GRADCFG: GRADIENT INVERSION OF CLASSIFIER-FREE GUIDANCE DIFFUSION MODELS

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Gradient inversion attacks, as a means of privacy theft, have been extensively studied and applied in classifier models, yet research on gradient inversion for diffusion models, particularly classifier-free guidance (CFG) diffusion models, remains relatively underdeveloped. CFG models such as Stable Diffusion present significant challenges for such attacks due to their complex training mechanisms, including the high-dimensional search space caused by multimodal variables, the non-uniqueness of the noise ϵ solution space, and the difficulty in optimizing discrete time steps t. To address these challenges, this paper proposes a novel joint inversion framework featuring two core algorithmic innovations: the GradCFG algorithm, which integrates a four-variable co-optimization mechanism for simultaneous reconstruction of image latent variables \mathbf{x}_0 , text embeddings C_0 , noise ϵ , and reparameterized continuous time steps t, alongside a periodic restart strategy for ϵ to enhance solution stability and generalization; and the **Inv-Sam** algorithm, a model-difference-based generation optimization method that leverages the generative capability disparities between pre-fine-tuning and post-fine-tuning models to restore high-resolution details through a reverse-forward diffusion editing process. Systematic experiments in CFG model fine-tuning scenarios demonstrate that the proposed method effectively achieves high-quality image-text joint reconstruction for various textual conditions ranging from concise descriptions to complex semantic combinations.

1 Introduction

Diffusion models (Ho et al., 2020) have achieved remarkable breakthroughs in the field of image generation by transforming random noise into high-fidelity images through a progressive denoising process. Classifier-Free Guidance (CFG) (Ho & Salimans, 2022) further introduces a text-conditioning mechanism, enabling semantically controllable image synthesis. Representative CFG models such as Stable Diffusion (SD) (Rombach et al., 2022) have been widely adopted in industry. With the growing demand for personalized generation, users often fine-tune pre-trained models on private data (Gal et al., 2022; Kumari et al., 2023; Hu et al., 2021; Bahmani et al., 2022). To protect privacy, users typically share only training gradients with the server instead of the original images (Sun et al., 2021). However, this process still faces severe privacy leakage risks: malicious attackers can reconstruct private training samples from the gradients via Gradient Inversion Attacks. Such attacks have been extensively demonstrated in classification models (Hatamizadeh et al., 2022; Zhu et al., 2019; Wei et al., 2020), yet the unique training mechanism of diffusion models provides inherent defense capabilities. In particular, for CFG models, the text-guided mechanism further increases the difficulty of gradient inversion.

The challenge of gradient inversion for CFG models stems from their complex multimodal training process: the client must independently sample Gaussian noise ϵ and a time step t, and compute gradients using the image latent \mathbf{x}_0 and text embedding C_0 generated by a pre-trained encoder. This expands the recovery target from a unimodal image to a bimodal image-text pair, significantly increasing the complexity of reconstruction. The key difficulty lies in the need to simultaneously optimize four sets of variables (\mathbf{x}_0 , C_0 , ϵ , and t) to achieve gradient alignment. However, our research reveals that the solution space of the sampled noise ϵ is non-unique (i.e., multiple solutions satisfy gradient alignment). Thus, we must adapt \mathbf{x}_0 , C_0 , and t to different ϵ within the solution space. Additionally, the discrete variable t cannot be directly updated via conventional optimization methods.

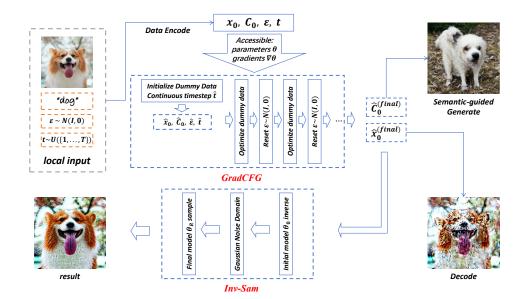


Figure 1: The GradCFG method reconstructs the original data through gradient matching, employing a time-step \hat{t} continuous strategy and a periodic $\hat{\epsilon}$ reset mechanism to mitigate variable discreteness and solution non-uniqueness. Furthermore, the Inv-Sam method refines the initially reconstructed image by leveraging the generative capability gap between the initial and final state models, thereby recovering richer image details.

Therefore, there is an urgent need to design a coordination mechanism: one that enables the recovered variables to adapt to the ϵ solution space while also facilitating effective optimization of the discrete time step. Furthermore, leveraging the inherent generation and editing capabilities of CFG models, without relying on external pre-trained models, to assist in high-resolution image detail reconstruction represents another critical issue.

To address these challenges, this paper proposes an innovative solution consisting of two interconnected algorithmic contributions, as shown in Figure 1. First, we propose the **GradCFG** algorithm, which constructs a four-variable joint optimization framework that reparameterizes the discrete time step t as a continuous variable \hat{t} , enabling simultaneous optimization of $\hat{\mathbf{x}}_0$, \hat{C}_0 , $\hat{\epsilon}$, and \hat{t} . Through periodic restart and re-optimization of $\hat{\epsilon}$, the GradCFG algorithm allows $\hat{\mathbf{x}}_0$, \hat{C}_0 , and \hat{t} to achieve gradient alignment across different $\hat{\epsilon}$ values. Second, we develop the **Inv-Sam** algorithm, a model-disparity optimization strategy that leverages the difference in generative capabilities between the pre-trained model at initial and final fine-tuning stages. This approach orchestrates a reverse-forward diffusion process guided by the pre-trained and fine-tuned models respectively, strategically injecting model-disparity information into the reconstruction pipeline to significantly enhance textual alignment and visual detail in recovered images.

Experiments are conducted under personalized fine-tuning scenarios for CFG models, employing a DREAMBOOTH-like (Ruiz et al., 2023) fine-tuning paradigm to systematically evaluate the method's effectiveness. Our comprehensive evaluation framework encompasses both general textual prompts and specific textual prompts during fine-tuning, enabling a thorough analysis of reconstruction performance across different semantic granularities. We simultaneously assess the recovery quality of both image data and textual embeddings, providing a holistic evaluation of multimodal privacy leakage. Additionally, we conduct ablation studies to systematically investigate the specific impact of the periodic restart mechanism on reconstruction fidelity, while also designing controlled experiments to validate the non-uniqueness of solutions for the noise variable ϵ . The proposed approach is rigorously validated on multiple fine-tuning datasets containing diverse semantic categories, demonstrating robust performance across varying image contents and textual descriptions.

