
Differentially Private Adaptation of Diffusion Models via Noisy Aggregated Embeddings

Pura Peetathawatchai¹ Wei-Ning Chen² Berivan Isik³ Sanmi Koyejo¹ Albert No⁴

Abstract

We introduce a novel method for adapting diffusion models under differential privacy (DP) constraints, enabling privacy-preserving style and content transfer without fine-tuning model weights. Traditional approaches to private adaptation, such as DP-SGD, incur significant computational and memory overhead when applied to large, complex models. In addition, when adapting to small-scale specialized datasets, DP-SGD incurs large amount of noise that significantly degrades the performance. Our approach instead leverages an embedding-based technique derived from Textual Inversion (TI) and adapted with differentially private mechanisms. We apply TI to Stable Diffusion for style adaptation using two private datasets: a collection of artworks by a single artist and pictograms from the Paris 2024 Olympics. Experimental results show that the TI-based adaptation achieves superior fidelity in style transfer, even under strong privacy guarantees.

1. Introduction

In recent years, diffusion models (Ho et al., 2020; Song et al., 2021b), particularly latent diffusion models (Rombach et al., 2022), have spearheaded high quality text-to-image generation, and have been widely adopted by researchers and the general public alike. Trained on massive datasets like LAION-5B (Schuhmann et al., 2022), these models have developed a broad understanding of visual concepts, enabling new creative and practical applications. Notably, tools such as Stable Diffusion (Rombach et al., 2022; Podell et al., 2023) have been made readily accessible for general use. Building on this foundation, efficient adaptation methods such as parameter efficient fine-tuning (PEFT) methods (Hu et al., 2022; von Platen et al., 2022; Ruiz et al., 2023), guidance based approaches (Ho & Salimans, 2021; Kim et al., 2022; Bansal et al., 2024), and pseudo-word generation (Gal et al., 2023) enable users to leverage this extensive pretraining for customizing models that can specialize on

downstream tasks with smaller datasets.

The rapid adoption of diffusion models has raised significant privacy and legal concerns. These models are vulnerable to privacy attacks, such as membership inference (Duan et al., 2023), where attackers determine if a specific data point was used for training, and data extraction (Carlini et al., 2023), which enables reconstruction of training data. This risk is amplified during fine-tuning on smaller, domain-specific datasets, where each record has a greater impact. Additionally, reliance on large datasets scraped without consent raises copyright concerns (Vyas et al., 2023), as diffusion models can reproduce original artworks without credit or compensation. These issues highlight the urgent need for privacy-preserving technologies and clearer ethical and legal guidelines for generative models.

Differential privacy (DP) (Dwork et al., 2006; Dwork, 2006) is a widely adopted approach to addressing these challenges, where controlled noise is added during training or inference to prevent information leakage from individual data points while still enabling the model to learn effectively from the overall dataset. One standard approach for ensuring DP in deep learning is Differentially Private Stochastic Gradient Descent (DP-SGD) (Abadi et al., 2016), which modifies traditional SGD by adding noise to clipped gradients.

However, applying DP-SGD to train diffusion models poses several challenges. It introduces significant computational and memory overhead due to per-sample gradient clipping (Hoory et al., 2021), which is essential for bounding gradient sensitivity (Dwork et al., 2006; Abadi et al., 2016). DP-SGD is also incompatible with batch-wise operations like batch normalization, as these link samples and hinder sensitivity analysis. Furthermore, training large models with DP-SGD often leads to substantial performance degradation, particularly under realistic privacy budgets since the required noise scales with the gradient norm. Consequently, existing diffusion models trained with DP-SGD are limited to relatively small-scale images. (Dockhorn et al., 2023; Ghalebikesabi et al., 2023).

As a result, recent research has focused on privacy-preserving strategies for adapting diffusion models without the need for full DP-SGD training. One approach adapts large, publicly pre-trained models to new domains under DP

¹Stanford University ²Microsoft ³Google ⁴Yonsei University. Correspondence to: Albert No <albertno@yonsei.ac.kr>.

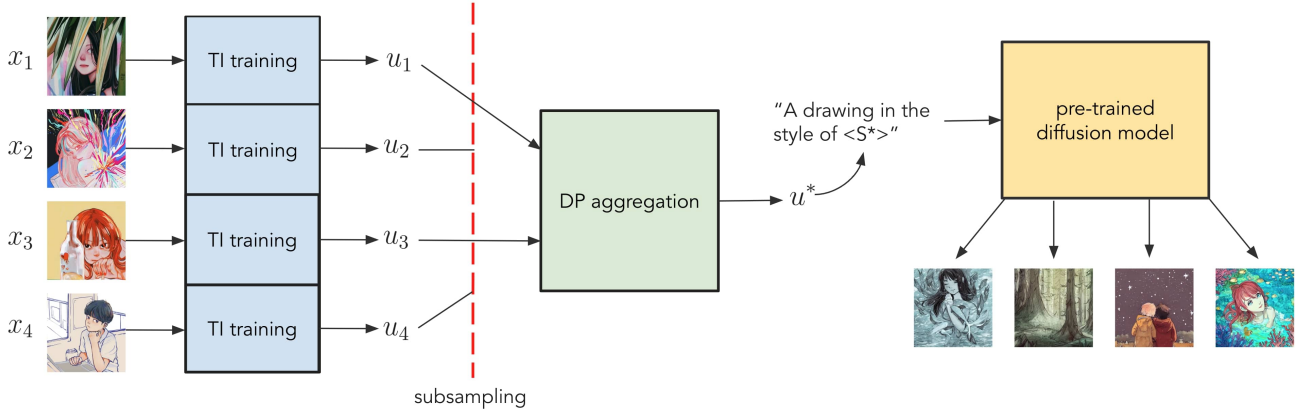


Figure 1. Overview of DPAGg-TI. We first apply Textual Inversion to extract embeddings for each image in the private dataset. These embeddings are then aggregated with differentially private mechanism, incorporating subsampling to produce a private embedding u^* . Finally, images are generated using the corresponding token $\langle S^* \rangle$.

constraints, leveraging their representational strengths while reducing computational and memory costs (Ghalebikesabi et al., 2023). Similarly, PEFT methods like DP-LoRA (Yu et al., 2022) fine-tune a small subset of parameters, enabling efficient adaptation with lower privacy costs. Methods like DP-RDM (Lebensold et al., 2024) avoid direct model updates by using retrieval mechanisms that condition image generation on private data retrieved during inference. However, these alternate approaches often fail to capture the detail of style, underscoring the challenges of balancing privacy, efficiency, and generative performance.

Independent of privacy concerns, Textual Inversion (TI) (Gal et al., 2023) provides an effective method for adapting diffusion models to specific styles or content without modifying the model. Instead, TI learns an external embedding vector that captures the style or content of a target image set, which is then incorporated into text prompts to guide the model’s outputs. A key advantage of TI is its ability to compress a style into a compact vector, reducing computational and memory demands while simplifying the application of privacy-preserving mechanisms, as privacy constraints can be applied directly to embeddings rather than the full model. Additionally, since TI avoids direct model optimization, it remains efficient and compatible with DP constraints on smaller datasets.

