

Appendix Overview

This appendix is organized as follows. Appendix A.1 represents Proposition 1 and its proof. Appendix A.3 provides detailed proofs of Theorem 1. Appendix B.1 provides additional details of our experimental setup. Appendix B.2 visualizes the synthetic-private image pairs under our method and baseline methods to intuitively assess reconstruction risk. Appendix B.3 provides additional details of our ablation study.

A. Theoretical Results

A.1. Proposition 1 and Its Proof

Proposition 1. *Under smoothness and strong convex on $\mathcal{L}(\theta, \gamma)$, the solution $(\gamma^*, \{\theta^*\})$ of the unrolled problem Eq. (7) converges to that of the original bilevel formulation Eq. (4b) as $T \rightarrow \infty$ and $\eta \rightarrow 0$ appropriately.*

Proof. For fixed γ , strong convexity and smoothness imply, we have

$$\|\theta^{(T)}(\gamma) - \theta^*(\gamma)\| \leq (1 - \eta\mu)^T \|\theta^{(0)} - \theta^*(\gamma)\|. \quad (13)$$

Because $0 < \eta < 2/L$, the factor $(1 - \eta\mu)^T$ decays exponentially; hence for any $\varepsilon > 0$ we can pick T_0 (independent of γ owing to boundedness of γ and continuity of θ^*) so that

$$\|\theta^{(T)}(\gamma) - \theta^*(\gamma)\| \leq \varepsilon \quad (14)$$

for all $T \geq T_0$ and all $\gamma \in \Gamma$. Denote

$$\begin{aligned} F_T(\gamma) &:= \mathcal{L}(\theta^{(T)}(\gamma), \gamma), & F(\gamma) &:= \mathcal{L}(\theta^*(\gamma), \gamma) \\ &\& \Delta_T(\gamma) &:= F_T(\gamma) - F(\gamma). \end{aligned} \quad (15)$$

Since $\gamma \in [0, 1]^M$ is bounded, we have

$$|\Delta_T(\gamma)| \leq G \|\theta^{(T)}(\gamma) - \theta^*(\gamma)\|, \quad (16)$$

Therefore, $|\Delta_T(\gamma)|$ converges uniformly to 0 as $T \rightarrow \infty$. Moreover, uniform convergence $F_T \rightarrow F$ on the compact set γ implies $\gamma_T^* \rightarrow \gamma^*$.

Finally, Combining Eq. (14) gives $(\gamma_T^*, \theta^{(T)}(\gamma_T^*)) \rightarrow (\gamma^*, \theta^*(\gamma^*))$, completing the proof. \square

A.2. Auxiliary Lemma

Lemma 1 (Standard Form of Pontryagin’s Maximum Principle [30]). *Consider the following optimization problem in a discrete dynamical system:*

$$\begin{aligned} \min_{\gamma_t} & \sum_{t=0}^{T-1} \mathcal{K}(\theta_t, \gamma_t) + K(\theta_T) \\ \text{s.t.} & \theta_{t+1} = f(\theta_t, \gamma_t), \quad \gamma_t \in U \end{aligned} \quad (17)$$

where the state variable $\theta_t \in \mathbb{R}^N$, the control variable $\gamma_t \in \mathbb{R}^D$, and $\mathcal{K} : \mathbb{R}^{N \times D} \mapsto \mathbb{R}$, $K : \mathbb{R}^{N \times D} \mapsto \mathbb{R}$, $f : \mathbb{R}^{N \times D} \mapsto \mathbb{R}^N$ are continuous in $\mathbb{R}^{N \times D}$.

