

Differentially Private Adaptation of Diffusion Models via Noisy Aggregated Embeddings

Pura Peetathawatchai^{1,2†}, Wei-Ning Chen³, Berivan Isik⁴, Sanmi Koyejo¹ and Albert No⁵

¹Stanford University, ²ETH Zurich, ³Microsoft, ⁴Google DeepMind, ⁵Yonsei University

albertno@yonsei.ac.kr

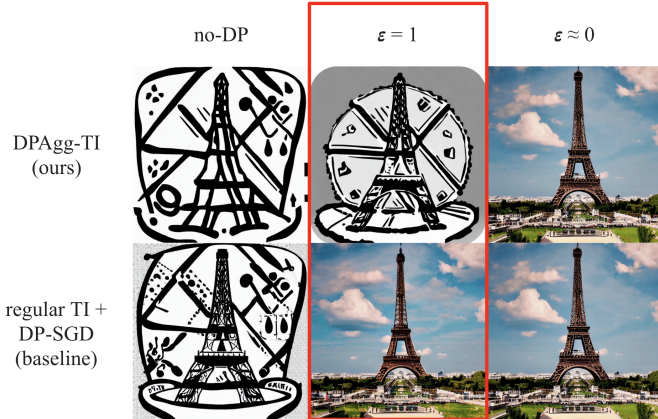


Fig. 1. We compare our method (DPAgg-TI, top) to a baseline applying DP-SGD to Textual Inversion (bottom), using the prompt “an icon of the Eiffel Tower in the style of the Paris 2024 Olympic Pictograms.” While the baseline learns a single embedding over the dataset, our method privately aggregates per-image embeddings. At privacy budget $\epsilon = 1$, DPAgg-TI preserves visual fidelity much better than the baseline, and closely matches the non-private output (left), demonstrating a superior privacy-utility tradeoff.

Abstract—Personalizing large-scale diffusion models poses serious privacy risks, especially when adapting to small, sensitive datasets. A common approach is to fine-tune the model using differentially private stochastic gradient descent (DP-SGD), but this suffers from severe utility degradation due to the high noise needed for privacy, particularly in the small data regime. We propose an alternative that leverages Textual Inversion (TI), which learns an embedding vector for an image or set of images, to enable adaptation under differential privacy (DP) constraints. Our approach, Differentially Private Aggregation via Textual Inversion (DPAgg-TI), adds calibrated noise to the aggregation of per-image embeddings to ensure formal DP guarantees while preserving high output fidelity. We show that DPAgg-TI outperforms DP-SGD finetuning in both utility and robustness under the same privacy budget, achieving results closely matching the non-private baseline on style adaptation tasks using private artwork from a single artist and Paris 2024 Olympic pictograms. In contrast, DP-SGD fails to generate meaningful outputs in this setting.

I. INTRODUCTION

The rapid adoption of diffusion models [1]–[3] has raised significant privacy and legal concerns. These models are vulnerable to privacy attacks, such as membership inference [4], where attackers determine if a specific data point was used for training, and data extraction [5], 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 [6], 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) [7] is a widely adopted framework for addressing these challenges. One standard approach for ensuring DP in deep learning is Differentially Private Stochastic Gradient Descent (DP-SGD) [8], 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 [9], which is essential for bounding gradient sensitivity [8], [10]. 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 small-scale images [11], [12].

Independent of privacy concerns, Textual Inversion (TI) [13] effectively adapts 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 privacy 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. Standard TI does not offer formal privacy guarantees, so we introduce a private variant, Differentially Private Aggregation via Textual Inversion (DPAgg-TI), summarized in Figure 2. 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.

[†]Work done while at Stanford University.

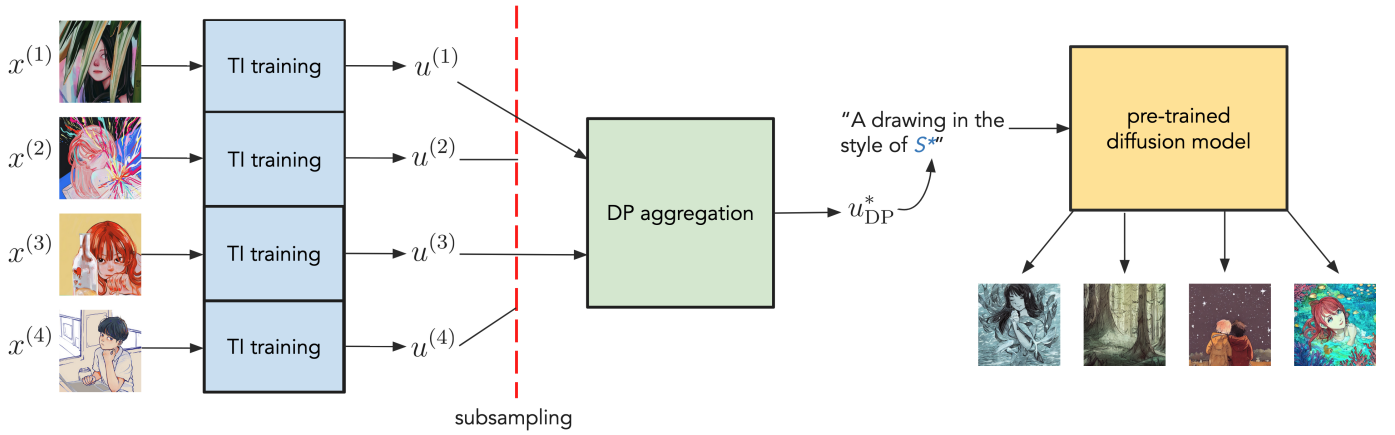


Fig. 2. Overview of DPagg-TI. We first apply Textual Inversion to extract embeddings for each image in the private dataset. These embeddings are then aggregated with differentially private mechanism, incorporating subsampling to produce a private embedding u_{DP}^* . Finally, images are generated using the corresponding token S^* .

Our experiments demonstrate the effectiveness of DPagg-TI, showing that TI remains robust in preserving stylistic fidelity even under privacy constraints (Figure 1). Applying our method to a private artwork collection by @eveismyname and Paris 2024 Olympics pictograms [14], 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.

II. BACKGROUND AND RELATED WORK

A. Diffusion Models

Diffusion models [1]–[3], [15] leverage an iterative denoising process to generate high-quality images that align with a given conditioning input from random noise. In text-to-image generation, this conditioning 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 conditioning 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 [16], outputting a single text-conditioning 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 conditioned on 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 between the true noise ϵ and the predicted noise $\epsilon_\theta(x_t, y, t)$. The denoising model is therefore trained over the optimization problem

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

B. Textual Inversion

Textual Inversion (TI) [13] 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 with respect to u , instead of the diffusion model parameters θ , which we keep fixed. By doing so, we iteratively refine this embedding to capture the unique characteristics of the training images. The resulting optimal embedding u^* is formalized as

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

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.

