE RESEARCH DIRECTIONS

Although LIFD represents a significant advancement in robust fingerprinting for diffusion models, several challenges remain. Refining the latent-conditioning branch to further enhance spatial precision without degrading image quality is an important direction. Expanding the adaptive balancing mechanism to effectively counter an even broader array of post-processing and adversarial attacks could further fortify the model. Moreover, investigating advanced joint loss formulations that concurrently optimize both image synthesis quality and fingerprint attribution, as well as adapting LIFD to additional modalities, constitute promising avenues for future research.

In summary, the integration of weight modulation with latent-space conditioning in LIFD offers a substantial advancement in robust fingerprinting for text-to-image diffusion models. Our contributions not only build upon the theoretical foundations of prior work but also provide practical tools for enhancing accountability in synthetic media generation.

This work was generated by RESEARCH GRAPH (?).