Let γ_t^* be the solution to this problem, and θ_t^* denote the corresponding state variable. For $0 \leq t < T$, there exists a co-state vector $\lambda_t^* \in \mathbb{R}^N$ such that

$$\theta_{t+1}^* = \nabla_{\lambda} H(\theta_t^*, \lambda_{t+1}^*, \gamma_t^*), \quad \theta_0^* = \theta_0 \quad (18)$$

$$\lambda_t^* = \nabla_{\theta} H(\theta_t^*, \lambda_{t+1}^*, \gamma_t^*), \quad \lambda_T^* = \nabla K(\theta_T) \quad (19)$$

$$\gamma_t^* = \arg \min_{\gamma_t} H(\theta_t^*, \lambda_{t+1}^*, \gamma_t), \quad \gamma_t \in U \quad (20)$$

where $H : \mathbb{R}^N \times \mathbb{R}^N \times \mathbb{R}^D \mapsto \mathbb{R}$ is the Hamiltonian function defined by

$$H(\theta, \lambda, \gamma) = \mathcal{K}(\theta, \gamma) + \lambda^\top f(\theta, \gamma) \quad (21)$$

Proof. As the standard form of discrete-time Pontryagin’s Maximum Principle in optimal control for time-variant control variables, the proof of Lemma 1 is available in textbooks [25]. \square

A.3. Proof of Theorem 1.

Theorem 1 (PMP Conditions for Data Selection). *Let γ^* solve the problem in Eq. (7), and θ_t^* denote the downstream task model trained with γ^* . For $0 \leq t < T$, there exists a vector $\lambda_t^* \in \mathbb{R}^N$ such that*

$$\begin{cases} \theta_{t+1}^* = \theta_t^* - \eta \nabla \mathcal{L}(\theta_t^*, \gamma^*), & \theta_0^* = \theta_0, \\ \lambda_t^* = \lambda_{t+1}^* + \nabla \ell_u(\theta_t^*) - \eta \nabla^2 \mathcal{L}(\theta_t^*, \gamma_0) \lambda_{t+1}^*, & \lambda_T^* = \nabla \ell_u(\theta_T^*), \\ \gamma^* = \arg \max_{\gamma \in U} \left\{ \sum_{i=1}^M \gamma_i \left[\sum_{t=0}^{T-1} (\lambda_{t+1}^*)^\top \nabla l(x_{g,i}, \theta_t^*) \right. \right. \\ \left. \left. - \frac{\alpha_1}{\eta} \mathcal{L}_{\text{pix}}(x_{g,i}) - \frac{\alpha_2}{\eta} \mathcal{L}_{\text{mem}}(x_{g,i}) \right] \right\}. \end{cases} \quad (22)$$

where $\nabla^2 \mathcal{L}(\theta_t^*, \gamma^*)$ denotes the Hessian matrix of $\mathcal{L}(\theta, \gamma^*)$ with respect to θ evaluated at $\theta = \theta_t^*$.

Proof. We leverage the standard form of discrete-time Pontryagin’s Maximum Principle [30] to prove the theorem.

Here, θ_t is the downstream task model parameters, \mathcal{K} is the utility loss and privacy leakage bi-objective optimization object. f is the gradient descent operation where the data weights γ_t changes with respect to the training steps t . As $\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t, \gamma_t)$, the Hamilton function is

$$H(\theta, \lambda, \gamma) = \mathcal{K}(\theta, \gamma) + \lambda^\top [\theta - \eta \nabla L(\theta, \gamma)]. \quad (23)$$

We assume the data weights keep consistent during the T steps training as $\gamma_1 = \gamma_0, \gamma_2 = \gamma_0, \dots, \gamma_{T-1} = \gamma_0$. We

have:

$$\begin{aligned} & \min_{\gamma_t \in \mathcal{U}} \sum_{t=1}^T \left[\ell_u(\theta_t) + \frac{\alpha}{T} \ell_p(\gamma) \right], \\ & \text{s.t. } \theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t, \gamma), \quad t = 0, \dots, T-1, \\ & \quad \gamma_0 = \gamma_1 = \dots = \gamma_{T-1}. \end{aligned} \quad (24)$$

We apply the method of Lagrange multipliers to solve Eq. (24) which can be converted to the following problem:

$$\begin{aligned} & \min_{\gamma_t \in \mathcal{U}} \sum_{t=0}^{T-1} \left[\ell_u(\theta_t) + \frac{\alpha}{T} \ell_p(\gamma) \right] + \sum_{t=1}^{T-1} \sum_{i=1}^M \mu_{i,t} (\gamma_{i,t} - \gamma_{i,0}) \\ & \quad + \left[\ell_u(\theta_T) + \frac{\alpha}{T} \ell_p(\gamma) \right], \\ & \text{s.t. } \theta_{t+1} = \theta_t - \eta \nabla \mathcal{L}(\theta_t, \gamma_t), \end{aligned} \quad (25)$$

where $(\mu_{i,t})_{1 \leq i \leq M, 0 \leq t \leq T-1}$ are Lagrange multipliers.