C. Differential Privacy

In this work, we adopt differential privacy (DP) [7], [10] 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 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*

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

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

To achieve DP, the Gaussian mechanism is commonly applied [17], [18], adding Gaussian noise scaled by the ℓ_2 -sensitivity Δ and privacy parameters (ϵ, δ) . We add zero-mean isotropic Gaussian noise with standard deviation

$$\sigma = \frac{\Delta \sqrt{2 \ln(1.25/\delta)}}{\epsilon}. \quad (3)$$

In practice, we calibrate σ using numerical privacy accountants (e.g., the analytic Gaussian mechanism and RDP) [18], [19]. This mechanism enables a smooth privacy-utility trade-off and is widely used in privacy-preserving machine learning, including DP-SGD [8], which adds Gaussian noise to model updates to achieve DP.

1) *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 (e.g. 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

$$\epsilon' = \log \left(1 + \frac{m}{n} (e^\epsilon - 1) \right) = O \left(\frac{m}{n} \epsilon \right). \quad (4)$$

This result is well-known [20, Theorem 29], with tighter amplification bounds available for Gaussian mechanisms [19].

D. 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. Early work on differentially private generative models, such as Chen et al.’s [21] investigation of DP-GANs with model inversion defenses, established foundational principles for protecting generative models from privacy breaches during training. Dockhorn et al. [11] 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 [12]. While effective in certain contexts, this approach raises privacy concerns, particularly around risks of information leakage during the fine-tuning phase [22].

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 Yu et al. [23], which demonstrates the feasibility of applying DP in later adaptation stages. Other approaches explore differentially private learning of feature representations [24], 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.

III. DIFFERENTIALLY PRIVATE ADAPTATION VIA TEXTUAL INVERSION

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 2. We can formalize the encoding process as

$$u^{(i)} = \arg \min_u \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), t} [\|\epsilon - \epsilon_\theta(x_t^{(i)}, y_u, t)\|^2]. \quad (5)$$

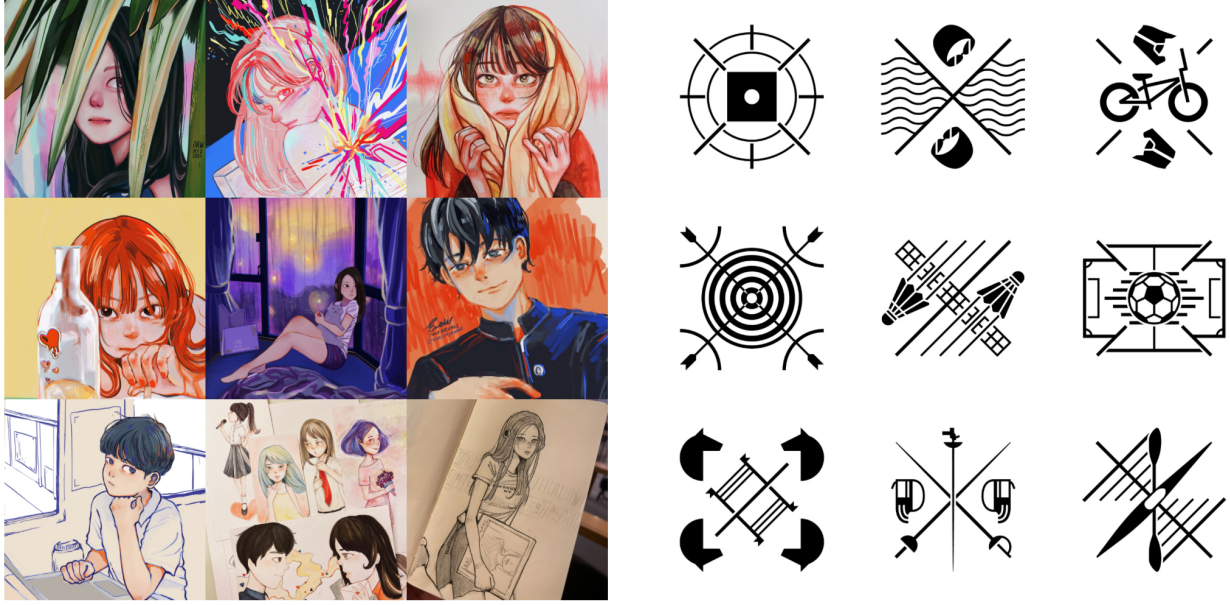


Fig. 3. Samples of images used in our style adaptation experiments. **Left:** artwork by @eveismynname ($n = 158$). **Right:** Paris 2024 Olympic pictograms ($n = 47$), © International Olympic Committee, 2023.

Then, we aggregate the embeddings $u^{(1)}, \dots, u^{(n)}$ by calculating their centroid and adding isotropic Gaussian noise. To ensure bounded sensitivity, we employ a purely directional token embedding (semantics depend only on direction), such as CLIP [16], and ℓ_2 -normalize each embedding vector prior to aggregation. We can therefore define the resulting centroid u_{DP}^* as

$$u_{\text{DP}}^* = \frac{1}{n} \sum_{i=1}^n \frac{u^{(i)}}{\|u^{(i)}\|} + \mathcal{N}(0, \sigma^2 I), \quad (6)$$

where σ is given by (3).

Under this normalization, the ℓ_2 -sensitivity of u_{DP}^* is

$$\Delta = \frac{1}{n} \sup_{u, u'} \left\| \frac{u}{\|u\|} - \frac{u'}{\|u'\|} \right\| = \frac{2}{n}, \quad (7)$$

The noisy centroid embedding u_{DP}^* can then be used to adapt the downstream image generation process. Similar to regular TI's u^* , we can use u_{DP}^* to represent a new placeholder token S^* that can be incorporated into prompts for personalized image generation. While u_{DP}^* 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 (4))

can be applied. Formally, we sample $D_{\text{sub}} \subseteq \{u^{(1)}, \dots, u^{(n)}\}$ where $|D_{\text{sub}}| = m$, and compute u_{DP}^* as

$$u_{\text{DP}}^* = \frac{1}{m} \sum_{u^{(i)} \in D_{\text{sub}}} \frac{u^{(i)}}{\|u^{(i)}\|} + \mathcal{N}(0, \sigma^2 I), \quad (8)$$

where σ can be computed numerically using (3) with ε' and δ' from Section II-C1.

IV. EXPERIMENTAL RESULTS

A. 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 @eveismynname, 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 [14], permitted strictly for non-commercial editorial use [25]. 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 3.