The core contributions of this paper are as follows:

• For the first time, we empirically demonstrate a joint image-text privacy leakage attack in text-guided diffusion models (CFG), opening a new attack surface.

- We propose the first joint optimization framework that simultaneously reconstructs images, text, noise, and time steps, overcoming the challenge of variable coupling under complex training mechanisms.
- We innovatively leverage the generative capability disparity of diffusion models during training to design a reconstruction optimization method, significantly enhancing the detail restoration and visual quality of high-resolution images.

2 RELATED WORK

Classifier-Free Guidance (CFG) Models and Privacy. Diffusion model variants integrating CFG have become mainstream architectures in text-to-image generation (Rombach et al., 2022; Dhariwal & Nichol, 2021; Saharia et al., 2022; Ramesh et al., 2022). Current research on CFG model privacy primarily focuses on Membership Inference Attacks (Shokri et al., 2017), Model Inversion Attacks (Zhou et al., 2024), and Training Data Extraction (Carlini et al., 2023), while privacy leakage risks from gradient inversion attacks during training remain systematically underexplored.

Gradient Inversion. Early pioneering work (e.g., DLG proposed by Zhu et al. (2019)) first demonstrated the feasibility of reconstructing training data from gradients, building upon earlier recognition of gradients as a primary leakage channel in collaborative learning (Podschwadt et al., 2022). The attack was primarily effective for small batches (e.g., size=1) and low-resolution images. Subsequent research significantly enhanced its practicality and scope: Geiping et al. (2020) improved reconstruction quality on complex datasets (e.g., ImageNet) by introducing a cosine similarity loss and critical regularizations like Total Variation (Zhu & Blaschko, 2021); Zhao et al. (2020) developed an analytical method to deduce labels from gradients exactly; Yin et al. (2021) proposed GradInversion, which leveraged Batch Normalization statistics and group consistency to tackle larger batches and higher-resolution images. Recently, methodologies have diversified: Generation-based attacks (GEN-GIA) (Wei et al., 2020; Jeon et al., 2021) employ pre-trained generative models (e.g., GANs, diffusion models) as strong priors to produce high-fidelity reconstructions, but their reliance on external data and sensitivity to architectures limit generality. Analytics-based attacks (ANA-GIA) (Gao et al., 2021; He et al., 2019) derive data through maliciously altering model parameters or analyzing model outputs, which is efficient but operates under a strong threat model. In addition to direct attacks, Tian et al. (2025) explored reconstructing data by analyzing the weight differences between pre-training and post-training states of models. Notably, Huang et al. (2025a) investigated gradient inversion in diffusion models, but their approach fails to address the fundamental challenges of sampling noise ϵ multiplicity and time step t discontinuity, while also being constrained by dependency on pre-trained generators. Currently, research on gradient inversion attacks targeting diffusion models, particularly their widely adopted CFG mechanism, remains notably scarce (Yu et al., 2024).

3 METHODOLOGY

This paper focuses on the fine-tuning scenario of CFG models, employing a training paradigm **similar to Dreambooth** (Ruiz et al., 2023). The objective is to reconstruct private training data—including images and their corresponding text embeddings—from gradient information.

3.1 CFG MODEL FINE-TUNING FRAMEWORK

In this distributed training framework, each client (user) maintains four private data elements during training round r: Raw image X, Text prompt P, Sampled time step $\mathbf{t} \sim \mathcal{U}(\{1,\ldots,T\})$, Gaussian noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0},\mathbf{I})$. None of these elements are directly shared with the server. The fine-tuning process is shown in Algorithm. 1. The training objective minimizes the following denoising loss function:

$$\mathcal{L}(\theta_r) = \mathbb{E}_{t, \boldsymbol{x}_0, \boldsymbol{\epsilon}} \left[\left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta_r} \left(\sqrt{\bar{\alpha}_t} \boldsymbol{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t, C_0 \right) \right\|^2 \right], \tag{1}$$

where the frozen pretrained image encoder $\mathrm{Vae}(\cdot)$ encodes the raw image X into latent representation $\boldsymbol{x}_0 = \mathrm{Vae}(X)$, and the text encoder $\mathrm{Encoder}(\cdot)$ maps the text prompt P to conditional embedding $C_0 = \mathrm{Encoder}(P)$; θ_r denotes the learnable parameters of the diffusion model at round $r, \bar{\alpha}_t$ is the hyperparameter controlling the noise schedule in the forward diffusion process, and $\epsilon_{\theta_r}(\cdot)$ represents the noise prediction network taking (\cdot) as input.

Algorithm 1 CFG Model Training Process

Input: Training rounds R, user dataset \mathcal{X} , text prompt P, diffusion steps T, learning rate η Pretrained model θ_0 , pretrained text encoder Encoder, pretrained VAE Vae(·)

for training round $r = 1, 2, \dots, R$ **do**

User execution:

Encode image: $x_0 \leftarrow \text{Vae}(X)$ Encode text: $C_0 \leftarrow \text{Encoder}(P)$ Sample time step: $t \sim \mathcal{U}(\{1, \dots, T\})$ Sample noise: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

Compute loss function: $\mathcal{L}(\theta_r) = \mathbb{E}_{t, \boldsymbol{x}_0, \boldsymbol{\epsilon}} \left[\left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta_r} \left(\sqrt{\bar{\alpha}_t} \boldsymbol{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t, C_0 \right) \right\|^2 \right]$

Compute gradient: $g^{(r)} = \nabla_{\theta_r} \mathcal{L}(\theta_r)$

Send $g^{(r)}$ to server

Server execution: Update model: $\theta_{r+1} \leftarrow \theta_r - \eta g^{(r)}$

Output: Optimized model parameters θ_R

Clients locally compute the gradient $g^{(r)} = \nabla_{\theta_r} \mathcal{L}(\theta_r)$ and transmit only this gradient to the server. The server coordinates R training rounds by aggregating gradients and updating global model parameters.