When $\mathcal{K}(\theta_t, \gamma_t)$ and $K(\theta_T)$ takes the following form, we can use Lemma 1 to solve Eq. (26).

$$\mathcal{K}(\theta_t, \gamma_t) = \begin{cases} \ell_u(\theta_0) - \sum_{t'=1}^{T-1} \sum_{i=1}^M \mu_{i,t'} \gamma_{i,0} + \frac{\alpha}{T} \ell_p(\gamma) & \text{if } t = 0 \\ \ell_u(\theta_t) + \sum_{i=1}^M \mu_{i,t} \gamma_{i,t} + \frac{\alpha}{T} \ell_p(\gamma) & \text{if } 1 \leq t \leq T-1 \end{cases} \quad (26)$$

Plugging the expressions from Eq. (23) and Eq. (26) into Eq. (18) and Eq. (19) yields:

$$\theta_{t+1}^* = \theta_t^* - \eta \nabla \mathcal{L}(\theta_t^*, \gamma_0^*), \quad \theta_0^* = \theta_0, \quad (27)$$

$$\begin{aligned} \lambda_t^* &= \lambda_{t+1}^* + \nabla \ell_u(\theta_t^*) - \eta \nabla^2 \mathcal{L}(\theta_t^*, \gamma_0^*) \lambda_{t+1}^*, \\ \lambda_T^* &= \nabla \ell_u(\theta_T), \end{aligned} \quad (28)$$

We have proved the first and second line in Eq. (12) in Theorem 1 when we set $\gamma_0^* = \gamma^*$.

Plugging Eq. (23) and Eq. (26) into Eq. (20) yields:

$$\gamma_t^* = \begin{cases} \arg \max_{\gamma_0} \left\{ \sum_{i=1}^M \gamma_{i,0} \left[\lambda_1^{*\top} \nabla l(x_{g,i}, \theta_0^*) + \sum_{t'=1}^{T-1} \mu_{i,t'} \right] - \frac{\alpha}{T} \ell_p(\gamma_0) \right\}, & \text{if } t = 0, \\ \arg \max_{\gamma_t} \left\{ \sum_{i=1}^M \gamma_{i,t} \left[\lambda_{t+1}^{*\top} \nabla l(x_{g,i}, \theta_t^*) - \mu_{i,t} \right] - \frac{\alpha}{T} \ell_p(\gamma_t) \right\}, & \text{if } 1 \leq t \leq T-1. \end{cases} \quad (29)$$

Having the time-invariant constraint assumption, we set $\gamma_0^* = \gamma_1^* = \dots = \gamma_{T-1}^* = \gamma^*$,

$$\begin{cases} \gamma^* = \arg \max_{\gamma} \left\{ \sum_{i=1}^M \gamma_{i,0} \left[\lambda_1^{*\top} \nabla l(x_{g,i}, \theta_0^*) + \sum_{t'=1}^{T-1} \mu_{i,t'} \right] - \frac{\alpha}{T} \ell_p(\gamma) \right\}, & \text{if } t = 0, \\ \gamma^* = \arg \max_{\gamma} \left\{ \sum_{i=1}^M \gamma_{i,t} \left[\lambda_{t+1}^{*\top} \nabla l(x_{g,i}, \theta_t^*) - \mu_{i,t} \right] - \frac{\alpha}{T} \ell_p(\gamma) \right\}, & \text{if } 1 \leq t \leq T-1. \end{cases} \quad (30)$$