B. Style Transfer Results

Using both the @eveismyname and Paris 2024 pictograms dataset, we trained TI [13] embeddings on Stable Diffusion v1.5 [3] 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 4 and 5 present generated images across two key configurations: (1) regular TI without DP, (2) DPagg-TI with DP at different values of m and ε . To ensure reproducibility and fair comparison across all experimental conditions, we fixed the random seed for the entire generation pipeline. This design choice allows us to isolate the effect of our style transfer method while holding other sources of randomness constant. As with common practice, we set $\delta = 1/n$.

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 4), while those of Paris 2024 pictograms captured the logo’s original structure (Figure 5). In contrast, DP configurations introduce a discernible degradation in image quality, with lower epsilon values and smaller subsampling sizes resulting in diminished stylistic fidelity.

As a no-learning baseline, we consider the limit $\varepsilon \rightarrow 0$, under which u_{DP}^* should convey zero information about the target dataset. Since $\sigma \propto 1/\varepsilon$ is undefined at $\varepsilon = 0$, we approximate this regime by setting $\varepsilon = 10^{-5} \approx 0$, which yields effectively infinite noise.

As $\varepsilon \rightarrow 0$, the resulting token embedding u_{DP}^* gradually loses its semantic meaning, leading to a loss of stylistic fidelity. In particular, $y_{u_{\text{DP}}^*}$ 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 (Figure 4), or a generic icon of a dragon (with color, as opposed to the black and white design of the pictograms, Figure 5). To verify this interpretation, we generated images with the same prompts but without the special token S^* and compared them to the $\varepsilon \approx 0$ generations. The images were visually identical, confirming that at $\varepsilon \approx 0$, the token becomes semantically meaningless and is ignored by the text encoder.

With this in mind, given a fixed seed, ε can be interpreted as a drift parameter, representing the progression from the optimal u_{DP}^* towards a semantically meaningless embedding, 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., $m = 16, \varepsilon \in [0.5, 1]$ in Figure 4 and intermediate ε values in Figure 5) which restores as u_{DP}^* drifts even further from its optimal value. We hypothesize that this is due to drifted u_{DP}^* capturing a different meaning unrelated to the prompt, before losing any meaning that could be interpreted by Stable Diffusion’s text encoder, causing u_{DP}^* to be disregarded from $y_{u_{\text{DP}}^*}$ and the prompt fidelity to be restored.

Meanwhile, reducing m also reduces the sensitivity of the generated image to with respect 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 4 and 5, 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.

C. User Study

To evaluate our approach under different DP and subsampling configurations, we conducted a user study with 25 participants. The goal was to assess whether DPagg-TI preserves perceptual quality while offering privacy guarantees.

1) *Study Design and Setup*: Participants were shown reference images from two datasets: the @eveismyname dataset of private artwork and the Paris 2024 Pictogram dataset. For each dataset, we used 10 prompts to generate images, resulting in 20 groups in total. Each group included images produced under different configurations, including regular TI, DPagg-TI with and without DP noise, with and without subsampling ($m = 8$), style guidance (see Appendix A), and a no-adaptation baseline.

2) *Survey Procedure*: Each participant evaluated two groups, one randomly selected from each dataset. For each group, participants were first shown the reference images, then asked to compare pairs of generated images produced with different configurations (see Figure 6). For each pair, they indicated which image better captured the reference style, or marked the choice as “unsure.”

3) *Results and Analysis*: Survey results are summarized in Table I. 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 smaller subsampling size degraded style fidelity, consistent with the trade-offs inherent in differential privacy. At $\varepsilon = 1$, preferences were split between configurations with and without subsampling, although the subsampling variant was generally favored.



Fig. 4. 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 .

TABLE I
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



Fig. 5. 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 .

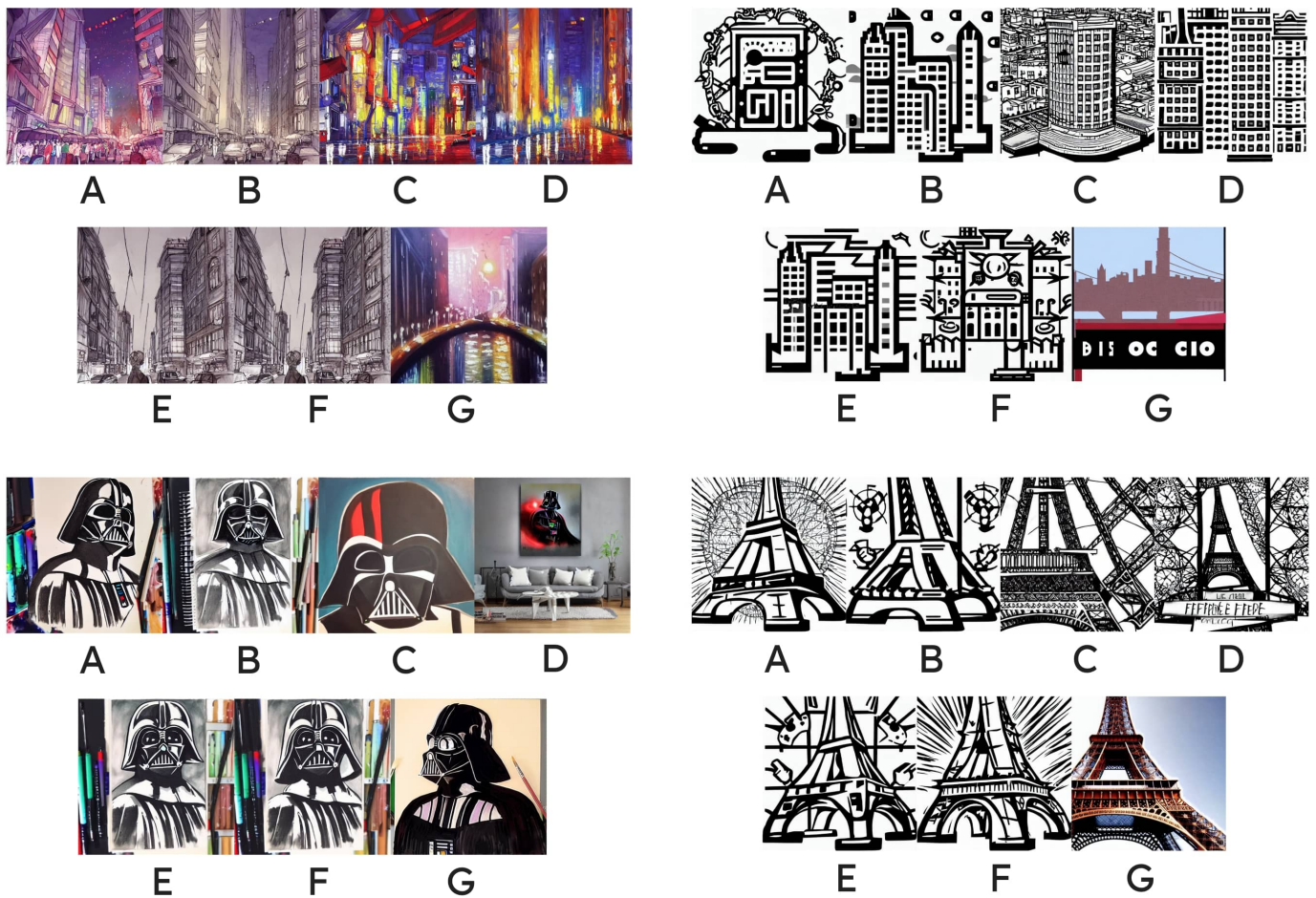


Fig. 6. Samples of image sets used in our user study. Participants were asked to compare 2 images at a time.