3.2 Gradient Inversion Attack Methodology (GradCFG)

3.2.1 ATTACK MODELING

A malicious server can inversely reconstruct users' private data \mathbf{x}_0 and text embeddings C_0 when only accessing model gradients $g^{(r)}$, where noise vector $\boldsymbol{\epsilon}$ and time step t remain private user information. This attack process is formalized as the following multivariate optimization problem:

$$\min_{\hat{\mathbf{x}}_0, \hat{C}_0, \hat{\boldsymbol{\epsilon}}, \hat{t}} \mathcal{D}\left(\nabla_{\theta_r} \mathcal{L}(\hat{\mathbf{x}}_0, \hat{C}_0, \hat{\boldsymbol{\epsilon}}, \hat{t}; \theta_r), g^{(r)}\right) + \eta(s) \mathcal{L}_{\text{mix}}(\hat{\mathbf{x}}_0)$$
 (2)

where $\mathcal{D}(\cdot,\cdot)$ is the gradient similarity metric function (cosine similarity is adopted in this paper), $g^{(r)} = \nabla_{\theta_r} \mathcal{L}(\mathbf{x}_0, C_0, \boldsymbol{\epsilon}, t; \theta_r)$ represents the observed true gradient, and $\nabla_{\theta_r} \mathcal{L}(\hat{\mathbf{x}}_0, \hat{C}_0, \hat{\boldsymbol{\epsilon}}, \hat{t}; \theta_r)$ is the gradient based on virtual parameters. $\mathcal{L}_{\text{mix}}(\hat{\mathbf{x}}_0)$ is defined as the disentanglement regularizer, implemented by computing the mean cosine similarity of all data pairs in the reconstructed image set $\{\hat{\mathbf{x}}_0^{(1)}, ..., \hat{\mathbf{x}}_0^{(k)}\}$:

$$\mathcal{L}_{\text{mix}}(\hat{\mathbf{x}}_0) = \frac{2}{k(k-1)} \sum_{i=1}^{k-1} \sum_{j=i+1}^{k} \frac{\langle \hat{\mathbf{x}}_0^{(i)}, \hat{\mathbf{x}}_0^{(j)} \rangle}{\|\hat{\mathbf{x}}_0^{(i)}\|_2 \cdot \|\hat{\mathbf{x}}_0^{(j)}\|_2}$$
(3)

The regularization strength is dynamically controlled by the iteration step scheduler $\eta(s)$, where s denotes the current optimization iteration count:

$$\eta(s) = \begin{cases} \eta_{\text{max}} & s < S_{\text{switch}} \\ 0 & s \ge S_{\text{switch}} \end{cases}$$
(4)

This scheduling strategy preserves the full regularizer $\eta_{\max}\mathcal{L}_{\min}(\hat{\mathbf{x}}_0)$ during the early training phase $(s < S_{\text{switch}})$ to enforce feature disentanglement and prevent feature mixing in reconstructed samples. During the later training phase $(s \geq S_{\text{switch}})$, the regularizer constraint is completely removed, allowing the optimization process to focus solely on minimizing the gradient difference $\mathcal{D}(\cdot,\cdot)$. Here, S_{switch} is a preset iteration count threshold controlling the transition from the feature disentanglement phase to the precision optimization phase.

3.2.2 QUADRUPLE COLLABORATIVE OPTIMIZATION ALGORITHM

Defining the pseudo-gradient objective function $\mathcal{G} = \mathcal{D}\left(\nabla_{\theta_r}\mathcal{L}(\hat{\mathbf{x}}_0, \hat{C}_0, \hat{\boldsymbol{\epsilon}}, \hat{t}; \theta_r), g^{(r)}\right)$ with its corresponding gradient denoted as $\nabla \mathcal{G}$, we propose a quadruple collaborative optimization framework:

Image Reconstruction (\mathbf{x}_0 optimization): To maintain latent space consistency, an initial point $\hat{X}_0 \sim \mathcal{N}(0, I)$ is sampled from image space and projected into latent space via the VAE encoder: $\hat{\mathbf{x}}_0^{(0)} = \text{VAE}(\hat{X}_0)$ Update rule: $\hat{\mathbf{x}}_0 \leftarrow \hat{\mathbf{x}}_0 - \eta_x \nabla_{\hat{\mathbf{x}}_0} \mathcal{G}$

Text Reconstruction (C_0 optimization): Initialized with empty text $\hat{P} = \phi$, projected through the text encoder: $\hat{C}_0^{(0)} = \text{Encoder}(\hat{P})$ Update rule: $\hat{C}_0 \leftarrow \hat{C}_0 - \eta_C \nabla_{\hat{C}_0} \mathcal{G}$

Time Step Reconstruction (t optimization): To optimize the originally discrete time step $\hat{t} \in \{1,\ldots,T\}$ in continuous space, we propose a function reparameterization strategy. Addressing the discrete nature of the noise scheduler α_t , we establish a continuous mapping through mathematical transformation. Specifically, considering the definition $\alpha_t = \prod_{i=1}^t (1-\beta_i)$ where $\beta_i = f(i)$ is a predefined discrete scheduling function (e.g., linear or cosine decay), when $\beta_i \ll 1$, we utilize the natural logarithm approximation $\ln(1-\beta_i) \approx -\beta_i$ to transform the discrete summation $\sum_{i=1}^t \ln(1-\beta_i)$ into integral form:

$$\ln \alpha_t \approx -\sum_{i=1}^t \beta_i \approx -\int_0^t f(x)dx \tag{5}$$

This derivation yields the continuous time step representation: $\alpha(t) \approx \exp\left(-\int_0^t f(x)dx\right)$.

This continuous representation enables time step \hat{t} to be updated via standard gradient descent: $\hat{t} \leftarrow \hat{t} - \eta_t \nabla_{\hat{t}} \mathcal{G}$

Noise Reconstruction with Dynamic Reset (ϵ optimization): Initialize $\hat{\epsilon} \sim \mathcal{N}(0, I)$. We incorporate a dynamic reset mechanism where $\hat{\epsilon}$ is randomly resampled from $\mathcal{N}(0, I)$ at fixed intervals S_{reset} :

$$\hat{\epsilon} \sim \mathcal{N}(0, I) \quad \text{when} \quad s \equiv 0 \pmod{S_{\text{reset}}}$$
 (6)

This mechanism continuously searches for new solutions within the solution space, allowing other optimization variables $(\hat{\mathbf{x}}_0, \hat{C}_0, \text{ and } \hat{t})$ to adapt to different solutions and achieve optimal performance. The update rule is: $\hat{\boldsymbol{\epsilon}} \leftarrow \hat{\boldsymbol{\epsilon}} - \eta_{\epsilon} \nabla_{\hat{\boldsymbol{\epsilon}}} \mathcal{G}$

3.3 BIDIRECTIONAL SAMPLING ENHANCEMENT ALGORITHM (INV-SAM)

The preliminary reconstruction results $(\hat{\mathbf{x}}_0, \hat{C}_0) \in \mathbb{R}^{m \times B} \times \mathbb{R}^{77 \times 768}$ obtained through gradient inversion at round r are mapped to natural image space via the pretrained VAE decoder: $\hat{\mathbf{X}} = \mathrm{VAE}^{-1}(\hat{\mathbf{x}}_0)$ However, $\hat{\mathbf{X}}$ suffers from insufficient visual fidelity and missing high-frequency details. We leverage the dynamic evolution of generative capabilities during fine-tuning: compared to the initial model θ_0 , the final model θ_R generates images with richer training-set features under text embedding C_0 guidance. This paper proposes using the generative capability difference $\theta_R - \theta_0$ as an optimization prior.

