

Re-thinking Model Inversion Attacks Against Deep Neural Networks

Ngoc-Bao Nguyen* Keshigeyan Chandrasegaran* Milad Abdollahzadeh Ngai-Man Cheung[†]
 Singapore University of Technology and Design (SUTD)

thibaongoc_nguyen@mymail.sutd.edu.sg, {keshigeyan, milad_abdollahzadeh, ngaiman_cheung}@sutd.edu.sg

Abstract

Model inversion (MI) attacks aim to infer and reconstruct private training data by abusing access to a model. MI attacks have raised concerns about the leaking of sensitive information (e.g. private face images used in training a face recognition system). Recently, several algorithms for MI have been proposed to improve the attack performance. In this work, we revisit MI, study two fundamental issues pertaining to all state-of-the-art (SOTA) MI algorithms, and propose solutions to these issues which lead to a significant boost in attack performance for all SOTA MI. In particular, our contributions are two-fold: 1) We analyze the optimization objective of SOTA MI algorithms, argue that the objective is sub-optimal for achieving MI, and propose an improved optimization objective that boosts attack performance significantly. 2) We analyze “MI overfitting”, show that it would prevent reconstructed images from learning semantics of training data, and propose a novel “model augmentation” idea to overcome this issue. Our proposed solutions are simple and improve all SOTA MI attack accuracy significantly. E.g., in the standard CelebA benchmark, our solutions improve accuracy by 11.8% and achieve for the first time over 90% attack accuracy. Our findings demonstrate that there is a clear risk of leaking sensitive information from deep learning models. We urge serious consideration to be given to the privacy implications. Our code, demo, and models are available at https://ngoc-nguyen-0.github.io/re-thinking_model_inversion_attacks/.

1. Introduction

Privacy of deep neural networks (DNNs) has attracted considerable attention recently [2, 3, 31, 41, 42]. Today, DNNs are being applied in many domains involving private and sensitive datasets, e.g., healthcare, and security. There is a growing concern of privacy attacks to gain knowledge of confidential datasets used in training DNNs. One

important category of privacy attacks is Model Inversion (MI) [7, 10, 13, 14, 22, 47, 49, 52, 53] (Fig. 1). Given access to a model, MI attacks aim to infer and reconstruct features of the private dataset used in the training of the model. For example, a malicious user may attack a face recognition system to reconstruct sensitive face images used in training. Similar to previous work [7, 47, 52], we will use face recognition models as the running example.

Related Work. MI attacks were first introduced in [14], where simple linear regression is the target of attack. Recently, there is a fair amount of interest to extend MI to complex DNNs. Most of these attacks [7, 47, 52] focus on the whitebox setting and the attacker is assumed to have complete knowledge of the model subject to attack. As many platforms provide downloading of entire trained DNNs for users [7, 52], whitebox attacks are important. [52] proposes Generative Model Inversion (GMI) attack, where generic public information is leveraged to learn a distributional prior via generative adversarial networks (GANs) [15, 45], and this prior is used to guide reconstruction of private training samples. [7] proposes Knowledge-Enriched Distributional Model Inversion (KEDMI), where an inversion-specific GAN is trained by leveraging knowledge provided by the target model. [47] proposes Variational Model Inversion (VMI), where a probabilistic interpretation of MI leads to a variational objective for the attack. KEDMI and VMI achieve SOTA attack performance (See Supplementary E for further discussion of related work).

In this paper, we revisit SOTA MI, study two issues pertaining to all SOTA MI and propose solutions to these issues that are complementary and applicable to all SOTA MI (Fig. 1). In particular, despite the range of approaches proposed in recent works, common and central to all these approaches is an *inversion step* which formulates reconstruction of training samples as an optimization. The optimization objective in the inversion step involves the *identity loss*, which is the *same* for all SOTA MI and is formulated as the negative log-likelihood for the reconstructed samples under the model being attacked. While ideas have been proposed to advance other aspects of MI, *effective design of the identity loss has not been studied*.