Hence, we obtain T equations for the T unknowns as $T-1$ number of $\mu_t = [\mu_{1,t}, \mu_{2,t}, \dots, \mu_{M,t}]$ plus one γ^* . The resulting solution is:

$$\begin{aligned} \mu_{i,t} &= \eta \lambda_{t+1}^{*\top} \nabla l(x_{g,i}, \theta_t^*) - \frac{\eta}{T} \sum_{t=0}^{T-1} \lambda_{t+1}^{*\top} \nabla l(x_{g,i}, \theta_t^*), \\ & 1 \leq i \leq M, \quad 0 \leq t \leq T-1, \end{aligned} \quad (31)$$

$$\begin{aligned} \gamma^* &= \arg \max_{\gamma} \left\{ \sum_{i=1}^M \gamma_i \frac{\eta}{T} \sum_{t=0}^{T-1} \lambda_{t+1}^{*\top} \nabla l(x_{g,i}, \theta_t^*) - \frac{\alpha}{T} \ell_p(\gamma) \right\} \\ &= \arg \max_{\gamma} \sum_{i=1}^M \gamma_i \left[\frac{\eta}{T} \sum_{t=0}^{T-1} \lambda_{t+1}^{*\top} \nabla l(x_{g,i}, \theta_t^*) - \frac{\alpha}{T} (\mathcal{L}_{\text{pix}}(x_{g,i}) + \mathcal{L}_{\text{mem}}(x_{g,i})) \right] \\ &= \arg \max_{\gamma \in \mathcal{U}} \left\{ \sum_{n=1}^M \gamma_n \left[\sum_{t=0}^{T-1} (\lambda_{t+1}^*)^{\top} \nabla l(x_{g,i}, \theta_t^*) - \frac{\alpha_1}{\eta} \mathcal{L}_{\text{pix}}(x_{g,i}) - \frac{\alpha_2}{\eta} \mathcal{L}_{\text{mem}}(x_{g,i}) \right] \right\}, \end{aligned} \quad (32)$$

where α_1 and α_2 are the parameters controlling MIA and RA respectively.

We have finished the proof of Theorem 1 with Eq. (27), Eq. (28), and Eq. (32). \square

B. Experimental Results

B.1. Experiment Setup

Datasets & Model We conduct experiments on three datasets, including two benchmark datasets as *ImageNet* [6] and *DomainNet* [29], and a medical datasets as *PathMNIST* [39]. *ImageNet* is a large-scale image classification benchmark containing over 1.2 million images across 1,000 object categories. In our experiment, we randomly select 20 classes with 200 images per class as: banana, butterfly, car,