Overall, the findings highlight that DPAgg-TI achieves perceptual quality comparable to regular TI, while subsampling serves as an effective mechanism to balance privacy and stylistic fidelity.

D. Kernel Inception Distance

The Kernel Inception Distance (KID) [26] 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 [27], and their distributional differences are estimated. Unlike the more commonly used Fréchet Inception Distance (FID) [28], KID is an unbiased estimator of the true divergence between the learned and target distributions [29], making it more suitable for smaller datasets, as in our case.

We report KID scores for different parameters in Tables II and III. Our results indicate that DPAgg-TI preserves the style transfer fidelity of TI while also ensuring differential privacy. Notably, for @eveismynname ($m = 4$) 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 DPAgg-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 7), 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 unusually high 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. For each parameter configuration, we generate 100 images and compute KID by repeatedly subsampling the larger of the real and generated sets to match the size of the smaller set 100 times, then averaging the resulting KID scores.



Fig. 7. Sample of generated images for KID evaluations with respect to the Paris 2024 dataset.

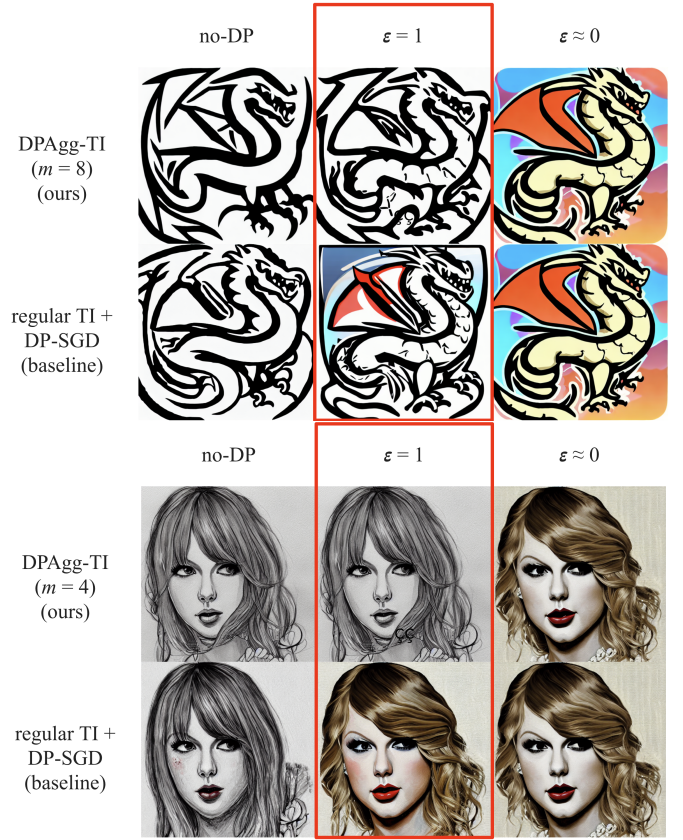


Fig. 8. Comparing our approach to applying DP-SGD to regular TI using prompts “an icon of a dragon in the style of the Paris 2024 Olympic Pictograms” and “a painting of Taylor Swift in the style of @eveismynname” respectively. Note that our method aggregates individual TI embeddings for each training image, whereas the baseline trains a single TI embedding over the entire dataset.

E. Ablation Study: Textual Inversion with DP-SGD

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 @eveismynname and Paris 2024 datasets. We found that in most cases, notably the @eveismynname dataset, the amount of noise required for DP-SGD to achieve a reasonable value of ϵ for DP is so high that the resulting 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 8. 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 C.

V. DISCUSSION

A. Copyright Protection Implications

Our proposed mechanism can also be interpreted through the lens of *copyright protection*. This connection is grounded in the framework of *Near Access-Freeness (NAF)* [6], which evaluates whether a model’s outputs reveal undue influence

TABLE II

KID SCORES OF DPAGG-TI ON @EVEISMNAME DATASET FOR VARIOUS ε VALUES RANGING FROM $\varepsilon = 10^{-5}, 0.1, 0.5, 1.0, 5.0$ (INCLUDING NO DP) UNDER DIFFERENT SUBSAMPLING LEVELS ($m = 4, 8, 16, 32$) AS WELL AS REGULAR TI (CTRL). REPORTED VALUES ARE THE MEAN \pm STANDARD DEVIATION OVER 100 RANDOM SUBSAMPLES.

m	No DP	$\varepsilon = 5.0$	$\varepsilon = 1.0$	$\varepsilon = 0.5$	$\varepsilon = 0.1$	$\varepsilon \approx 0$
–	0.0441 \pm 0.0027	0.0798 \pm 0.0032	0.0526 \pm 0.0022	0.0688 \pm 0.0020	0.1114 \pm 0.0032	0.0654 \pm 0.0027
32	0.0753 \pm 0.0047	0.0836 \pm 0.0042	0.1166 \pm 0.0037	0.0295 \pm 0.0019	0.0644 \pm 0.0021	0.0650 \pm 0.0025
16	0.0350 \pm 0.0020	0.0381 \pm 0.0018	0.0663 \pm 0.0025	0.1303 \pm 0.0033	0.0438 \pm 0.0030	0.0660 \pm 0.0029
8	0.0359 \pm 0.0018	0.0364 \pm 0.0017	0.0366 \pm 0.0019	0.0394 \pm 0.0025	0.0527 \pm 0.0033	0.0654 \pm 0.0024
4	0.0246 \pm 0.0013	0.0251 \pm 0.0016	0.0249 \pm 0.0014	0.0256 \pm 0.0012	0.0313 \pm 0.0017	0.0653 \pm 0.0023
ctrl	0.0314 \pm 0.0010	–	–	–	–	–