Inspired by reverse diffusion and forward sampling in image editing (Miyake et al., 2024; Huang et al., 2025b), we design a text-guided latent optimization method (Algorithm 2). Using the recovered condition \hat{C}_0 , we first perform inverse diffusion with θ_0 to project $\hat{\mathbf{x}}_0$ into noise space, then execute sampling using the generative difference $\theta_R - \theta_0$ to reproject to latent space. Analysis of the inverse and sampling path relationship reveals that under path proximity (Miyake et al., 2024), their difference is approximately proportional to the model prediction difference:

$$\mathbf{x}_t^{\text{sam}} - \mathbf{x}_t^{\text{inv}} \propto \omega_{\text{sam}} \cdot \underbrace{\left(\epsilon_{\theta_R}(\mathbf{x}_{t+1}^{\text{sam}}, t+1, \hat{C}_0) - \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{sam}}, t+1, \hat{C}_0)\right)}_{\Delta \epsilon_{\theta}}$$

where $\Delta \epsilon_{\theta}$ quantifies the directional correction from fine-tuning. A formal proof is provided in Appendix D.

Algorithm 2 Inv-Sam Optimization

Input: Initial latent state $\hat{\mathbf{x}}_0 \in \mathbb{R}^m$, reconstructed text embedding $\hat{C}_0 \in \mathbb{R}^{77 \times 768}$ Reverse step guidance factor $\omega_{\text{inv}} = 1$, sampling step guidance factor ω_{sam} Noise schedule $\{\bar{\alpha}_t\}_{t=0}^T$

Fine-tuned model parameters θ_R , initial model parameters θ_0

Phase I: Inverse Diffusion $\mathbf{x}_0^{\text{inv}} \leftarrow \hat{\mathbf{x}}_0$

$$\begin{array}{l} \text{for } t = 0 \text{ to } T - 1 \text{ do} \\ & \tilde{\epsilon}^{\text{inv}} \leftarrow \epsilon_{\theta_0}(\mathbf{x}_t^{\text{inv}}, t, \hat{C}_0) \\ & \mathbf{x}_{t+1}^{\text{inv}} \leftarrow \sqrt{\bar{\alpha}_{t+1}} \left(\frac{\mathbf{x}_t^{\text{inv}} - \sqrt{1 - \bar{\alpha}_t} \tilde{\epsilon}^{\text{inv}}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t+1}} \tilde{\epsilon}^{\text{inv}} \end{array}$$

Phase II: Conditional Sampling $\mathbf{x}_T^{\text{sam}} \leftarrow \mathbf{x}_T^{\text{inv}}$

for t = T - 1 to 0 do

$$\begin{split} & \epsilon_{\text{empty}} \leftarrow \epsilon_{\theta_0}(\mathbf{x}^{\text{sam}}_{t+1}, t+1, \hat{C}_0) \\ & \epsilon_{\text{text}} \leftarrow \epsilon_{\theta_R}(\mathbf{x}^{\text{sam}}_{t+1}, t+1, \hat{C}_0) \\ & \tilde{\epsilon}^{\text{sam}} \leftarrow \epsilon_{\text{empty}} + \omega_{\text{sam}}(\epsilon_{\text{text}} - \epsilon_{\text{empty}}) \\ & \mathbf{x}^{\text{sam}}_t \leftarrow \sqrt{\bar{\alpha}_t} \left(\frac{\mathbf{x}^{\text{sam}}_{t+1} - \sqrt{1 - \bar{\alpha}_{t+1}} \tilde{\epsilon}^{\text{sam}}}{\sqrt{\bar{\alpha}_{t+1}}} \right) + \sqrt{1 - \bar{\alpha}_t} \tilde{\epsilon}^{\text{sam}} \end{split}$$

Output: Optimized latent state $\hat{\mathbf{x}}_{0}^{opt} \leftarrow \mathbf{x}_{0}^{sam}$

4 EXPERIMENTAL EVALUATION

4.1 EXPERIMENTAL SETUP

Experiments utilize the standard DREAMBOOTH (Ruiz et al., 2023) training dataset containing image samples at 512×512 resolution. All models are fine-tuned using the TinySD (Kim et al., 2023) framework to ensure parameter-efficient optimization. Two experimental scenarios are designed:

Gradient inversion for generic text prompt fine-tuning: Constructs a category-uniform scenario where images of similar objects (e.g., backpacks in different environments) are fine-tuned using unified text prompts (e.g., "backpack").

Gradient inversion for specific text prompt fine-tuning: Constructs a fine-grained control scenario where each image is paired with a dedicated granular prompt (e.g., "a red backpack on a mountain trail") to evaluate reconstruction performance under complex text conditions.

To the best of our knowledge, this is the first work to investigate gradient inversion attacks on CFG-based diffusion models. Currently, no existing methods can jointly reconstruct both images and text prompts from gradients in this context. Therefore, our evaluation focuses on comparing the performance of our complete method (GradCFG + Inv-Sam) against its ablated variants to isolate the contribution of each component, rather than comparing against external baselines.

4.2 EVALUATION METRICS

To comprehensively assess reconstruction quality, we employ both quantitative and qualitative evaluation frameworks. Image reconstruction quality is measured using three complementary metrics: Peak Signal-to-Noise Ratio (PSNR) (Wang et al., 2004) for pixel-level fidelity assessment, Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018) for human visual perception similarity (with lower values indicating better performance), and the Structural Similarity Index Measure (SSIM) (Wang et al., 2004) for assessing the perceptual quality related to structural information, luminance, and contrast.

For semantic recovery evaluation of generic prompts, we implement a dual-modality analysis framework that computes cosine similarity between reconstructed and original prompts in embedding space to quantify semantic consistency, while also generating images guided by original and reconstructed prompts under identical initial noise conditions using the TinySD model to compute PSNR between paired images as an indirect measure of semantic recovery effectiveness.

5 RESULTS

5.1 IMAGE RECONSTRUCTION ANALYSIS

We conducted gradient inversion experiments under two distinct fine-tuning scenarios: generic text prompts and detailed text prompts. The reconstruction process employed a two-stage approach: initial image reconstruction was performed using the GradCFG method, followed by detail enhancement via the Inv-Sam algorithm.