*Equal Contribution †Corresponding Author

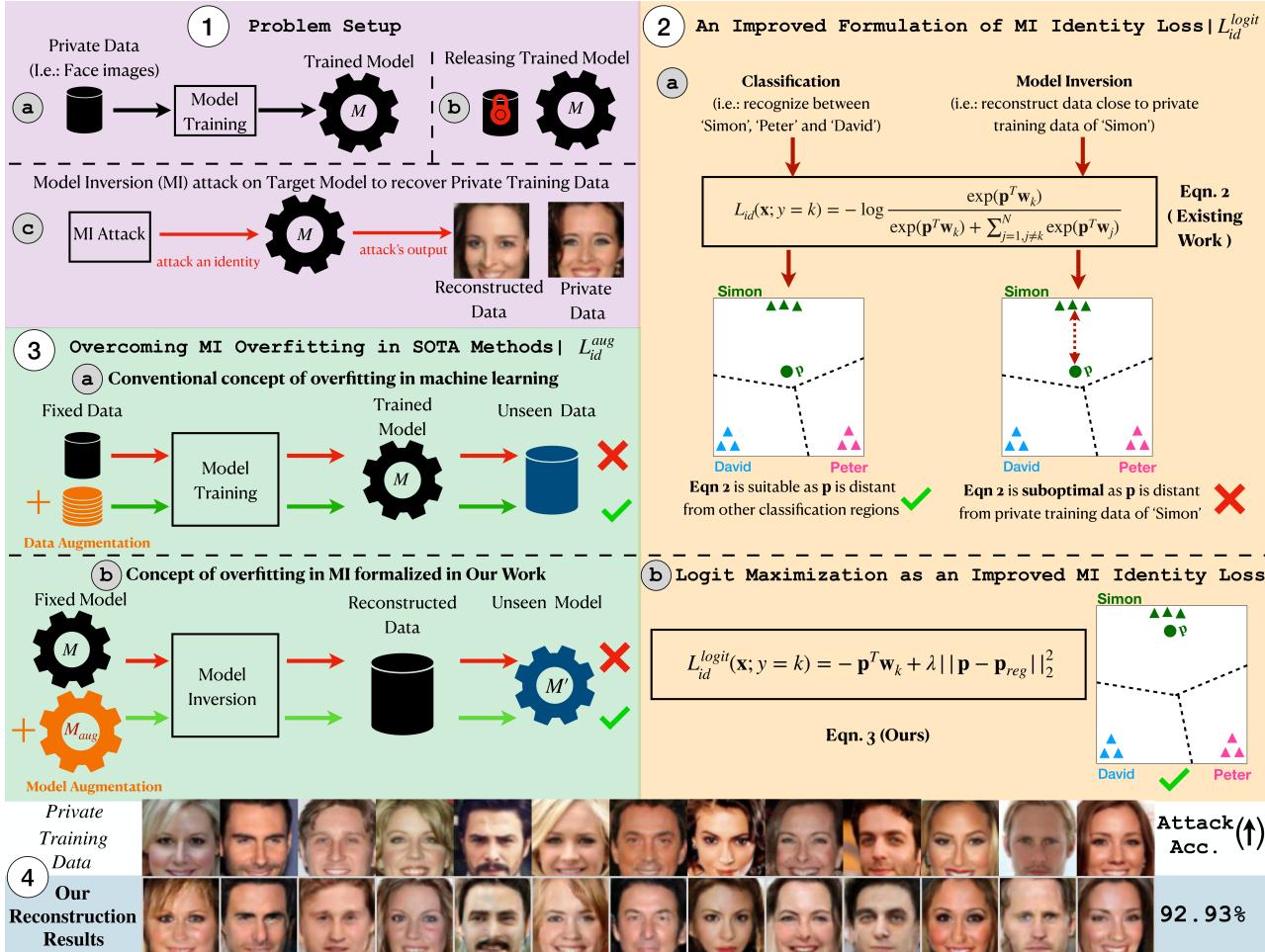


Figure 1. Overview and our contributions. ① We consider the problem of the Model Inversion (MI) attack to reconstruct private training data based on model parameters. Our work makes two foundational contributions to MI attacks. ② First, we analyse the optimization objective of existing SOTA MI algorithms and show that they are sub-optimal. Further, we propose an improved optimization objective that boosts MI attack performance significantly (Sec 3.1). ③ Second, we formalize the concept of “MI overfitting” showing that it prevents reconstructed images from learning identity semantics of training data. Further, we propose a novel “model augmentation” idea to overcome this issue (Sec 3.2). ④ Our proposed method significantly boosts MI attack accuracy. *E.g.* In the standard CelebA benchmark, our method boosts attack accuracy by 11.8%, achieving above 90% attack accuracy for the first time in contemporary MI literature.