TABLE III

KID SCORES OF DPAGG-TI ON PARIS DATASET FOR VARIOUS ε VALUES RANGING FROM $\varepsilon = 10^{-5}, 0.1, 0.5, 1.0, 5.0$ (INCLUDING NO DP) UNDER DIFFERENT SUBSAMPLING LEVELS ($m = 4, 8, 16, 32$) AS WELL AS REGULAR TI (CTRL). REPORTED VALUES ARE THE MEAN \pm STANDARD DEVIATION OVER 100 RANDOM SUBSAMPLES.

m	No DP	$\varepsilon = 5.0$	$\varepsilon = 1.0$	$\varepsilon = 0.5$	$\varepsilon = 0.1$	$\varepsilon \approx 0$
–	0.1153 \pm 0.0055	0.1194 \pm 0.0054	0.1306 \pm 0.0046	0.1395 \pm 0.0057	0.1201 \pm 0.0053	0.1274 \pm 0.0055
32	0.1222 \pm 0.0066	0.1036 \pm 0.0065	0.1375 \pm 0.0047	0.1311 \pm 0.0048	0.1248 \pm 0.0060	0.1258 \pm 0.0054
16	0.1321 \pm 0.0057	0.1411 \pm 0.0077	0.1309 \pm 0.0061	0.1380 \pm 0.0047	0.1359 \pm 0.0060	0.1273 \pm 0.0057
8	0.1303 \pm 0.0084	0.1303 \pm 0.0074	0.1112 \pm 0.0062	0.1311 \pm 0.0064	0.1318 \pm 0.0052	0.1267 \pm 0.0056
4	0.1158 \pm 0.0057	0.1085 \pm 0.0056	0.1184 \pm 0.0068	0.1194 \pm 0.0065	0.1592 \pm 0.0065	0.1268 \pm 0.0055
ctrl	0.1383 \pm 0.0066	–	–	–	–	–

from specific data points by comparing them to those from a safe model trained without access to the same data.

Modern generative models typically produce outputs via randomized sampling. Leveraging this inherent randomness, Vyas et al. [6] introduced NAF as a metric to quantify the similarity between a model’s output and copyrighted content. The key idea is to compare the output distribution of a potentially infringing model to that of a *safe* model—one trained without access to the target content. A canonical example is the *leave-one-out-safe* model, trained on the full dataset excluding x . Since $\text{safe}(x)$ lacks access to x , the probability that it generates content resembling x is expected to be small; any such resemblance is considered fortuitous.

Definition 2 (Near Access-Freeness [6]). *Let \mathcal{C} be a set of copyrighted samples and \mathcal{W} a set of generative models. Given a mapping $\text{safe} : \mathcal{C} \rightarrow \mathcal{W}$ and a divergence measure Δ , we say a model $w \in \mathcal{W}$ is k_y -near access-free (or k_y -NAF) on prompt $y \in \mathcal{Y}$ if for every $x \in \mathcal{C}$,*

$$\Delta(P_w(\cdot|y) \| P_{\text{safe}(x)}(\cdot|y)) \leq k_y.$$

If $k_y = 0$, the model is indistinguishable from a safe model, meaning any resemblance to copyrighted material is by random chance. More generally, a small k_y suggests the model is unlikely to generate outputs resembling x with higher probability than a model that has never seen x .

NAF is closely related to concepts in DP. Depending on the divergence measure Δ , NAF resembles different DP variants—for example, ε -DP when $\Delta = \Delta_{\max}$ [10], and $(1, \varepsilon)$ -Rényi DP when $\Delta = \Delta_{\text{KL}}$. Translating DP to generative models yields this definition:

Definition 3 (Differentially Private Generation (DPG)). *Let S and S' be neighboring datasets. Denote by $P_S(\cdot|y)$ the distribution over outputs generated by a model trained on, or adapted from, S with algorithm \mathcal{A} , where randomness includes both training and generation. The generation satisfies ε -Differentially Private Generation (ε -DPG) if for every $y \in \mathcal{Y}$,*

$$\Delta(P_S(\cdot|y) \| P_{S'}(\cdot|y)) \leq \varepsilon.$$

Here, *neighboring datasets* differ by a single data point (or privacy unit). If the training process is ε -DP, then the outputs naturally satisfy ε -DPG via the data processing inequality. One benefit of DPG is the flexibility to add noise during generation rather than training, potentially improving the utility–privacy tradeoff. However, there are notable distinctions: ε -DP offers protection under arbitrary post-processing and multiple outputs, whereas ε -DPG only guarantees privacy for single outputs. Also, under DP, the trained model can be released, but under DPG, only the outputs are safe to share. Elkin-Koren et al. [30] highlight further differences: NAF is *one-sided*—comparing a model to a fixed safe reference—whereas DPG is *symmetric*. This asymmetry in NAF can enable better utility. Additionally, NAF allows more flexibility in choosing the safe model, which can be exploited in algorithm design.

Since DPAGG-TI satisfies ε -DP, it also satisfies ε -NAF under the leave-one-out-safe model. Within the NAF framework, this means the adapted model behaves similarly to one that never saw the private images. Importantly, this guarantee is meaningful only *within NAF*; it does not imply broader legal immunity or empirical indistinguishability from the original content. However, it allows us to argue that any close resem-

balance between outputs and private training data is no more likely than would be expected from a model with no access to that data.

Finally, the goal of DPAGg-TI is to adapt to the *style* of private image sets, not their precise content. This distinction matters: pure style imitation (without reproducing protectable expression or “substantive elements”) is often argued to be non-infringing in many creative contexts, though the legal status is jurisdiction- and fact-dependent; particularly in artistic and creative contexts. As discussed in Elkin-Koren et al. [30] and legal analyses such as Carlini et al. [5], generating new content in the style of a work, without reproducing its substantive elements, is generally not considered copyright infringement. Therefore, the use of DPAGg-TI to learn and reproduce stylistic attributes does not contradict the spirit or intent of the NAF framework. Instead, it offers a promising direction for responsibly fine-tuning generative models on private or copyrighted sources while respecting both privacy and intellectual property boundaries.

Remark 1 (Scope of Protection and Artist-Level Extension). *We provide record-level DP: it limits leakage or reconstruction of any individual private image, not an artist’s entire style. Consequently, it yields a corresponding NAF-style guarantee at the per-image level (under the chosen safe reference). This interpretation should be understood strictly within NAF and does not constitute a general copyright compliance claim. The DP guarantee continues to apply under targeted prompts: conditioning on detailed descriptions can increase the likelihood of reproducing a specific private work by at most an e^ϵ factor (up to δ). Artist-level (user-level) DP is conceptually possible by treating each artist as one unit and privately aggregating artist-level embeddings (e.g., via a DP mean mechanism), but typically requires stronger noise and may reduce utility; we leave a full exploration to future work.*

B. Limitations

DPAGg-TI is designed for the low-data, strong-privacy regime, where the number of private images is small ($n \approx 100$) and per-record protection with $\epsilon < 5$ matters. For large datasets with moderate subsampling, DP-SGD on the full model may become more efficient and could provide better utility. We explicitly position our method for scenarios where DP-SGD is known to struggle: strong privacy guarantees with limited training data.