Table 1 presents comparative reconstruction results under different text prompts. Experimental results demonstrate that baseline gradient inversion achieves preliminary reconstruction for both prompt types, with comparable quantitative metrics (structural similarity SSIM~0.12, peak signal-to-noise ratio PSNR~10.6 dB, perceptual similarity LPIPS~0.78). Inv-Sam optimization significantly enhances reconstruction quality: under generic prompts, SSIM increases by 77% to 0.219 while LPIPS decreases by 24% to 0.591. For specific prompts, SSIM improvement reaches 55% with a 20% reduction in LPIPS.

Table 1: Experimental results of GradCFG and Inv-Sam under different prompt settings (Values in parentheses indicate metric changes after applying Inv-Sam)

Setting	GradCFG		+ Inv-Sam			
2311 g	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓
Generic prompts	0.1240	10.60	0.7778	0.2189 (+77%)	11.65 (+10%)	0.5911 (-24%)
Specific prompts	0.1213	10.65	0.7789	0.1875 (+55%)	11.51 (+8%)	0.6226(-20%)

Visual results are presented in Figure 2 (generic prompts) and Figure 3 (specific prompts). The gradient inversion stage effectively captures object contours and base colors, establishing structural frameworks. Subsequently, Inv-Sam refinement enhances textural details, enabling reconstruction of fine-grained features.



Figure 2: Reconstruction workflow under generic prompts: GradCFG results (left), Inv-Sam optimized results (middle), Original picture (right)

Comprehensive analysis confirms the effective synergy: gradient inversion recovers structural frameworks, while Inv-Sam enhances texture details. This mechanism achieves greater metric improvements under generic prompts while maintaining robustness for specific prompt scenarios.

5.2 TEXT EMBEDDING RECOVERY ANALYSIS

In the fine-tuning scenario with generic text prompts, we conducted a systematic analysis of text encoding embeddings. The evaluation procedure consisted of three sequential steps: First, we computed the cosine similarity between the recovered text embeddings and the original text embeddings.

A backpack sitting on top of a rock with mountains in the background

A red backpack sitting on a tree branch

A woman with a backpack looking up at the sky

A red backpack sitting on the ground in the woods

A red backpack hanging on a tree branch

GradCFG +Inv-Sam Original

Figure 3: Reconstruction results for backpack category images under various scenes.

Subsequently, both the recovered and original embeddings were separately employed as conditional guidance inputs to a pre-trained CFG model to generate corresponding images. Finally, we quantified the similarity between these generated image pairs using PSNR, thereby validating the semantic consistency of the recovered embeddings.

Quantitative results presented in Table 2 demonstrate that the gradient inversion method effectively recovers generic text embeddings while maintaining high semantic fidelity (cosine similarity: 0.7953; PSNR: 16.28 dB). Visual comparisons in Figure 4 further confirm that images generated from recovered embeddings exhibit strong semantic.

Table 2: Recovery Performance for Generic Embeddings

Similarity	PSNR
0.7953	16.28

prompt

alignment with those produced from original embeddings, validating the method's effectiveness in preserving semantic information integrity.

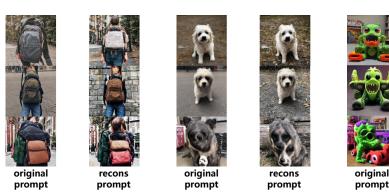


Figure 4: Generated image comparison based on text embeddings (left to right: backpack, dog, monster toy). For each group: left image generated from original prompt embedding, right image from reconstructed text embedding.

5.3 Impact of ϵ Reset Cycle on Reconstruction Quality

This experiment systematically investigates the influence of different ϵ reset cycles on image reconstruction quality in gradient inversion attacks. As shown in Figure 5, under a fixed total optimization budget of 4000 iterations, we evaluated reconstruction performance at reset cycles of 1, 10, 100, 1000, and 4000 (no reset). Experimental results reveal significant performance variations based on reset

cycle selection. When the reset cycle is too short (e.g., 1 or 10), the ϵ parameter cannot sufficiently optimize to convergence regions, causing estimation bias in latent variables. This manifests as reconstructed images with clear pixel-level details but noticeable structural misalignments and semantic inconsistencies. Conversely, excessively long reset cycles (e.g., 1000) or no reset (4000) prevent adequate exploration of diversity within the ϵ solution space, causing optimization to stagnate in local minima. This results in reconstructed images with blurred details and lacking texture features.

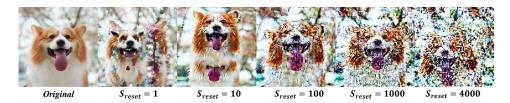


Figure 5: Impact of ϵ reset cycles on image reconstruction quality. From left to right: reset cycles = 1, 10, 100, 1000, 4000 (no reset). Experimental results demonstrate that a moderate reset cycle (100) optimally balances optimization stability and solution space exploration capability.

In conclusion, the ϵ reset cycle requires careful balancing between optimization stability and solution space exploration capacity. Our experimental findings indicate that a reset cycle of 100 achieves optimal reconstruction quality, ensuring sufficient ϵ parameter optimization while maintaining effective solution space exploration.

5.4 Analysis of Non-Uniqueness in ϵ Solutions

In this experiment, we attempt to recover the noise ϵ under known conditions of \mathbf{x}_0 , C_0 , and t, while simultaneously monitoring the similarity evolution between the simulated gradient g and the original gradient g_0 , as well as the similarity between the recovered noise $\hat{\epsilon}$ and the ground-truth noise ϵ .

As shown in Figure 6, we observe that even when g and g_0 exhibit high similarity (exceeding 95%), significant discrepancies persist between $\hat{\epsilon}$ and ϵ . Notably, $\hat{\epsilon}$ fails to converge to the original noise even during later optimization stages. This demonstrates the existence of multiple distinct ϵ configurations capable of achieving high gradient alignment, indicating solution non-uniqueness.

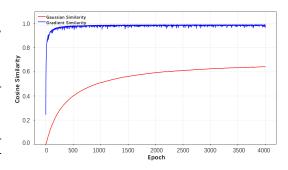


Figure 6: Optimization of randomly sampled noise ϵ under fixed conditions of \mathbf{x}_0 , C_0 , and t, while monitoring gradient alignment and noise similarity relative to original data.

Therefore, our reconstruction methodology incorporates a strategy of periodically resetting and re-optimizing ϵ . This approach ensures the recovered variables remain robust across diverse valid ϵ solutions, thereby enhancing the generalization capability of the reconstruction outcomes.