In this work, we propose a novel privacy-preserving adaptation method for smaller datasets, leveraging TI to avoid the extensive model updates required by DP-SGD or DP-PEFT approaches. Standard TI inherently compresses the target dataset into a single, low-dimensional vector, providing some obfuscation benefits, but it does not offer formal privacy guarantees. To address this limitation, we introduce a private variant of TI, called Differentially Private Aggregation via Textual Inversion (DPAGg-TI) and summarize

it in Figure 1. Our method decouples interactions among samples by learning a separate embedding for each target image, which are then aggregated into a noisy centroid. This approach ensures efficient and secure adaptation to private datasets.

Our experiments demonstrate the effectiveness of DPAGg-TI, showing that TI remains robust in preserving stylistic fidelity even under privacy constraints. Applying our method to a private artwork collection by @eveismyname and Paris 2024 Olympics pictograms (Paris 2024), we show that DPAGg-TI captures nuanced stylistic elements while ensuring privacy. We observe a trade-off between privacy (controlled by DP parameter ϵ) and image quality: lower ϵ reduces fidelity but maintains the target style under moderate noise. Subsampling further amplifies privacy by reducing sensitivity to individual data points, mitigating noise impact on image quality. This framework enables privacy-preserving adaptation of diffusion models to new styles and domains while protecting sensitive data.

Our contributions can be summarized as follows:

- We propose DPAGg-TI that ensures privacy by learning separate embeddings for individual images and aggregating them into a noisy centroid.
- Our approach enables style adaptation without extensive model updates, reducing computational overhead while preserving privacy.
- We analyze the trade-off between privacy and image quality, showing that moderate noise maintains stylistic fidelity while protecting sensitive data.
- We validate our method on diverse datasets, demonstrating its effectiveness in capturing stylistic elements under privacy constraints.

2. Preliminaries and Related Work

2.1. Diffusion Models

Diffusion models (Ho et al., 2020; Song et al., 2021b;a; Rombach et al., 2022) leverage an iterative denoising process to generate high-quality images that align with a given conditional input from random noise. In text-to-image generation, this conditional input is based on a textual description (a prompt) that guides the model in shaping the image to reflect the content and style specified by the text. To convert the text prompt into a suitable conditional format, it is first broken down into discrete tokens, each representing a word or sub-word unit. These tokens are then converted into a sequence of embedding vectors v_i that encapsulate the meaning of each token within the model’s semantic space. Next, these embeddings pass through a transformer text encoder, such as CLIP (Radford et al., 2021), outputting a single text-conditional vector y that serves as the conditioning input. This vector y is then incorporated at each denoising step, guiding the model to align the output image with the specific details outlined in the prompt.

The image generation process, also known as the reverse diffusion process, comprises of T discrete timesteps and starts with pure Gaussian noise x_T . At each decreasing timestep t , the denoising model, which often utilizes a U-Net structure with cross-attention layers, takes a noisy image x_t and text conditioning y as inputs and predicts the noise component $\epsilon_\theta(x_t, y, t)$, where θ denotes the denoising model’s parameters. The predicted noise is then used to make a reverse diffusion step from x_t to x_{t-1} , iteratively refining the noisy image closer to a coherent output x_0 that aligns with the text conditional y .

The objective function for a text-conditioned diffusion model, given both the noisy image x_t and the text conditioning y , is typically a mean squared error (MSE) between the true noise ϵ and the predicted noise $\epsilon_\theta(x_t, y, t)$. The denoising model is therefore trained over the following optimization problem:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0, I), t \sim [T]} [\|\epsilon - \epsilon_\theta(x_t, y, t)\|^2]. \quad (1)$$

Textual Inversion. Textual Inversion (TI) (Gal et al., 2023) is an adaptation technique that enables personalization using a small dataset of typically 3-5 images. This approach essentially learns a new token that encapsulates the semantic meaning of the training images, allowing the model to associate specific visual features with a custom token.

To achieve this, TI trains a new token embedding, denoted as u , representing a placeholder token, denoted as S . During training, images are conditioned on phrases such as “A photo of S ” or “A painting in the style of S ”. However,

unlike the fixed embeddings of typical tokens v_i , u is a learnable parameter. Let y_u denote the text conditioning vector resulting from a prompt containing the token S . Through gradient descent, TI minimizes the diffusion model loss given in (1) with respect to u , while keeping the diffusion model parameters θ fixed, iteratively refining this embedding to capture the unique characteristics of the training images. The resulting optimal embedding u_* is formalized as follows:

$$u_* = \arg \min_u \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0, I), t \sim [T]} [\|\epsilon - \epsilon_\theta(x_t, y_u, t)\|^2]. \quad (2)$$

Hence, u_* represents an optimized placeholder token S_* , which can be employed in prompts such as “A photo of S_* floating in space” or “A drawing of a capybara in the style of S_* ”, enabling the generation of personalized images that reflect the learned visual characteristics.

2.2. Differential Privacy

In this work, we adopt differential privacy (DP) (Dwork et al., 2006; Dwork, 2006) as our privacy framework. Over the past decade, DP has become the gold standard for privacy protection in both research and industry. It measures the stability of a randomized algorithm with respect to changes in an input instance, thereby quantifying the extent to which an adversary can infer the existence of a specific input based on the algorithm’s output.

Definition 2.1 ((Approximate) Differential Privacy). For $\epsilon, \delta \geq 0$, a randomized mechanism $\mathcal{M} : \mathcal{X}^n \rightarrow \mathcal{Y}$ satisfies (ϵ, δ) -DP if for all neighboring datasets $\mathcal{D}, \mathcal{D}' \in \mathcal{X}^n$ which differ in a single record (i.e., $\|\mathcal{D} - \mathcal{D}'\|_H \leq 1$ where $\|\cdot\|_H$ is the Hamming distance) and all measurable \mathcal{S} in the range of \mathcal{M} , we have that

$$\Pr(\mathcal{M}(\mathcal{D}) \in \mathcal{S}) \leq e^\epsilon \Pr(\mathcal{M}(\mathcal{D}') \in \mathcal{S}) + \delta.$$

When $\delta = 0$, we say \mathcal{M} satisfies ϵ -pure DP or (ϵ) -DP.

To achieve DP, the Gaussian mechanism is often applied (Dwork et al., 2014; Balle & Wang, 2018), adding Gaussian noise scaled by the sensitivity of the function f and privacy parameters ϵ and δ . Specifically, noise with standard deviation $\sigma = \frac{\Delta_f \sqrt{2 \ln(1.25/\delta)}}{\epsilon}$ is added to the output¹ (Balle & Wang, 2018), where Δ_f represents ℓ_2 -sensitivity of the target function $f(\cdot)$. When the context is clear, we may omit the subscript f . This mechanism enables a smooth privacy-utility tradeoff and is widely used in privacy-preserving machine learning, including in DP-SGD (Abadi et al., 2016), which applies Gaussian noise during model updates to achieve DP.

¹In practice, we use numerical privacy accountant such as (Balle & Wang, 2018; Mironov, 2017) to calibrate the noise.