We acknowledge that in moderate privacy regimes ($5 \leq \epsilon \leq 10$) with larger batch sizes and careful tuning, DP-SGD with parameter-efficient fine-tuning methods might perform better than our approach. However, in our experiments, applying DP-SGD to regular TI even at $\epsilon = 20$ with carefully tuned hyperparameters (learning rate, batch size, scheduler) still required prohibitively high noise levels to satisfy the privacy accountant, preventing meaningful learning. Moreover, DP-SGD typically involves multiple training iterations, so the effective noise further accumulates due to composition across epochs, making convergence extremely difficult in the strong-privacy regime ($\epsilon < 5$) with approximately 100 images.

VI. CONCLUSION

We presented a differentially private adaptation method for diffusion models based on Textual Inversion, enabling privacy-preserving style transfer without the need for full model fine-tuning. By learning per-image embeddings and aggregating them with calibrated noise, our method, DPAGg-TI, achieves strong formal privacy guarantees while maintaining high output fidelity. Experiments on private artwork and Paris 2024 pictograms show that DPAGg-TI consistently outperforms DP-SGD, which fails to produce meaningful results under comparable privacy budgets. These results highlight the effectiveness of embedding-level adaptation as an efficient and scalable alternative to traditional gradient-based approaches, especially in low-data regimes. Unlike DP-SGD, which introduces significant computational overhead and utility degradation, DPAGg-TI is lightweight, modular, and compatible with existing diffusion backbones. Our findings suggest that embedding-centric approaches offer a promising direction for privacy-aware personalization, and motivate further research into cross-modal extensions, improved aggregation techniques, and integration with broader privacy-preserving frameworks.

ETHICAL 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.

ACKNOWLEDGMENTS

We sincerely thank @eveismynname 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. This work was supported in part by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. RS-2025-23525649).

REFERENCES

- [1] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *Advances in Neural Information Processing Systems*, 2020.
- [2] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based generative modeling through stochastic differential equations,” in *International Conference on Learning Representations*, 2021.
- [3] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [4] J. Duan, F. Kong, S. Wang, X. Shi, and K. Xu, “Are diffusion models vulnerable to membership inference attacks?” in *International Conference on Machine Learning*, 2023, pp. 8717–8730.
- [5] N. Carlini, J. Hayes, M. Nasr, M. Jagielski, V. Schwag, F. Tramèr, B. Balle, D. Ippolito, and E. Wallace, “Extracting training data from diffusion models,” in *USENIX Security Symposium*, 2023.
- [6] N. Vyas, S. M. Kakade, and B. Barak, “On provable copyright protection for generative models,” in *International Conference on Machine Learning*, 2023, pp. 35 277–35 299.
- [7] C. Dwork, “Differential privacy,” in *International Colloquium on Automata, Languages, and Programming*, 2006.
- [8] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in *ACM SIGSAC Conference on Computer and Communications Security*, 2016.
- [9] S. Hoory, A. Feder, A. Tendler, S. Erell, A. Peled-Cohen, I. Laish, H. Nakhost, U. Stemmer, A. Benjamini, A. Hassidim *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.
- [10] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Theory of Cryptography Conference*, 2006.
- [11] T. Dockhorn, T. Cao, A. Vahdat, and K. Kreis, “Differentially private diffusion models,” *Transactions on Machine Learning Research*, 2023.
- [12] S. Ghalebikesabi, L. Berrada, S. Gowal, I. Ktena, R. Stanforth, J. Hayes, S. De, S. L. Smith, O. Wiles, and B. Balle, “Differentially private diffusion models generate useful synthetic images,” *arXiv preprint arXiv:2302.13861*, 2023.
- [13] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-or, “An image is worth one word: Personalizing text-to-image generation using textual inversion,” in *International Conference on Learning Representations*, 2023.
- [14] Paris 2024, “Paris 2024 - pictograms,” <https://olympics.com/en/paris-2024/the-games/the-brand/pictograms>.
- [15] J. Song, C. Meng, and S. Ermon, “Denoising diffusion implicit models,” in *International Conference on Learning Representations*, 2021.
- [16] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning*, 2021.
- [17] C. Dwork, A. Roth *et al.*, “The algorithmic foundations of differential privacy,” *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- [18] B. Balle and Y.-X. Wang, “Improving the gaussian mechanism for differential privacy: Analytical calibration and optimal denoising,” in *International Conference on Machine Learning*, 2018.
- [19] I. Mironov, “Rényi differential privacy,” in *IEEE Computer Security Foundations Symposium*, 2017.
- [20] T. Steinke, “Composition of differential privacy & privacy amplification by subsampling,” *arXiv preprint arXiv:2210.00597*, 2022.
- [21] D. Chen, S.-c. S. Cheung, C.-N. Chuah, and S. Ozonoff, “Differentially private generative adversarial networks with model inversion,” in *IEEE International Workshop on Information Forensics and Security*, 2021.
- [22] F. Tramèr, G. Kamath, and N. Carlini, “Position: Considerations for differentially private learning with large-scale public pretraining,” in *International Conference on Machine Learning*, 2024.
- [23] Y. Yu, M. Sanjabi, Y. Ma, K. Chaudhuri, and C. Guo, “Vip: A differentially private foundation model for computer vision,” in *International Conference on Machine Learning*, 2024.
- [24] T. Sander, Y. Yu, M. Sanjabi, A. O. Durmus, Y. Ma, K. Chaudhuri, and C. Guo, “Differentially private representation learning via image captioning,” in *International Conference on Machine Learning*, 2024.
- [25] International Olympic Committee, “Olympic properties,” <https://olympics.com/ioc/olympic-properties>.
- [26] M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton, “Demystifying mmd gans,” in *International Conference on Learning Representations*, 2018.
- [27] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.
- [28] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” in *Advances in Neural Information Processing Systems*, 2017.
- [29] S. Jayasumana, S. Ramalingam, A. Veit, D. Glasner, A. Chakrabarti, and S. Kumar, “Rethinking fid: Towards a better evaluation metric for image generation,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9307–9315.
- [30] N. Elkin-Koren, U. Hachohen, R. Livni, and S. Moran, “Can copyright be reduced to privacy?” in *Symposium on Foundations of Responsible Computing*, 2024, pp. 3:1–3:18.
- [31] A. Bansal, H.-M. Chu, A. Schwarzschild, S. Sengupta, M. Goldblum, J. Geiping, and T. Goldstein, “Universal guidance for diffusion models,” in *International Conference on Learning Representations*, 2024.
- [32] Innat, “Van gogh paintings,” <https://www.kaggle.com/datasets/ipythonx/van-gogh-paintings>.