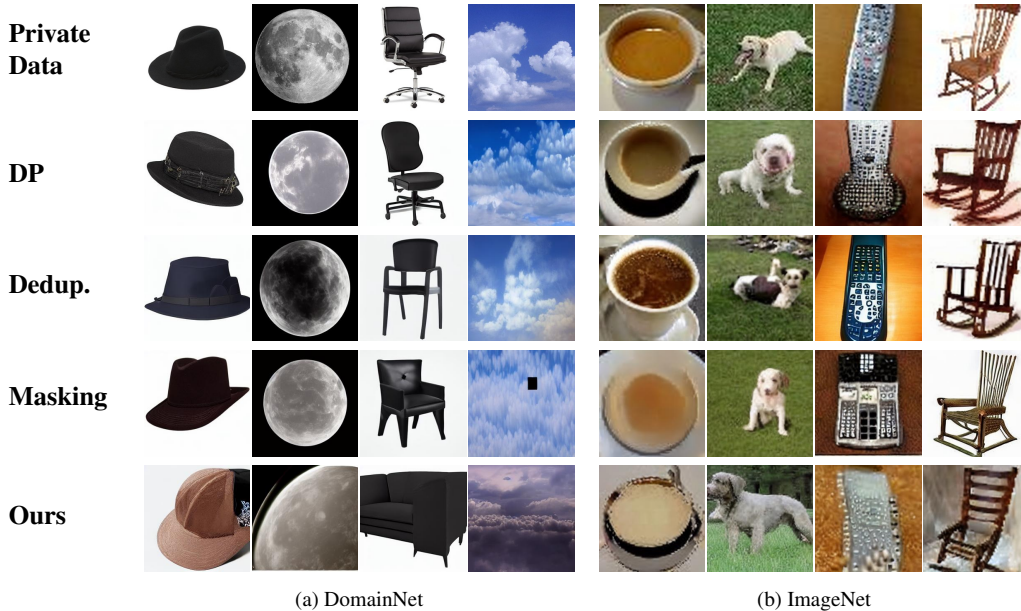


Figure 5. **Most similar private–synthetic pairs under different scenarios.** Top row: private data. Second Row: original synthetic data without defense. Rows 3–6: Differential Privacy (DP), De-duplication (Dedup.), Masking, and Ours. Columns: representative examples from (a) DomainNet and (b) ImageNet. Our method yields the least visually similar pairs, indicating the strongest mitigation of privacy leakage.

cat, chair, coffee, dog, elephant, fish, girl, lobster, nut, orange, pepper, pizza, remote, scorpion, ship, train and van. In each class, we equally split the images and train and test data. *DomainNet* is a widely used benchmark, comprising approximately 600,000 images from 6 distinct domains (e.g., clipart, real, sketch), each annotated with 345 shared categories. In our experiment, we randomly select 20 classes in “real” domain with 200 images per class as: bird, book, butterfly, car, chair, cloud, cow, dog, duck, hat, horse, house, map, moon, pear, pig, pizza, train, tree, van. In each class, we equally split the images into train and test data. *PathMNIST* We utilize the PathMNIST subset provided by the MedMNIST benchmark suite [39]. In our experiment, we use 8 classes with 60 images per class as: adipose, background, debris, lymphocytes, muscle, normal colon mucosa, cancer-associated stroma and tumor epithelial. In each class, we equally split the images and train and test data. We clarify that all references to the dataset in this paper correspond to PathMNIST, which does not affect any experimental results.

Evaluation Metrics. For classification task, we measure downstream utility by top-1 accuracy on a held-out test set (utility loss as $1 - \text{acc.}$) For object detection task, we use the Faster R-CNN model [31] with a ResNet-50 backbone [15] and Feature Pyramid Network (FPN) [27]. We measure the downstream task utility by using the average precision at an intersection-over-union threshold of 0.5 (AP_{50}). Annotation of both private and synthetic datasets was performed automatically using a pretrained detector. The generative

model, synthetic data generation, and baselines are consistent with the main experiment. All experiments are repeated with three random seeds, and results are reported as mean \pm standard deviation.

Baselines. We compare our proposed PrivSynth against Differential Privacy [1], De-Duplication [23], Masking [16], Filtering [11], and training with private data (No Defense). First, **Differential Privacy (DP)** [1] injects Gaussian noise into each class embedding, with the noise standard deviation controlling the strength of privacy protection. Second, **De-Duplication** [23] computes LPIPS distances across the private dataset and removes one image from each most-similar pair, varying the fraction of images dropped to adjust privacy. Third, **Masking** [16] divides each image into an 8×8 grid and randomly blackens a specified number of blocks per image—more masked blocks correspond to stronger privacy. Fourthly, **Filtering** [11] trains a membership inference attack on the target model and excludes those synthetic samples that the attacker deems least likely to be training members, with the exclusion threshold modulating privacy strength. Finally, **Anti-Memorization Guidance (AMG)** [5] steers the sampling trajectory away from training examples when high similarity is detected, with its guidance scale controlling privacy strength.

B.2. Privacy Preserving Visualization

Figure 5 visualizes synthetic–private image pairs under different methods to intuitively assess reconstruction risk. DP, De-Duplication, and Masking introduce slight varia-

tion, with key visual features, such as color, shape, or layout—remain similar. In contrast, our method produces synthetic images that differ significantly from their private counterparts. We draw the conclusion that our method effectively mitigates this risk—synthetic images show strong visual divergence, highlighting the success of our method in reducing memorization.

B.3. Ablation Study

Effect of Data Selection Ratio. We further investigate how the proportion of selected synthetic data affects the trade-off between utility and privacy. Table 3 reports results on ImageNet and DomainNet datasets under varying selection ratios. As expected, increasing the number of selected samples generally improves downstream utility. However, this utility gain comes at the cost of increased privacy risk. In parallel, LPIPS scores, which reflect resistance to reconstruction attacks, show a consistent decline as more data is retained. This suggests that although larger subsets offer better task performance, they also include more memorized samples, thus weakening privacy protection.