Privacy Amplification by Subsampling. Subsampling is a standard technique in DP, where a full dataset of size n is first subsampled to m records without replacement (typically with $m \ll n$) before the privatization mechanism (such as the Gaussian mechanism) is applied. Specifically, if a mechanism provides (ϵ, δ) -DP on a dataset of size m , it achieves (ϵ', δ') -DP on the subsampled dataset, where $\delta' = \frac{m}{n}\delta$ and

$$\epsilon' = \log \left(1 + \frac{m}{n} (e^\epsilon - 1) \right) = O \left(\frac{m}{n} \epsilon \right). \quad (3)$$

This result is well-known (Steinke (2022, Theorem 29)), with tighter amplification bounds available for Gaussian mechanisms (Mironov, 2017).

2.3. Differentially Private Adaptation of Diffusion Models

Recent advancements in applying DP to diffusion models have aimed to balance privacy preservation with the high utility of generative outputs. Dockhorn et al. (Dockhorn et al., 2023) proposed a Differentially Private Diffusion Model (DPDM) that enables privacy-preserving generation of realistic samples, setting a foundational approach for adapting diffusion processes using DP-SGD. Another common strategy involves training a model on a large public dataset, followed by differentially private fine-tuning on a private dataset, as explored by Ghalebikesabi et al. (2023). While effective in certain contexts, this approach raises privacy concerns, particularly around risks of information leakage during the fine-tuning phase (Tramèr et al., 2024).

In response to these limitations, various adaptation techniques have emerged. Although not specific to diffusion models, some methods focus on training models on synthetic data followed by DP-constrained fine-tuning, as in the VIP approach (Yu et al., 2024), which demonstrates the feasibility of applying DP in later adaptation stages. Other approaches explore differentially private learning of feature representations (Sander et al., 2024), aiming to distill private information into a generalized embedding space while maintaining DP guarantees. Although these adaptations are not yet implemented for diffusion models, they lay essential groundwork for developing secure and efficient privacy-preserving generative models.

3. Differentially Private Adaptation via Textual Inversion

TI is inherently parameter-efficient and offers certain privacy benefits, as information from an entire dataset of images is compressed into a single token embedding vector. This compression limits the model’s capacity to memorize specific images, making data extraction attacks difficult. However, this privacy is merely heuristic and yet to be proven, so TI

may still be vulnerable to privacy attacks such as membership inference. A similar adaptation technique with privacy guarantees may therefore be desirable.

Let $x^{(1)}, \dots, x^{(n)}$ represent a target dataset of images whose characteristics we wish to privately adapt our image generation towards. Instead of training a single token embedding on the entire dataset as in regular TI, we train a separate embedding $u^{(i)}$ on each $x^{(i)}$ to obtain a set of embeddings $u^{(1)}, \dots, u^{(n)}$, as illustrated in Figure 1. We can formalize the encoding process as follows:

$$u^{(i)} = \arg \min_u \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), t} [\|\epsilon - \epsilon_\theta(x_t^{(i)}, y_u, t)\|^2]. \quad (4)$$

Then, we can aggregate the embeddings $u^{(1)}, \dots, u^{(n)}$ by calculating the centroid. The purpose of this aggregation is to limit the sensitivity of the final output to each $x^{(i)}$. In order to provide DP guarantees, we also add isotropic Gaussian noise to the centroid. We can therefore define the resulting embedding vector u^* as follows:

$$u^* = \frac{1}{n} \sum_{i=1}^n u^{(i)} + \mathcal{N}(0, \sigma^2 I), \quad (5)$$

where the minimum σ required to provide (ϵ, δ) -DP is given by the following expression based on Balle & Wang (2018, Theorem 1):

$$\sigma = \frac{\Delta}{n} \cdot \frac{\sqrt{2 \ln(1.25/\delta)}}{\epsilon}. \quad (6)$$

In the context of our problem, $\Delta = \sup_{i,j} \|u^{(i)} - u^{(j)}\|$. Since our embedding vectors are directional, we can normalize each $u^{(i)}$, allowing us to set $\Delta = 2$.

The noisy centroid embedding u^* can then be used to adapt the downstream image generation process. Similar to regular TI’s u_* , we can use u^* to represent a new placeholder token S^* that can be incorporated into prompts for personalized image generation. While u^* may not fully solve the TI optimization problem presented in (2), it provides provable privacy guarantees, with only a minimal trade-off in accurately representing the style of the target dataset.

To reduce the amount of noise needed to provide the same level of DP, we employ subsampling: instead of computing the centroid over all n embedding vectors, we randomly sample $m \leq n$ embedding vectors without replacement and compute the centroid over only the sampled vectors. Then the standard privacy amplification by subsampling bounds (such as (3)) can be applied. Formally, we sample $D_{\text{sub}} \subseteq \{u^{(1)}, \dots, u^{(n)}\}$ where $|D_{\text{sub}}| = m$, and compute the output embedding as follows:

$$u^* = \frac{1}{m} \sum_{u^{(i)} \in D_{\text{sub}}} u^{(i)} + \mathcal{N}(0, \sigma^2 I), \quad (7)$$

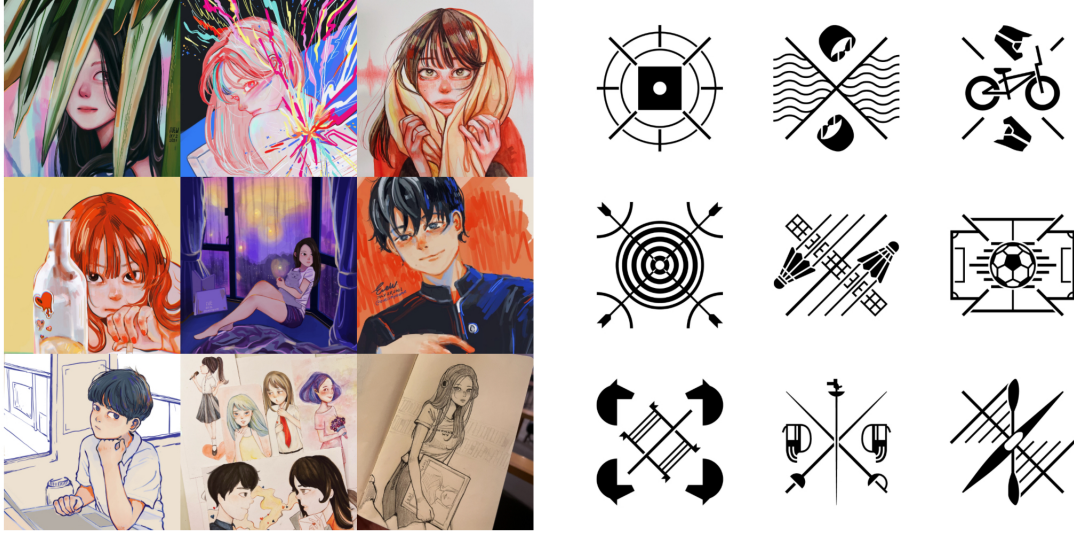


Figure 2. Samples of images used in our style adaptation experiments. **Left:** artwork by @eveismyname ($n = 158$). **Right:** Paris 2024 Olympic pictograms ($n = 47$), © International Olympic Committee, 2023.

where σ can be computed numerically for any target ε, δ and subsampling rate $\frac{m}{n}$.