6 CONCLUSION

This study introduces a novel gradient inversion technique for reconstructing high-resolution images and their corresponding text prompts during CFG-based model training. The proposed GradCFG method enables, for the first time, simultaneous recovery of both training images and associated text prompts, overcoming a key limitation of prior approaches that struggle with joint visual and semantic reconstruction. An enhancement module, Inv-Sam, leverages the generative gap between fine-tuned and initial models as prior knowledge, substantially improving image quality and semantic accuracy. Experiments conducted under a DREAMBOOTH-like fine-tuning setup using TinySD models demonstrate high-fidelity reconstruction of 512×512 complex scenes and accurate text recovery. The method performs robustly across both generic and specific prompts, regardless of complexity.

REFERENCES

- Shervin Bahmani, Andrea Park, Daniel Kappler, Bernt Schiele, and Vladislav Golyanik. Prompt tuning for text-based image editing. *arXiv preprint arXiv:2212.09608*, 2022.
- Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models, 2023. URL https://arxiv.org/abs/2301.13188.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021. URL https://arxiv.org/abs/2105.05233.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations*, 2022.
- Wei Gao, Shangwei Guo, Tianwei Zhang, Han Qiu, Yonggang Wen, and Yang Liu. Privacy-preserving collaborative learning with automatic transformation search, 2021. URL https://arxiv.org/abs/2011.12505.
- Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients how easy is it to break privacy in federated learning?, 2020. URL https://arxiv.org/abs/2003.14053.
- Ali Hatamizadeh, Hongxu Yin, Holger Roth, Wenqi Li, Jan Kautz, Daguang Xu, and Pavlo Molchanov. Gradvit: Gradient inversion of vision transformers, 2022. URL https://arxiv.org/abs/2203.11894.
- Zecheng He, Tianwei Zhang, and Ruby B. Lee. Model inversion attacks against collaborative inference. In *Proceedings of the 35th Annual Computer Security Applications Conference*, AC-SAC '19, pp. 148–162, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450376280. doi: 10.1145/3359789.3359824. URL https://doi.org/10.1145/3359789.3359824.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. Introduces the Classifier-Free Guidance (CFG) technique for diffusion models.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *CoRR*, abs/2006.11239, 2020. URL https://arxiv.org/abs/2006.11239.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeb Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- Jiyue Huang, Chi Hong, Stefanie Roos, and Lydia Y. Chen. Gidm: Gradient inversion of federated diffusion models. In Mila Dalla Preda, Sebastian Schrittwieser, Vincent Naessens, and Bjorn De Sutter (eds.), *Availability, Reliability and Security*, pp. 380–401, Cham, 2025a. Springer Nature Switzerland. ISBN 978-3-032-00624-0.
- Yi Huang, Jiancheng Huang, Yifan Liu, Mingfu Yan, Jiaxi Lv, Jianzhuang Liu, Wei Xiong, He Zhang, Liangliang Cao, and Shifeng Chen. Diffusion model-based image editing: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(6):4409–4437, June 2025b. ISSN 1939-3539. doi: 10.1109/tpami.2025.3541625. URL http://dx.doi.org/10.1109/TPAMI.2025.3541625.
- Jinwoo Jeon, Jaechang Kim, Kangwook Lee, Sewoong Oh, and Jungseul Ok. Gradient inversion with generative image prior, 2021. URL https://arxiv.org/abs/2110.14962.
- Bo-Kyeong Kim, Hyoung-Kyu Song, Thibault Castells, and Shinkook Choi. Bk-sdm: Architecturally compressed stable diffusion for efficient text-to-image generation. *ICML Workshop on Efficient Systems for Foundation Models (ES-FoMo)*, 2023. URL https://openreview.net/forum?id=bOVydU0XKC.

- Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9461–9471, 2023.
 - Daiki Miyake, Akihiro Iohara, Yu Saito, and Toshiyuki Tanaka. Negative-prompt inversion: Fast image inversion for editing with text-guided diffusion models, 2024. URL https://arxiv.org/abs/2305.16807.
 - Robert Podschwadt, Daniel Takabi, and Peizhao Hu. Sok: Privacy-preserving deep learning with homomorphic encryption, 2022. URL https://arxiv.org/abs/2112.12855.
 - Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. URL https://arxiv.org/abs/2204.06125.
 - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 10674–10685. IEEE, 2022.
 - 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. URL https://arxiv.org/abs/2208.12242.
 - Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. URL https://arxiv.org/abs/2205.11487.
 - Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models, 2017. URL https://arxiv.org/abs/1610.05820.
 - Tao Sun, Dongsheng Li, and Bao Wang. Decentralized federated averaging, 2021. URL https://arxiv.org/abs/2104.11375.
 - Hanling Tian, Yuhang Liu, Mingzhen He, Zhengbao He, Zhehao Huang, Ruikai Yang, and Xiaolin Huang. Simulating training dynamics to reconstruct training data from deep neural networks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=ZJftXKy12x.
 - Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. doi: 10.1109/TIP.2003.819861.
 - Wenqi Wei, Ling Liu, Margaret Loper, Ka-Ho Chow, Mehmet Emre Gursoy, Stacey Truex, and Yanzhao Wu. A framework for evaluating gradient leakage attacks in federated learning, 2020. URL https://arxiv.org/abs/2004.10397.
 - Hongxu Yin, Arun Mallya, Arash Vahdat, Jose M. Alvarez, Jan Kautz, and Pavlo Molchanov. See through gradients: Image batch recovery via gradinversion, 2021. URL https://arxiv.org/ abs/2104.07586.
 - Lei Yu, Meng Han, Yiming Li, Changting Lin, Yao Zhang, Mingyang Zhang, Yan Liu, Haiqin Weng, Yuseok Jeon, Ka-Ho Chow, and Stacy Patterson. A survey of privacy threats and defense in vertical federated learning: From model life cycle perspective, 2024. URL https://arxiv.org/abs/2402.03688.
 - Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric, 2018. URL https://arxiv.org/abs/1801.03924.
 - Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. idlg: Improved deep leakage from gradients, 2020. URL https://arxiv.org/abs/2001.02610.

Zhanke Zhou, Jianing Zhu, Fengfei Yu, Xuan Li, Xiong Peng, Tongliang Liu, and Bo Han. Model inversion attacks: A survey of approaches and countermeasures, 2024. URL https://arxiv. org/abs/2411.10023. Junyi Zhu and Matthew Blaschko. R-gap: Recursive gradient attack on privacy, 2021. URL https://arxiv.org/abs/2010.07733. Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients, 2019. URL https: //arxiv.org/abs/1906.08935.