APPENDIX A
DIFFERENTIALLY PRIVATE ADAPTATION VIA STYLE GUIDANCE

A. Background: Denoising Diffusion Implicit Models

Denoising Diffusion Implicit Models (DDIM) sampling [15] 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

$$\hat{x}_0 = \frac{x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t, y, t)}{\sqrt{\alpha_t}} \quad (9)$$

$$x_{t-1} = \sqrt{\alpha_{t-1}} \hat{x}_0 + \sqrt{1 - \alpha_{t-1}} \epsilon_\theta(x_t, y, t). \quad (10)$$

B. Implementation

We extend our approach to style guidance (SG) by leveraging the framework of Universal Guidance [31]. 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_{DP}^* using (6) or (8), depending on whether subsampling is applied. The aggregated embedding u_{DP}^* is then incorporated into the reverse diffusion process as a style guide.

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

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

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_{DP}^* using the centroid method. The aggregated u_{DP}^* guides the reverse diffusion process:

$$\hat{\epsilon}_\theta(x_t, y, t) = \epsilon_\theta(x_t, y, t) + w \sqrt{1 - \alpha_t} \nabla_{x_t} \ell_{\cos}(u_{\text{DP}}^*, \mathcal{E}(\hat{x}_0)). \quad (12)$$

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

C. Style Transfer Results



Fig. 9. 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 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 9, 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. [31] 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.

D. Ablation

To better understand the limited effectiveness of style guidance in our experiments, despite its success in Bansal et al. [31], we applied our approach to a dataset of 143 paintings from Van Gogh’s Saint-Paul Asylum, Saint-Rémy collection [32] (Figure 10). Unlike the @eveismynname and Paris 2024 datasets, it is highly likely that Stable Diffusion has been trained on these images. Additionally, Bansal et al. [31] 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 11. 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.



Fig. 10. Sample of paintings by Van Gogh used to generate style guidance embeddings.



Fig. 11. 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)”.

APPENDIX B COMPUTATIONAL COST COMPARISONS

Direct comparisons of computational cost across methods are challenging due to differing training paradigms (per-image optimization vs. dataset-level training), optimization procedures, and privacy accounting. Nonetheless, to provide a concrete sense of scale, we report representative costs measured using Stable Diffusion v1.5 on a single NVIDIA A100 GPU (Tables IV and V). For each method, we tuned the number of optimization steps to reach its best utility under the target privacy budget.

a) Sequential vs. batched execution: The per-image runtime reported in Table IV corresponds to a sequential implementation that optimizes each textual-inversion embedding independently. In practice, we can optimize multiple embeddings jointly by batching several images at once (and optionally subsampling the private set per update), which reduces the effective wall-clock cost and avoids the naive linear scaling implied by “minutes per image $\times n$.”

TABLE IV
TRAINING COST COMPARISON ACROSS METHODS. OVERHEAD FROM DP-SGD IS RELATIVELY MODEST DUE TO THE LOW-DIMENSIONAL EMBEDDING BEING OPTIMIZED. N/A FOR SG MEANS NOTHING IS TRAINED ASIDE FROM THE BASE MODEL.

Method	Steps	Batch Size	Time	Memory Usage
TI (no DP)	10,000 (for 150 images)	1	25 min	7 GB
		8	2.5 hours	20 GB
TI (DP-SGD)	30,000 (for 150 images)	1	80 min	7 GB
		8	7 hours	20 GB
DPAgg-TI	2,000 per image	N/A	~ 5 min/image	7 GB
SG	N/A	N/A	N/A	N/A

b) When DP-SGD can be faster: DP-SGD amortizes computation across the dataset and can be faster wall-clock-wise for larger n and/or moderate privacy budgets. In contrast, our method is designed for low-data personalization with strong per-record guarantees, and it offers a practical advantage in dynamic settings: it supports incremental updates to the private set (e.g. adding or removing images) without retraining a large set of model parameters, whereas DP-SGD-style training typically requires rerunning optimization to reflect such changes. We view these approaches as complementary, targeting different operating regimes.

TABLE V
INFERENCE COST COMPARISON ACROSS METHODS.

Method	Steps	Batch Size	Time	Memory Usage
TI (no DP, DP-SGD, DPAGg-TI)	50	1	1–2 sec	4 GB
	100	1	1–2 min	58 GB
SG (no DP, DPAGg-SG)	500	1	~30 min	17 GB