To address this research gap, our work studies subtleties of identity loss in all SOTA MI, analyzes the issues and proposes improvements that boost the performance of all SOTA significantly. In summary, our contributions are as follows:

- We analyze existing identity loss, argue that it could be sub-optimal for MI, and propose an improved identity loss that aligns better with the goal of MI (Fig. 1 ②).
 - We formalize the concept of *MI overfitting*, analyze its impact on MI and propose a novel solution based on *model augmentation*. Our idea is inspired by the conventional issue of overfitting in model training and data augmentation as a solution to alleviate the issue (Fig. 1 ③).
 - We conduct extensive experiments to demonstrate that

our solutions can improve SOTA MI algorithms (GMI [52], KEDMI [7], VMI [47]) significantly. Our solutions achieve for the first time over 90% attack accuracy under standard CelebA benchmark (Fig. 1 ④).

Our work sounds alarm over the rising threats of MI attacks, and urges more attention on measures against the leaking of private information from DNNs.

2. General Framework of SOTA MI Attacks

Problem Setup. In MI, an attacker abuses access to a model M trained on a private dataset \mathcal{D}_{priv} . The attacker can access M , but \mathcal{D}_{priv} is not intended to be shared. The goal of MI is to infer information about private samples in \mathcal{D}_{priv} . In existing work, for the desired class (label) y ,

Table 1. Categorizing SOTA MI attacks based on their difference in latent code distribution and prior loss. $p_{\text{GAN}}(\mathbf{z})$ is a GAN prior. G and D are generator and discriminator of a GAN.

Method	Latent distribution $q(\mathbf{z})$	Prior loss L_{prior}
GMI [52]	Point estimate $\delta(\mathbf{z} - \mathbf{z}_0)$	$-D(G(\mathbf{z}))$
KEDMI [7]	Gaussian $\mathcal{N}(\mu, \Sigma)$	$-\log D(G(\mathbf{z}))$
VMI [47]	Gaussian $\mathcal{N}(\mu, \Sigma)$ or Normalizing Flow [26]	Distance w.r.t. GAN prior $D_{\text{KL}}(q(\mathbf{z}) p_{\text{GAN}}(\mathbf{z}))$

MI is formulated as the reconstruction of an input \mathbf{x} which is most likely classified into y by the model M . For instance, if the problem involves inverting a facial recognition model, given the desired identity, MI is formulated as the reconstruction of facial images that are most likely to be recognized as the desired identity. The model subject to MI attacks is called *target model*. Following previous works [7, 47, 52], we focus on *whitebox* MI attack, where the attacker is assumed to have complete access to the target model. For high-dimensional data such as facial images, this reconstruction problem is ill-posed. Consequently, various SOTA MI methods have been proposed recently to constrain the search space to the manifold of meaningful and relevant images using a GAN: using a GAN trained on some public dataset \mathcal{D}_{pub} [52], using an inversion-specific GAN [7], and defining variational inference in latent space of GAN [47].

Despite the differences in various SOTA MI, common and central to all these methods is an *inversion step* –called *secret revelation* in [52]–, which performs the following optimization:

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z})} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \{L_{\text{id}}(\mathbf{z}; y, M) + \lambda L_{\text{prior}}(\mathbf{z})\} \quad (1)$$