A ETHICS STATEMENT

This work investigates potential privacy vulnerabilities in classifier-free guidance (CFG) diffusion models during personalized fine-tuning. We strictly adhere to the ICLR Code of Ethics and are fully aware of the dual-use implications of our research. To mitigate potential risks, all experiments exclusively utilize publicly available benchmark datasets (DREAMBOOTH), ensuring no private or sensitive data is involved. The fundamental objective of this research is to raise community awareness about privacy leakage threats in distributed learning scenarios, with the ultimate goal of contributing to more secure federated learning frameworks. We emphasize that the defensive value of understanding these vulnerabilities significantly outweighs any offensive potential. While our method demonstrates reconstruction capabilities, we strongly oppose any malicious application of this technique and believe transparent analysis of such attack vectors is essential for developing robust defenses.

B REPRODUCIBILITY STATEMENT

To ensure the reproducibility of this research, comprehensive efforts have been made to provide complete implementation details and resources. Full implementation code, pretrained model weights, and evaluation scripts are available in the supplementary materials. The core methodologies (GradCFG and Inv-Sam algorithms) are thoroughly described in Sec. 3, including optimization objectives, loss functions, and key reparameterization strategies. Complete hyperparameter configurations (learning rates, reset period S_{reset} , etc.) are documented in Appendix. E.

C ON THE USE OF LARGE LANGUAGE MODELS

Large language models (LLMs) were employed solely for writing assistance, including surface-level text editing (grammar correction, clarity improvement), document formatting, and experimental code comment generation. LLMs did not contribute to originating research ideas, claims, or conclusions. The authors take full responsibility for all intellectual content. All LLM-assisted text was carefully reviewed and rewritten by the authors to ensure accurate expression of the research.

D PROOF OF INV-SAM

Why path proximity implies proportional correction. We formalize this relationship as the following Prop. D.1 and provide a dimension-free proof based on linear operator analysis.

Proposition D.1 (Proportionality of sampling-inverse path difference to model prediction difference). Consider the DDIM framework with linear update function $g_t(\mathbf{x}, \epsilon) = A_t \mathbf{x} + B_t \epsilon$, where

$$A_t = \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{\bar{\alpha}_{t+1}}}, \quad B_t = \sqrt{1 - \bar{\alpha}_t} - \frac{\sqrt{\bar{\alpha}_t}\sqrt{1 - \bar{\alpha}_{t+1}}}{\sqrt{\bar{\alpha}_{t+1}}}.$$
 (7)

Define the path difference $\delta_t = \mathbf{x}_t^{sam} - \mathbf{x}_t^{inv}$, model prediction difference $\Delta \epsilon_\theta = \epsilon_{\theta_R}(\mathbf{x}_{t+1}^{sam}, t+1, \hat{C}_0) - \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{sam}, t+1, \hat{C}_0)$, and sampling noise $\tilde{\epsilon}^{sam} = \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{sam}, t+1, \hat{C}_0) + \omega_{sam}\Delta \epsilon_{\theta}$. Under the path proximity assumption $\mathbf{x}_{t+1}^{sam} \approx \mathbf{x}_{t+1}^{inv}$ (implying $\delta_{t+1} \approx 0$ and similar denoising outputs), we have

$$\delta_t \approx B_t \omega_{sam} \Delta \epsilon_{\theta}. \tag{8}$$

In particular, the path difference is proportional to the guidance-weighted model prediction difference:

$$\mathbf{x}_t^{sam} - \mathbf{x}_t^{inv} \propto \omega_{sam} \cdot \Delta \epsilon_{\theta} \tag{9}$$

with proportionality constant B_t depending only on the noise schedule.

Proof of Prop. D.1. Starting from the definition of the path difference and substituting the update operations:

$$\delta_t = g_t(\mathbf{x}_{t+1}^{\text{sam}}, \hat{\epsilon}^{\text{sam}}) - g_t(\mathbf{x}_{t+1}^{\text{inv}}, \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{inv}}, t+1, \hat{C}_0))$$
(10)

Exploiting the linearity of the update function g_t :

$$\delta_t = A_t(\mathbf{x}_{t+1}^{\text{sam}} - \mathbf{x}_{t+1}^{\text{inv}}) + B_t(\tilde{\epsilon}^{\text{sam}} - \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{inv}}, t+1, \hat{C}_0))$$

$$\tag{11}$$

Substituting δ_{t+1} and applying the path proximity assumption ($\delta_{t+1} \approx 0$):

$$\delta_t \approx B_t(\tilde{\epsilon}^{\text{sam}} - \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{inv}}, t+1, \hat{C}_0))$$
(12)

Expanding the sampling noise $\tilde{\epsilon}^{\text{sam}}$:

$$\delta_t \approx B_t \left[\epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{sam}}, t+1, \hat{C}_0) + \omega_{\text{sam}} \Delta \epsilon_{\theta} - \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{inv}}, t+1, \hat{C}_0) \right]$$
(13)

By path proximity, the denoising outputs are similar:

$$\epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{sam}}, t+1, \hat{C}_0) \approx \epsilon_{\theta_0}(\mathbf{x}_{t+1}^{\text{inv}}, t+1, \hat{C}_0)$$

$$\tag{14}$$

Thus the terms cancel, yielding the final result:

$$\delta_t \approx B_t \omega_{\text{sam}} \Delta \epsilon_{\theta} \tag{15}$$

This establishes the proportionality with schedule-dependent constant B_t .

The interpretation reveals that B_t represents the time-dependent scaling from the noise schedule, while $\omega_{\rm sam}$ directly controls the amplification of model corrections. This proportionality demonstrates that our latent optimization algorithm effectively translates fine-tuning improvements into controlled path deviations, maintaining the delicate balance between faithfulness to the input and incorporation of desired model enhancements.

E FURTHER EXPERIMENTAL DETAILS

All experiments were conducted using the pre-trained TinySD model as the base architecture (a lightweight version of Stable Diffusion with 3×10^8 parameters) on NVIDIA A800 80GB PCIe GPU platforms. The optimization process employed the Adam optimizer ($\beta_1=0.8,\beta_2=0.9$) with the following learning rate configuration: 0.1 for image latent variables $\hat{\mathbf{x}}_0$, 0.001 for text embeddings \hat{C}_0 , and 0.1 for both noise parameters $\hat{\epsilon}$ and timesteps \hat{t} . The total number of iterations was fixed at 4000, with periodic reset of noise parameters $\hat{\epsilon}$ implemented every 100 steps. To prevent vanishing gradients during optimization, we constrained the temporal range of \hat{t} to values between 400 and 600 throughout the reconstruction process. A default batch size of B=5 was used throughout, with gradient alignment measured using a cosine similarity-based distance function $\mathcal{D}(\cdot)=1-\cos(\cdot)$. The feature decoupling regularization term $\mathcal{S}(\cdot)$ was activated during the first 100 iterations to enhance initial convergence stability. All experiments were performed under identical configuration environments to ensure result comparability and reproducibility.