4. Experimental Results

4.1. Datasets

We compiled two datasets to evaluate our style adaptation method, specifically selecting content unlikely to be recognized by Stable Diffusion v1.5, our base model.

The first dataset consists of 158 artworks by the artist @eveismyname, who has granted consent for non-commercial use. This dataset allows us to assess whether models can capture artistic styles without memorizing individual works. While some of these artworks may have been publicly accessible on social media, making incidental inclusion in Stable Diffusion’s pretraining possible, the artist’s limited recognition and relatively small portfolio reduce the likelihood that the model has internalized her unique style. This dataset serves as a controlled test for privacy-preserving style transfer on individual artistic collections.

The second dataset contains 47 pictograms from the Paris 2024 Olympics (Paris 2024), permitted strictly for non-commercial editorial use (International Olympic Committee). These pictograms were officially released in February 2023, several months after the release of Stable Diffusion v1.5, ensuring they were absent from the model’s pretraining data. This dataset allows us to assess how well our approach adapts to newly introduced visual styles that the base model has never encountered.

Both datasets are used to test the ability of our method to extract and transfer stylistic elements while preserving privacy. Representative samples are shown in Figure 2.

4.2. Style Transfer Results

Using both the @eveismyname and Paris 2024 pictograms dataset, we trained TI (Gal et al., 2023) embeddings on Stable Diffusion v1.5 (Rombach et al., 2022) using DPagg-TI. Our primary goal is to investigate how DP configurations, specifically the privacy budget ε and subsampling size m , affect the generated images quality and privacy resilience. For regular TI, we utilize the default process to embed the private dataset without any additional noise. For the DPagg-TI, we test multiple configurations of m and ε to analyze the trade-off between image fidelity and privacy.

Figures 3 and 4 present generated images across two key configurations: (1) regular TI without DP, (2) DPagg-TI with DP at different values of m and ε . We used the same random seed to generate embeddings, subsample images, and sample DP noise for ease of visual comparison between different configurations. As with common practice, we set $\delta = 1/n$. Since σ is undefined for $\varepsilon = 0$, we demonstrate the results of $\varepsilon \approx 0$, in other words, infinite noise, by setting $\varepsilon = 10^{-5}$. The purpose of this parameter value is to demonstrate the image generated when u^* contains zero information about the target dataset. Images generated without DP closely resemble the unique stylistic elements of the target dataset. In particular, images adapted using @eveismyname images displayed crisp details and nuanced color gradients characteristic of the artist’s work,



Figure 3. Images generated by Stable Diffusion v1.5 using the prompt “A painting of Taylor Swift in the style of <@eveismyname>”, with the embedding <@eveismyname> trained using different values of m and ϵ .

while those of Paris 2024 pictograms captured the logo’s original structure. In contrast, DP configurations introduce a discernible degradation in image quality, with lower epsilon values and smaller subsampling sizes resulting in more noticeable noise and diminished stylistic fidelity.

As $\epsilon \rightarrow 0$, the resulting token embedding u^* gradually loses its semantic meaning, leading to a loss of stylistic fidelity. In particular, y_{u^*} tends towards y (a conditioning vector independent of the learnable embedding). In our results, this manifests as a painting of Taylor Swift devoid of the artist-specific stylistic elements, or a generic icon of a dragon (with color, as opposed to the black and white design of the pictograms). With this in mind, ϵ can be interpreted as a drift parameter, representing the progression from the optimal u^* towards infinity, gradually steering the generated image away from the target style in exchange for stronger privacy guarantees. We also observe instances where there is a temporary drop in prompt fidelity (e.g., $m = 16, \epsilon \in [0.5, 1]$ in Figure 3 and intermediate ϵ values in Figure 4) which restores as u^* drifts even further from its optimal value. We hypothesize that this is due to drifted u^* capturing a different meaning unrelated to the prompt, before losing any meaning that could be interpreted by Stable Diffusion’s text encoder, causing u^* to be disregarded from y_{u^*} and the prompt fidelity to be restored. Another possible explanation is that the temporary drop in prompt fidelity is due to the drift path of u^* passing through non-linear regions within embedding space. We leave further investigations into this observation for future work.

Meanwhile, reducing m also reduces the sensitivity of the generated image to ϵ , as evident by the observation that, on both datasets at $m = 4$, (subsampling rate below 0.1) image generation can tolerate ϵ as low as 0.5 without significant changes in visual characteristics, and retaining stylistic elements of the target dataset at ϵ as low as 0.1. This strong boost in robustness comes at a small price of base style capture fidelity. As observed in Figures 3 and 4, we can also treat subsampling as an introduction of noise. Mathematically, the subsample centroid is an unbiased estimate of the true centroid, and so the subsampling process itself defines a distribution centered at the true centroid. However, the amount of noise introduced by the subsampling process is limited by the individual image embeddings, as a subsample centroid can only stray from the true centroid as much as the biggest outlier in the dataset.

4.3. Quantitative Evaluation

4.3.1. USER STUDY

To evaluate the utility of our approach under different DP and subsampling configurations, we conducted a user study with 25 participants. Each participant was shown reference images from the target dataset and asked to compare pairs of generated images, selecting the one that better captured the style of the reference images. Images were generated using 10 prompts and adapted TI embeddings for the @eveismyname and Paris 2024 Pictogram datasets, resulting in 20 groups of images. Each participant evaluated



Figure 4. Images generated by Stable Diffusion v1.5 using the prompt “Icon of a dragon in the style of <Paris 2024 Pictograms>”, with the embedding <Paris 2024 Pictograms> trained using different values of m and ϵ .

two groups, one randomly selected from each dataset, with comparisons focusing on model configurations differing by DP noise and subsampling size.

Survey results, summarized in Table 3 in Appendix A, align with our design goals. Participants showed no clear preference between regular TI and DPagg-TI, suggesting that our privacy-preserving approach maintains perceptual quality. As expected, both DP noise and reduced subsampling size degraded style fidelity, consistent with the trade-offs inherent in differential privacy. Preferences at $\epsilon = 1$ were split, but subsampling was generally favored, reinforcing its role in reducing noise impact while preserving style.

4.3.2. KERNEL INCEPTION DISTANCE

The Kernel Inception Distance (KID) (Bińkowski et al., 2018) is a metric for evaluating generative models by measuring the difference between the distributions of generated and training images in an embedding space. To compute KID, images generated by the model and real training images are passed through an Inception network (Szegedy et al., 2015), and their distributional differences are estimated. Unlike the more commonly used Fréchet Inception Distance (FID) (Heusel et al., 2017), KID is an unbiased estimator of the true divergence between the learned and target distributions (Jayasumana et al., 2024), making it more suitable for smaller datasets, as in our case.

We report KID scores for different parameters in Tables 1 and 2, showing that DPagg-TI maintains the style transfer fidelity of TI while ensuring differential privacy. Further

discussion of these results is provided in Appendix B.