APPENDIX C
ADDITIONAL STYLE TRANSFER AND ABLATION RESULTS

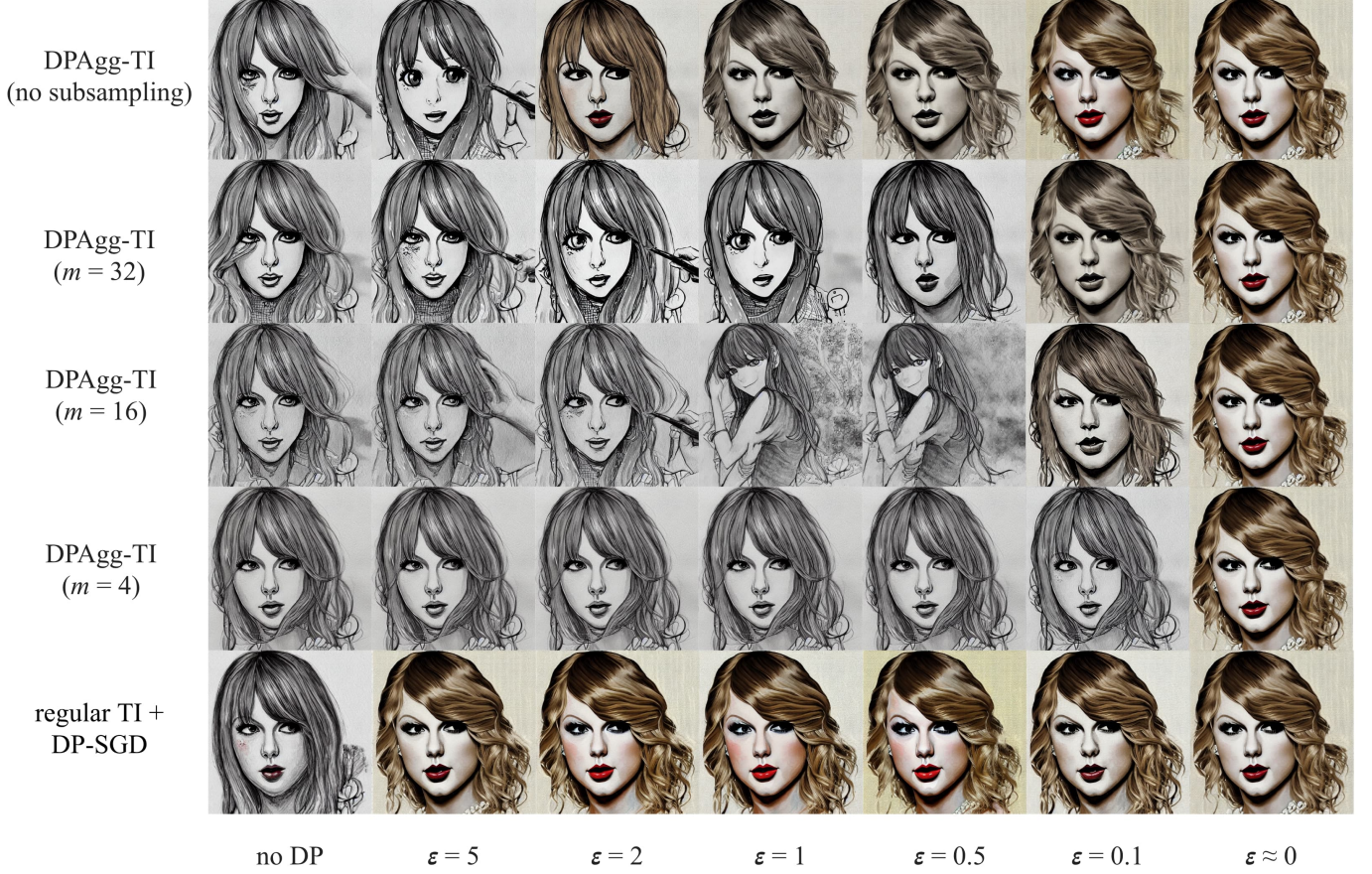


Fig. 12. Images generated by Stable Diffusion v1.5 using the prompt “A painting of Taylor Swift in the style of <@eveismyname>”, with the embedding <@eveismyname> trained using DPAGg-TI (with different subsample sizes m) and TI with DP-SGD using different values of ϵ .

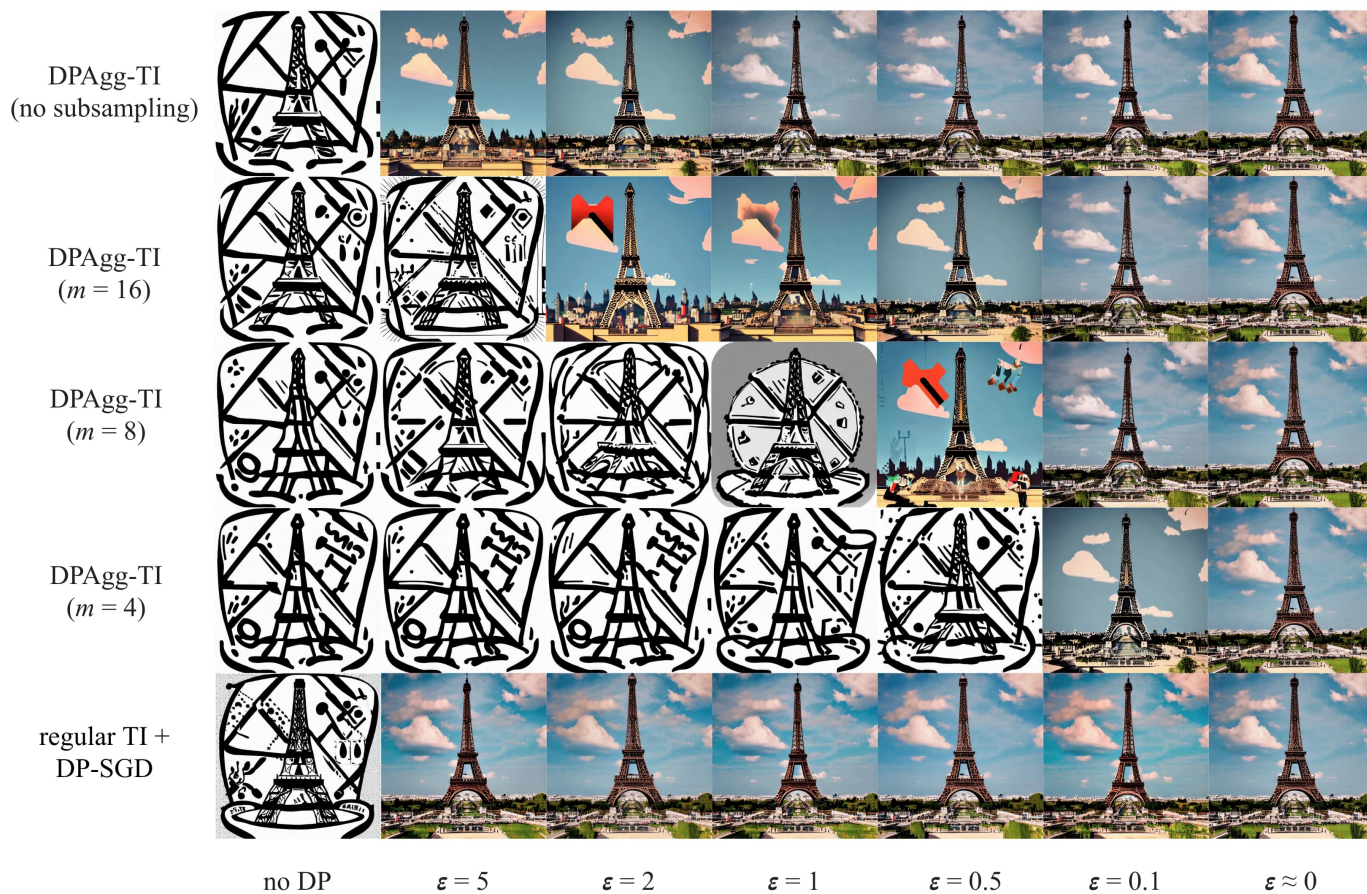


Fig. 13. Images generated by Stable Diffusion v1.5 using the prompt “An icon of the Eiffel Tower in the style of <Paris 2024 Pictograms>”, with the embedding <Paris 2024 Pictograms> trained using DPAgg-TI (with different subsample sizes m) and TI with DP-SGD using different values of ε .



Fig. 14. 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 .