F DETAILED EXPERIMENTAL RESULTS OF GRADCFG

F.1 DETAILED RESULTS FOR FINETUNING EXPERIMENTS WITH GENERIC TEXT PROMPTS

This section provides the complete experimental results for the finetuning experiments using generic text prompts, serving as supplementary data to Section 5.1. Table 3 presents the comprehensive performance comparison between baseline GradCFG and our Inv-Sam enhanced approach across all object categories.

Table 4 provides the detailed semantic recovery metrics, including embedding similarity scores and image similarity measurements for each category.

F.2 DETAILED RESULTS FOR FINETUNING EXPERIMENTS WITH SPECIFIC TEXT PROMPTS

This section provides the complete experimental results for finetuning experiments using specific text prompts, serving as supplementary data to Section 5.1. Table 5 presents the comprehensive performance comparison between baseline GradCFG and our Inv-Sam enhanced approach under complex textual conditions.

Table 3: GradCFG results using generic text prompts for fine-tuning.

Category	GradCFG			+ Inv-Sam		
	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓
Backpack	0.1267	10.84	0.7222	0.2441	11.97	0.5672
Can	0.1706	10.10	0.8073	0.3984	11.70	0.5813
Candle	0.0684	10.43	0.8096	0.1251	12.16	0.5640
Cat	0.1710	12.20	0.6745	0.2259	13.76	0.5321
Sneaker	0.0555	9.97	0.9142	0.1024	11.48	0.6376
Dog	0.2095	10.98	0.7446	0.3271	12.53	0.5690
Monster Toy	0.1298	10.43	0.7705	0.2706	11.78	0.5770
Robot Toy	0.0692	9.28	0.7801	0.1165	10.20	0.6102
Race Car	0.1170	10.15	0.7702	0.1623	10.87	0.6635
Average Improv.	0.1240	10.60 Baseline	0.7778	0.2189 +77%	11.65 +10%	0.5911 -24%

Table 4: Evaluation of semantic recovery for text prompts

Category	Embedding Sim.	Image Sim. (PSNR)
Backpack	0.7538	14.94
Can	0.9228	16.69
Candle	0.8000	17.52
Cat	0.8457	18.08
Sneaker	0.7010	13.06
Dog	0.8766	19.49
Monster Toy	0.7330	13.10
Robot Toy	0.7368	15.30
Race Car	0.7878	17.32
Avg.	0.7953	16.28

For complex text recovery tasks, this experiment first analyzes the text reconstruction performance using backpack-related prompts as a detailed case study, followed by a comprehensive analysis across all text categories. It is important to note that complete reconstruction of all textual information is not our primary objective, as full recovery of complex semantic content presents significant challenges. Instead, we focus on evaluating the improvement in similarity between recovered text embeddings and original data compared to null text embeddings. As demonstrated in Table 6, our method achieves significant improvements in this measured similarity for specific prompt examples. Table 7 further presents the overall performance across all categories, showing that our approach substantially enhances text embedding recovery quality by this metric, effectively demonstrating the utility of our method without requiring complete semantic reconstruction.

F.3 VISUAL COMPARISON OF INV-SAM GUIDANCE SCALES

In this experiment, we systematically investigate the optimization effects of different guidance scales ω_{sam} in the Inv-Sam Algorithm. 2 on preliminary reconstruction results. As shown in Figure 7, the visual quality of reconstructed images progressively improves with increasing ω_{sam} values. Notably, when $\omega_{\text{sam}}=0$, the reconstruction maintains the structural integrity and content fidelity of the initial recovery without introducing distortion or artifacts. This demonstrates that our algorithm can effectively enhance local details while preserving the fundamental framework of preliminary reconstructions, highlighting its robustness and controllability during detail refinement.

Table 5: GradCFG results using specific text prompts for fine-tuning

Category	GradCFG		+ Inv-Sam			
	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓
Backpack	0.1381	10.81	0.7252	0.2277	11.52	0.6253
Can	0.1626	10.03	0.8099	0.2037	10.62	0.7403
Candle	0.0743	10.28	0.8037	0.1222	11.73	0.5836
Cat	0.1553	11.78	0.7028	0.1840	13.28	0.5539
Sneaker	0.0627	10.35	0.9303	0.1272	12.43	0.6155
Dog	0.1861	10.43	0.7619	0.2583	11.33	0.6520
Monster Toy	0.1134	10.20	0.7913	0.2346	11.28	0.6251
Robot Toy	0.1142	10.61	0.7786	0.2230	11.36	0.6168
Race Car	0.0836	10.38	0.7546	0.1215	11.07	0.5652
Avg. Improv.	0.1213	10.65 Baseline	0.7789	0.1875 +55%	11.51 +8%	0.6226 -20%

Table 6: Semantic similarity comparison between reconstructed text embeddings and original prompts

Original Prompt	Null Text Sim.	Recon. Sim.	Improv. (%)
A backpack sitting on top of a rock with mountains in the background	0.244	0.375	53.7
A red backpack sitting on a tree branch	0.269	0.477	77.2
A woman with a backpack looking up at the sky	0.302	0.351	16.5
A red backpack sitting on the ground in the woods	0.326	0.480	47.5
A red backpack hanging on a tree branch	0.274	0.463	69.1
Avg.	0.283	0.429	52.8

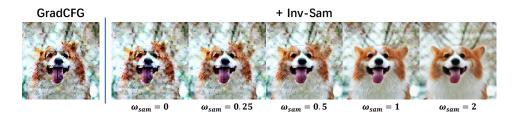


Figure 7: Visual comparison of reconstructed images using different ω_{sam} values. From left to right: Original image, GradCFG baseline and Inv-Sam with $\omega_{\text{sam}}=0,0.25,0.5,1.0,2.0$. Higher guidance scales generally produce sharper details and better semantic alignment with the original prompt.

Table 7: Overall text embedding recovery performance

Null Text Sim.	Recon. Sim.	Improv. (%)
0.3178	0.4536	42.7