4.4. Ablation Study

4.4.1. TEXTUAL INVERSION WITH DP-SGD

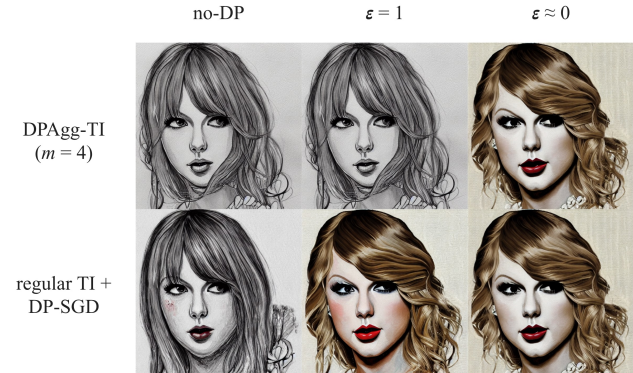


Figure 5. Comparing our approach to applying DP-SGD to regular TI using prompt “a painting of Taylor Swift in the style of @eveismyname”.

A natural question that arises is how well our approach compares to the naive method of applying DP-SGD to regular TI training. We therefore integrated DP-SGD into the TI codebase using the Opacus library and trained similar embeddings on the @eveismyname and Paris 2024 datasets. We found that in most cases, notably the @eveismyname dataset, the amount of noise required for DP-SGD to achieve a reasonable value of ϵ for DP is so high that the resulting

Table 1. KID scores of DPAGg-TI on @eveismyname dataset for various ϵ values ranging from $\epsilon = 10^{-5}, 0.1, 0.5, 1.0, 2.0, 5.0$ (including no DP) under different subsampling levels ($m = 4, 8, 16, 32$) as well as regular TI (ctrl).

m	No DP	$\epsilon = 5.0$	$\epsilon = 2.0$	$\epsilon = 1.0$	$\epsilon = 0.5$	$\epsilon = 0.1$	$\epsilon \approx 0$
–	0.0444	0.0794	0.0422	0.0529	0.0690	0.1117	0.0654
32	0.0752	0.0845	0.0865	0.1167	0.0300	0.0649	0.0657
16	0.0351	0.0379	0.0430	0.0663	0.1309	0.0438	0.0658
8	0.0359	0.0366	0.0350	0.0366	0.0396	0.0530	0.0658
4	0.0245	0.0250	0.0249	0.0250	0.0258	0.0314	0.0653
ctrl	0.0318	–	–	–	–	–	–

Table 2. KID scores of DPAGg-TI on Paris dataset for various ϵ values ranging from $\epsilon = 1e - 5, 0.1, 0.5, 1.0, 2.0, 5.0$ (including no DP) under different subsampling levels ($m = 4, 8, 16, 32$) as well as regular TI (ctrl).

m	No DP	$\epsilon = 5.0$	$\epsilon = 2.0$	$\epsilon = 1.0$	$\epsilon = 0.5$	$\epsilon = 0.1$	$\epsilon \approx 0$
–	0.1146	0.1202	0.1368	0.1314	0.1389	0.1209	0.1274
32	0.1220	0.1036	0.1258	0.1377	0.1307	0.1245	0.1259
16	0.1311	0.1424	0.1170	0.1311	0.1381	0.1335	0.1278
8	0.1317	0.1307	0.1220	0.1117	0.1295	0.1313	0.1272
4	0.1141	0.1094	0.1137	0.1190	0.1194	0.1583	0.1259
ctrl	0.1388	–	–	–	–	–	–

embedding contains negligible information about the training dataset. In particular, the results for $\epsilon = 1$ are almost indistinguishable to $\epsilon \approx 0$, as shown in Figure 5. We believe that this is simply because DP-SGD is not designed to handle such small datasets in the order of 100 images. Additional results can be found in Appendix D.

4.4.2. DIFFERENTIALLY PRIVATE ADAPTATION USING STYLE GUIDANCE



Figure 6. Attempts of using universal guidance to generate drawings of Taylor Swift and icons of the Eiffel Tower in the styles of @eveismyname and Paris 2024 Pictograms respectively. Here, we apply no subsampling or DP-noise.

We extend our approach to style guidance (SG) by leveraging the framework of Universal Guidance (Bansal et al., 2024). Specifically, we focus on CLIP-based style guidance, which optimizes the similarity between the CLIP embeddings of a target image and the generated image.

We encode each target image $x^{(i)}$ as $u^{(i)}$ via a CLIP image encoder, then aggregate the embeddings $u^{(1)}, \dots, u^{(n)}$ into u^* using (5) or (7), depending on whether subsampling is applied. The aggregated embedding u^* is then incorporated into the reverse diffusion process as a style guide. Further implementation details are provided in Appendix C.

We apply our SG-based approach to both datasets. While it provides privacy protection by obfuscating embedding details, the resulting images captured only generalized stylistic elements and lack the detailed fidelity and coherence achieved with the TI-based method. As shown in Figure 6, this highlights the superiority of TI in balancing privacy and high-quality image generation.

The reduced effectiveness of SG for style transfer may stem from its sensitivity to hyperparameters such as the guidance weight w , leading to instability. Although Bansal et al. (2024) proposed remedies, namely backward guidance and per-step self-recurrence, these proved insufficient for our application. Additionally, the CLIP embeddings may not retain enough stylistic detail after the aggregation.

5. Conclusion

We presented a differentially private adaptation method for diffusion models using Textual Inversion for privacy-preserving style transfer. Experiments on private artwork and Paris 2024 pictograms showed TI preserves stylistic fidelity and outperforms Style Guidance. Our results demonstrate embedding-driven methods as efficient, scalable alternatives to DP-SGD, balancing style quality and privacy.

Impact Statement

The use of images without owner consent raises significant ethical concerns, particularly regarding the exploitation of intellectual property. This work introduces a method for visual generative models to adapt to new styles and classes while ensuring privacy and copyright protection for data owners. By providing a framework for privacy-preserving adaptation, this technology aims to respect intellectual property and address ethical challenges in generative AI. While it does not eliminate the need for consent from data owners, we hope that it represents a step toward balancing innovation with ethical considerations in AI development. Beyond creative applications, the proposed method has broader potential uses, including synthetic data generation, privacy-preserving personalization, and fine-tuning diffusion models for private or domain-specific tasks.

Acknowledgements

We sincerely thank Tatchamon Wongworakul (@eveismyname) for providing her artwork for use in this study. We are also grateful to Anwar Hithnawi and Varun Chandrasekaran for their insightful discussions and feedback, as well as to all participants in our user study. Sanmi Koyejo acknowledges support by NSF 2046795 and 2205329, IES R305C240046, the MacArthur Foundation, Stanford HAI, OpenAI, and Google.