Dataset	Ratio	Utility	MIA %	LPIPS
ImageNet (1000 Data)	20%	53.82 ± 0.30	1.73 ± 0.37	0.5480 ± 0.01
	50%	64.13 ± 0.41	2.08 ± 0.74	0.5300 ± 0.01
	80%	68.47 ± 0.45	19.90 ± 2.40	0.5265 ± 0.00
	100%	69.43 ± 0.58	31.00 ± 0.51	0.4676 ± 0.02
DomainNet (2000 Data)	12.5%	63.80 ± 0.48	0.33 ± 0.31	0.5190 ± 0.02
	25%	70.80 ± 0.59	1.91 ± 0.71	0.5133 ± 0.02
	50%	75.30 ± 0.93	12.76 ± 1.85	0.4208 ± 0.03
	100%	78.37 ± 0.04	48.08 ± 1.06	0.3328 ± 0.01

Table 3. **Utility–Privacy Trade-off Under Varying Data Selection Ratios.** More selected samples lead to higher utility but increase privacy leakage, as evidenced by rising MIA% and decreasing LPIPS scores.

Effect of Privacy Penalties α_1 and α_2 . We conduct an ablation study to examine the effect of the penalty hyper coefficients α_1 and α_2 in Eq. (12). Table 4 presents results on ImageNet and DomainNet. On ImageNet, increasing α_1 from 1.0 to 50.0 leads to a notable decrease in MIA success rate. This demonstrates that stronger penalization of MIA risk is effective in reducing privacy leakage. On DomainNet, we vary α_2 while keeping α_1 fixed at 50.0. We observe that larger α_2 values lead to higher LPIPS scores, indicating better resistance to reconstruction attacks.

Effectiveness of Data Selection. We conduct an ablation study to assess the impact of data selection on utility. As shown in Table 5, directly filtering out privacy sensitive data leads to a substantial utility drop. For instance, in the Im-

Dataset	α'_1	α'_2	Utility	MIA %	LPIPS
ImageNet (1000 Data)	1.0	1.0	64.13 ± 0.41	2.08 ± 0.54	0.5300 ± 0.01
	5.0	1.0	62.97 ± 0.60	0.83 ± 0.25	0.5313 ± 0.00
	10.0	1.0	62.47 ± 0.66	0.59 ± 0.30	0.5289 ± 0.01
	50.0	1.0	60.98 ± 0.58	0.15 ± 0.16	0.5333 ± 0.01
DomainNet (2000 Data)	50.0	1.0	71.30 ± 0.59	2.29 ± 0.58	0.4648 ± 0.02
	50.0	5.0	71.67 ± 0.60	2.86 ± 0.46	0.4975 ± 0.01
	50.0	10.0	70.88 ± 1.15	2.47 ± 0.83	0.5228 ± 0.01
	50.0	50.0	68.30 ± 0.87	22.05 ± 1.13	0.5626 ± 0.01

Table 4. **Effect of Privacy Penalty Weights α_1 and α_2 on Utility and Privacy.** On ImageNet (top), increasing α_1 steadily lowers MIA success rate; on DomainNet (bottom), raising α_2 boosts LPIPS (reconstruction resistance) but can slightly reduce utility.

geNet dataset, accuracy decreases by 5.4% when filtering is applied without selection.

Dataset	Method	Epoch							
		10	20	50	100	150	200	250	300
ImageNet	Filter Out	67.6	58.9	49.7	45.8	44.6	39.7	40.5	40.7
	Our Method	61.9	53.3	43.9	43.1	37.7	36.2	35.4	35.3
DomainNet	Filter Out	58.6	49.8	41.2	35.1	33.8	33.6	33.7	33.5
	Our Method	55.4	47.9	35.0	33.3	31.9	28.4	27.9	28.3

Table 5. **Utility Loss over Training Epochs.** Comparison between our data selection method and MIA-based filtering method.