References

- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. In *ACM SIGSAC*, 2016.
- Balle, B. and Wang, Y.-X. Improving the gaussian mechanism for differential privacy: Analytical calibration and optimal denoising. In *ICML*, 2018.
- Bansal, A., Chu, H.-M., Schwarzschild, A., Sengupta, S., Goldblum, M., Geiping, J., and Goldstein, T. Universal guidance for diffusion models. In *ICLR*, 2024.
- Bińkowski, M., Sutherland, D. J., Arbel, M., and Gretton, A. Demystifying mmd gans, 2018.
- Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwag, V., Tramèr, F., Balle, B., Ippolito, D., and Wallace, E. Extracting training data from diffusion models. In *USENIX Security*, 2023.
- Dockhorn, T., Cao, T., Vahdat, A., and Kreis, K. Differentially private diffusion models. *TMLR*, 2023.
- Duan, J., Kong, F., Wang, S., Shi, X., and Xu, K. Are diffusion models vulnerable to membership inference attacks? In *ICML*, 2023.
- Dwork, C. Differential privacy. In *ICALP*, 2006.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *TCC*, 2006.
- Dwork, C., Roth, A., et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A. H., Chechik, G., and Cohen-or, D. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2023.
- Ghalebikesabi, S., Berrada, L., Goyal, S., Ktena, I., Stanforth, R., Hayes, J., De, S., Smith, S. L., Wiles, O., and Balle, B. Differentially private diffusion models generate useful synthetic images. *arXiv preprint arXiv:2302.13861*, 2023.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/8a1d694707eb0fefe65871369074926d-Paper.pdf.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. In *NeurIPS Workshop on Deep Generative Models and Downstream Applications*, 2021.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *NeurIPS*, 2020.
- Hoory, S., Feder, A., Tendler, A., Erell, S., Peled-Cohen, A., Laish, I., Nakhost, H., Stemmer, U., Benjamini, A., Hassidim, A., et al. Learning and evaluating a differentially private pre-trained language model. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 1178–1189, 2021.
- Hu, E. J., yelong shen, Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022.
- Innat. Van gogh paintings. <https://www.kaggle.com/datasets/ipythonx/van-gogh-paintings>.
- International Olympic Committee. Olympic properties. <https://olympics.com/ioc/olympic-properties>.

- Jayasumana, S., Ramalingam, S., Veit, A., Glasner, D., Chakrabarti, A., and Kumar, S. Rethinking fid: Towards a better evaluation metric for image generation. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9307–9315. IEEE, June 2024. doi: 10.1109/cvpr52733.2024.00889. URL <http://dx.doi.org/10.1109/CVPR52733.2024.00889>.
- Kim, G., Kwon, T., and Ye, J. C. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *CVPR*, 2022.
- Lebensold, J., Sanjabi, M., Astolfi, P., Romero-Soriano, A., Chaudhuri, K., Rabbat, M., and Guo, C. Dp-rdm: Adapting diffusion models to private domains without fine-tuning. *arXiv preprint arXiv:2403.14421*, 2024.
- Mironov, I. Rényi differential privacy. In *CSF*, 2017.
- Paris 2024. Paris 2024 - pictograms. <https://olympics.com/en/paris-2024/the-games/the-brand/pictograms>.
- Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., and Aberman, K. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023.
- Sander, T., Yu, Y., Sanjabi, M., Durmus, A. O., Ma, Y., Chaudhuri, K., and Guo, C. Differentially private representation learning via image captioning. In *ICML*, 2024.
- Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 2022.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. In *ICLR*, 2021a.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021b.
- Steinke, T. Composition of differential privacy & privacy amplification by subsampling. *arXiv preprint arXiv:2210.00597*, 2022.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision, 2015. URL <https://arxiv.org/abs/1512.00567>.
- Tramèr, F., Kamath, G., and Carlini, N. Position: Considerations for differentially private learning with large-scale public pretraining. In *ICML*, 2024.
- von Platen, P., Patil, S., Lozhkov, A., Cuenca, P., Lambert, N., Rasul, K., Davaadorj, M., Nair, D., Paul, S., Berman, W., Xu, Y., Liu, S., and Wolf, T. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022.
- Vyas, N., Kakade, S. M., and Barak, B. On provable copyright protection for generative models. In *ICML*, 2023.
- Yu, D., Naik, S., Backurs, A., Gopi, S., Inan, H. A., Kamath, G., Kulkarni, J., Lee, Y. T., Manoel, A., Wutschitz, L., et al. Differentially private fine-tuning of language models. In *ICLR*, 2022.
- Yu, Y., Sanjabi, M., Ma, Y., Chaudhuri, K., and Guo, C. Vip: A differentially private foundation model for computer vision. In *ICML*, 2024.

A. User Study

A.1. Study Design and Objective

The user study aimed to assess the utility of our approach under different DP and subsampling configurations by evaluating the models' ability to adapt to novel styles. The study involved 25 participants, each of whom was tasked with comparing images generated using various configurations and selecting the one that better captured the style of reference images.

A.2. Experimental Setup

Participants were shown reference images from two datasets:

- The [@eveismyname](#) dataset of private artwork.
- The Paris 2024 Pictogram dataset.

For each dataset, 10 prompts were used to generate images, resulting in 20 groups of images (10 prompts per dataset). Each group included images generated using the same prompt and dataset but with different model configurations. Configurations varied in the addition of DP noise and the size of subsampling.

- Original Textual Inversion (TI)
- DPAGg-TI ($\epsilon = \infty$, no DP) w/o subsampling
- DPAGg-TI ($\epsilon = 1$) without subsampling
- No Adaptation
- DPAGg-TI ($\epsilon = \infty$, no DP) with subsampling ($m = 8$)
- DPAGg-TI ($\epsilon = 1$) with subsampling ($m = 8$)
- Style Guidance (SG)

A.3. Survey Procedure

Participants were asked to evaluate two groups of images: one randomly selected from the [@eveismyname](#) dataset and one from the Paris 2024 Pictogram dataset. For each group:

1. Participants were shown reference images from the target dataset.
2. They were presented with pairs of images generated using different model configurations for the same prompt.
3. Participants selected the image they felt better captured the style of the reference images.

A.4. Evaluation Metrics

The study focused on assessing:

- Participants' preference between regular TI and DPAGg-TI for style adaptation.
- The impact of DP noise and subsampling size on the perceived utility of style transfer.

A.5. Results and Analysis

The results are summarized in Table 3. Key observations include:

- Participants showed no clear preference between regular TI and DPAGg-TI in capturing styles for either dataset.

- Both DP noise and reduced subsampling size decreased the perceived quality of style transfer.
- Preferences were split between configurations with $\epsilon = 1$ with and without subsampling, though subsampling generally had favorable outcomes.

These findings highlight the trade-off between increased DP robustness and reduced utility, suggesting that the optimal configuration may depend on subjective preferences and specific application requirements.

Table 3. Survey Results.

	regular TI	No Adaptation	Unsure
@eveismyname	19	4	2
Paris 2024	16	6	3

	DPAgg-TI (no DP, no subsampling)	No Adaptation	Unsure
@eveismyname	16	9	0
Paris 2024	15	4	6

	regular TI	DPAgg-TI (no DP, no subsamp.)	Unsure
@eveismyname	12	13	0
Paris 2024	9	10	6

	regular TI	DPAgg-TI (no DP, subsamp. $m = 8$)	Unsure
@eveismyname	16	6	3
Paris 2024	7	13	5

	DPAgg-TI (no DP, no subsampling)	DPAgg-TI (no DP, subsamp. $m = 8$)	Unsure
@eveismyname	18	4	3
Paris 2024	10	8	7

	DPAgg-TI ($\epsilon = 1$) no subsampling	DPAgg-TI ($\epsilon = 1$, subsamp. $m = 8$)	Unsure
@eveismyname	14	10	1
Paris 2024	3	16	6

	DPAgg-TI (no DP, no subsampling)	Style Guidance	Unsure
@eveismyname	16	8	1
Paris 2024	20	2	3

	DPAgg-TI ($\epsilon = 1$, subsamp. $m = 8$)	Style Guidance	Unsure
@eveismyname	16	8	1
Paris 2024	19	2	4

	DPAgg-TI (no DP, subsamp. $m = 8$)	DPAgg-TI ($\epsilon = 1$, subsamp. $m = 8$)	Unsure
@eveismyname	8	5	12
Paris 2024	15	4	6

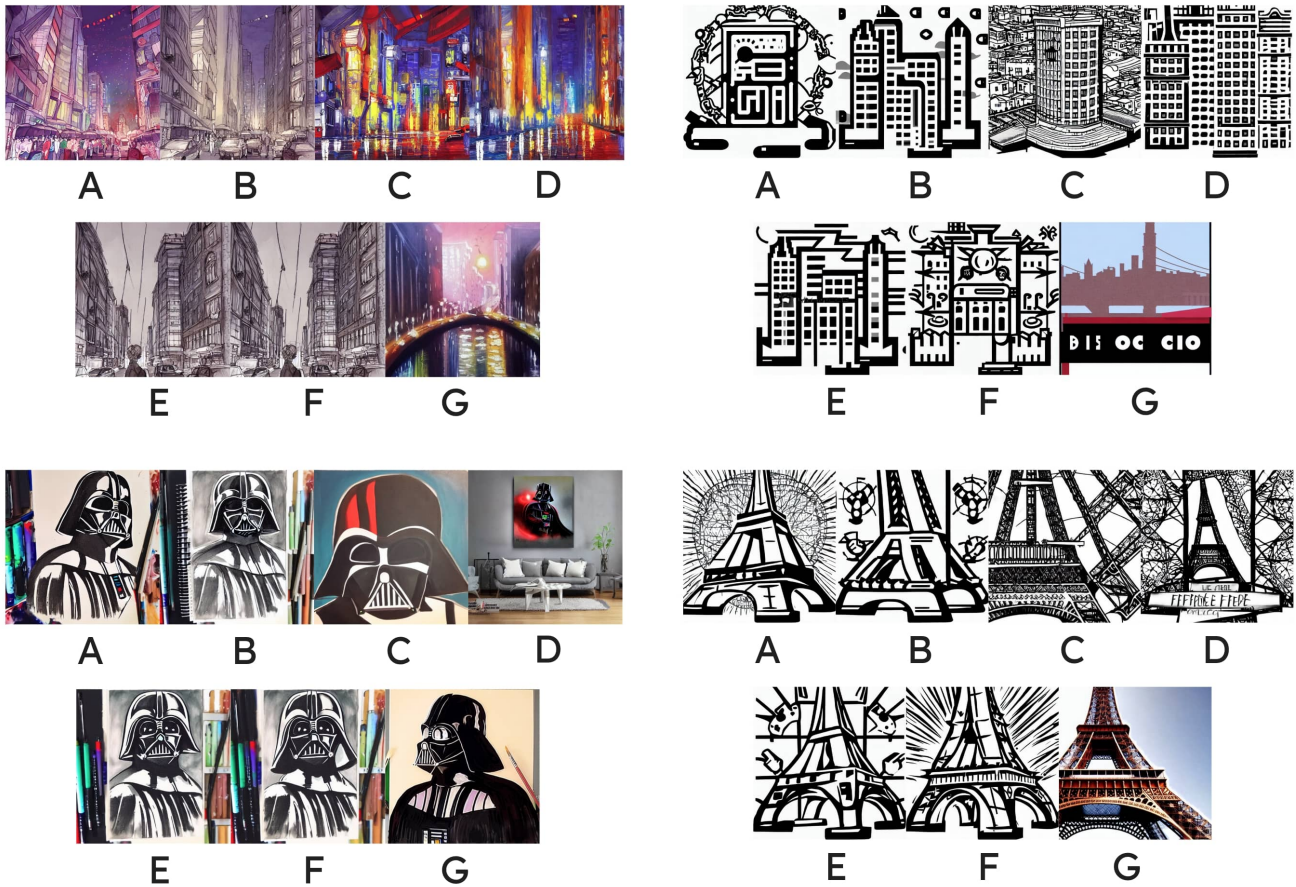


Figure 7. Samples of image sets used in our user study. Participants are asked to compare 2 images at a time.

B. Kernel Inception Distance Discussion

Our results indicate that DP_{Agg}-TI preserves the style transfer fidelity of TI while also ensuring differential privacy. Notably, for `@eveismyname` ($m = 4$) at low privacy budgets, we observe even lower KID values than standard TI, suggesting enhanced style alignment. Similarly, results for the Paris 2024 dataset follow a comparable trend, with DP_{Agg}-TI achieving KID scores similar to TI at low privacy budgets. However, the overall KID scores for this dataset remain high within the context of diffusion model style transfer.

Upon inspecting the generated images (Figure 8), we hypothesize that the abstract and out-of-distribution nature of the Paris 2024 images poses a challenge for the Inception network, leading to less meaningful feature embeddings. This likely inflates the measured embedding distances between generated and reference images, resulting in higher-than-expected KID values.

For KID evaluations, we used prompts similar to those employed during TI training: “A painting/icon in the style of S^* ”. Consistent with the training image captions, these prompts do not specify a subject.



Figure 8. Sample of generated images for KID evaluations with respect to the Paris 2024 dataset.

C. Style Guidance

C.1. Background: Denoising Diffusion Implicit Models

Denoising Diffusion Implicit Models (DDIM) sampling (Song et al., 2021a) uses the predicted noise $\epsilon_\theta(x_t, y, t)$ and a noise schedule represented by an array of scalars $\{\alpha_t\}_{t=1}^T$ to first predict a clean image \hat{x}_0 , then makes a small step in the direction of \hat{x}_0 to obtain x_{t-1} . The reverse diffusion process for DDIM sampling can be formalized as follows:

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t, y, t)}{\sqrt{\alpha_t}} \quad (8)$$

$$x_{t-1} = \sqrt{\alpha_{t-1}} \hat{x}_0 + \sqrt{1 - \alpha_{t-1}} \epsilon_\theta(x_t, y, t). \quad (9)$$

C.2. Implementation

We follow the style guidance process introduced by Bansal et al. (2024), modifying it to include differential privacy mechanisms. Let x_c denote the target style image, x_t the noisy image at step t , and $\mathcal{E}(\cdot)$ the CLIP image encoder. The forward guidance process is defined as follows:

$$\hat{\epsilon}_\theta(x_t, y, t) = \epsilon_\theta(x_t, y, t) + w \sqrt{1 - \alpha_t} \nabla_{x_t} \ell_{\cos}(\mathcal{E}(x_c), \mathcal{E}(\hat{x}_0)), \quad (10)$$

where w is a guidance weight and ℓ_{\cos} is the negative cosine similarity loss. For a detailed description of Universal Guidance, including the backward guidance process and per-step self-recurrence, we refer the reader to the original paper. The reverse diffusion step replaces $\epsilon_\theta(x_t, y, t)$ with $\hat{\epsilon}_\theta(x_t, y, t)$, generating an image x_0 that aligns with the text conditioning y while incorporating the stylistic characteristics of x_c .

To integrate differential privacy, we encode each target image $x^{(i)}$ into $u^{(i)} = \mathcal{E}(x^{(i)})$ and aggregate these embeddings into u^* using the centroid method. The aggregated u^* guides the reverse diffusion process:

$$\hat{\epsilon}_\theta(x_t, y, t) = \epsilon_\theta(x_t, y, t) + w \sqrt{1 - \alpha_t} \nabla_{x_t} \ell_{\cos}(u^*, \mathcal{E}(\hat{x}_0)). \quad (11)$$

This ensures privacy-preserving style transfer while maintaining high stylistic fidelity.

C.3. Ablation



Figure 9. Sample of paintings by Van Gogh used to generate style guidance embeddings.

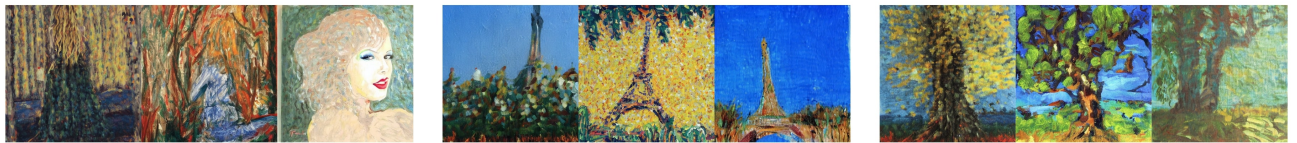


Figure 10. Images generated by Stable Diffusion v1.5 with style guidance towards Van Gogh’s *Saint-Paul Asylum*, *Saint-Rémy* collection using prompts “A painting of Taylor Swift (left) / the Eiffel Tower (center) / a tree (right)”.

To better understand the limited effectiveness of style guidance in our experiments, despite its success in (Bansal et al., 2024), we applied our approach to a dataset of 143 paintings from Van Gogh’s *Saint-Paul Asylum*, *Saint-Rémy* collection (Innat)

(Figure 9). Unlike the `@eveismyname` and Paris 2024 datasets, it is highly likely that Stable Diffusion has been trained on these images. Additionally, Bansal et al. (2024) demonstrated successful adaptation towards the style of Van Gogh’s Starry Night as a single reference image, making this dataset a reasonable interpolation between their successful results and our more limited findings.

Without DP noise or subsampling, we obtained reasonable style transfer results, as shown in Figure 10. This suggests that style guidance struggles when applied to previously unseen target styles, and that its effectiveness may depend on prior exposure within the pre-training data.

D. Additional Style Transfer and Ablation Results



Figure 11. Images generated by Stable Diffusion v1.5 using the prompt “A painting of Taylor Swift in the style of `<@eveismyname>`”, with the embedding `<@eveismyname>` trained using DPAgg-TI (with different subsample sizes m) and TI with DP-SGD using different values of ϵ .

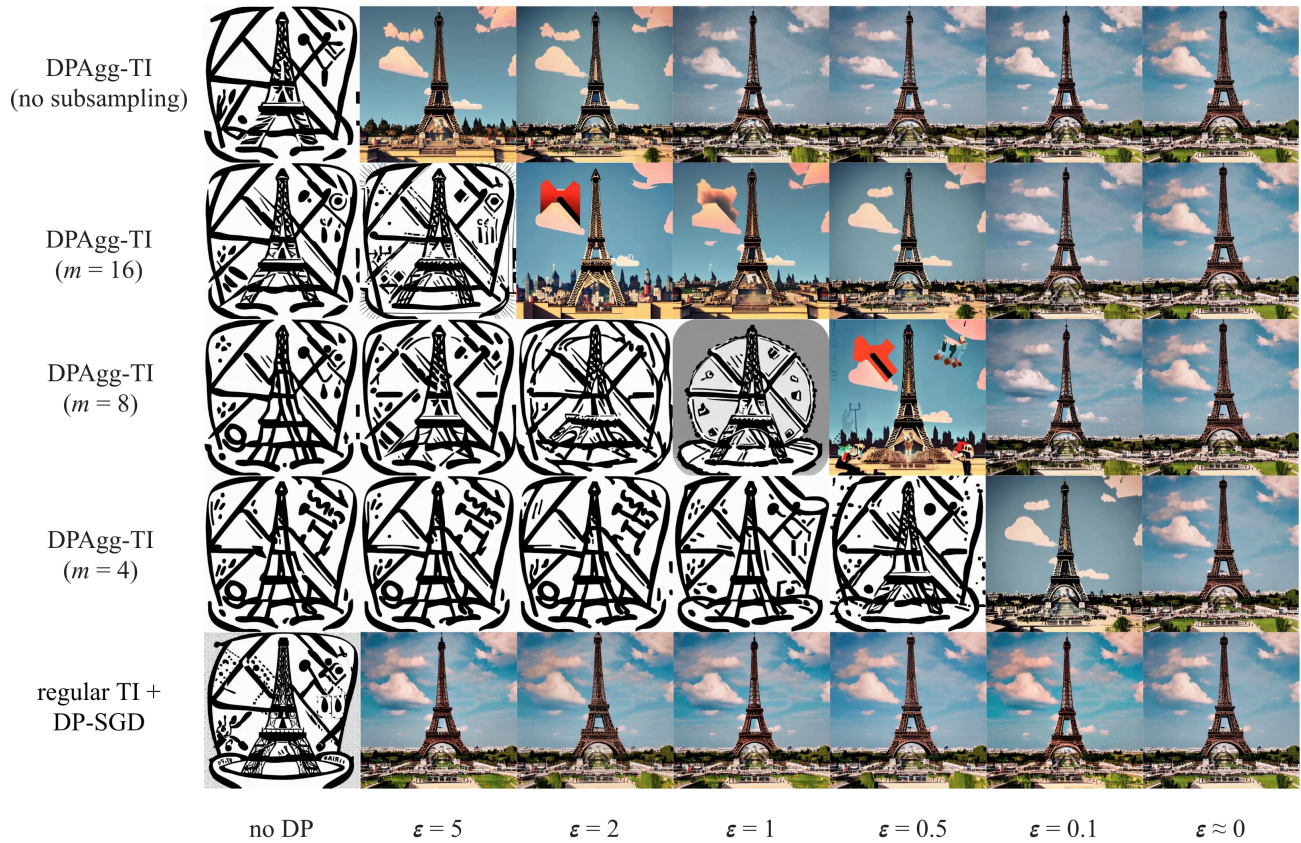


Figure 12. Images generated by Stable Diffusion v1.5 using the prompt “An icon of the Eiffel Tower in the style of <Paris 2024 Pictograms>”, with the embedding <Paris 2024 Pictograms> trained using DPAgg-TI (with different subsample sizes m) and TI with DP-SGD using different values of ϵ .



Figure 13. Images generated by Stable Diffusion v1.5 using the prompt “An icon of a dragon in the style of <Paris 2024 Pictograms>”, with the embedding <Paris 2024 Pictograms> trained using DPAgg-TI (with different subsample sizes m) and TI with DP-SGD using different values of ϵ .