Here $L_{\text{id}}(\mathbf{z}; y, M) = -\log \mathbb{P}_M(y|G(\mathbf{z}))$ is referred to as *identity loss* in MI [52], which guides the reconstruction of $\mathbf{x} = G(\mathbf{z})$ that is most likely to be recognized by model M as identity y , and L_{prior} is some prior loss, and $q^*(\mathbf{z})$ is the optimal distribution of latent code used to generate inverted samples by GAN ($\mathbf{x} = G(\mathbf{z})$; $\mathbf{z} \sim q^*(\mathbf{z})$). Importantly, all SOTA MI methods use the *same* identity loss $L_{\text{id}}(\mathbf{z}; y, M)$, although they have different assumption about $q(\mathbf{z})$ and the prior loss L_{prior} (see Table 1 and Supplementary for more details on each algorithm). While advances observed by improving $q(\mathbf{z})$ and L_{prior} , **the design of more effective L_{id} has been left unnoticed** in all SOTA MI algorithms. Therefore, our work instead focuses on L_{id} , analyzes issues, and proposes improvement for L_{id} that can lead to a performance boost in all SOTA MI. To simplify notations, we denote $L_{\text{id}}(\mathbf{z}; y, M)$ by $L_{\text{id}}(\mathbf{x}; y)$ when appropriate, where

$\mathbf{x} = G(\mathbf{z})$ is the reconstructed image.

3. A Closer Look at Model Inversion Attacks

3.1. An Improved Formulation of MI Identity Loss

In this section, we discuss our first contribution and take a closer look at the optimization objective of *identity loss*, $L_{\text{id}}(\mathbf{x}; y)$. Existing SOTA MI methods, namely GMI [52], KEDMI [7] and VMI [47] formulate the identity loss as an optimization to minimize the negative log likelihood of an identity under model parameters (*i.e.* cross-entropy loss). Particularly, the $L_{\text{id}}(\mathbf{x}; y)$ introduced in Eqn. 1 for an inversion targeting class k can be re-written as follows:

$$L_{\text{id}}(\mathbf{x}; y = k) = -\log \frac{\exp(\mathbf{p}^T \mathbf{w}_k)}{\exp(\mathbf{p}^T \mathbf{w}_k) + \sum_{j=1, j \neq k}^N \exp(\mathbf{p}^T \mathbf{w}_j)} \quad (2)$$

where \mathbf{p} refers to penultimate layer activations [6, 35] for sample \mathbf{x} and \mathbf{w}_i refers to the last layer weights for the i^{th} class ¹ in target model M .

Existing identity loss (Eqn. 2) used in SOTA MI methods [7, 47, 52] is sub-optimal for MI (Fig. 1 ②). Although the optimization in Eqn. 2 accurately captures the essence of a classification problem (*e.g.* face recognition), we postulate that such formulation is sub-optimal for MI. We provide our intuition through the lens of penultimate layer activations, \mathbf{p} (Fig. 1 ②). In a classification setting, the main expectation for \mathbf{p} is to be sufficiently discriminative for class k (*e.g.* recognize between ‘Peter’, ‘Simon’ and ‘David’). This objective can be achieved by both maximizing $\exp(\mathbf{p}^T \mathbf{w}_k)$ and/or minimizing $\sum_{j=1, j \neq k}^N \exp(\mathbf{p}^T \mathbf{w}_j)$ in Eqn. 2. On the contrary, *the goal of MI is to reconstruct training data*. That is, in addition to \mathbf{p} being sufficiently discriminative for class k , successful inversion also requires \mathbf{p} to be close to the training data representations for class k represented by \mathbf{w}_k (*i.e.* an inversion targeting ‘Simon’ needs to reconstruct a sample close to the private training data of ‘Simon’; Fig. 1 ②). Specifically, we argue that MI requires a lot more attention on maximizing $\exp(\mathbf{p}^T \mathbf{w}_k)$ compared to minimizing $\sum_{j=1, j \neq k}^N \exp(\mathbf{p}^T \mathbf{w}_j)$ in Eqn. 2.

Motivated by this hypothesis, we conduct an analysis to investigate the proximity between private training data and reconstructed data in SOTA MI methods using penultimate layer representations [6, 34, 35, 39]. Particularly, our analysis using KEDMI [7] (SOTA) shows several instances where using Eqn. 2 for identity loss is unable to reconstruct data close to the private training data. We show this in Fig. 2 (top row). Consequently, our analysis motivates the search for an improved identity loss focusing on maximizing $\exp(\mathbf{p}^T \mathbf{w}_k)$ for MI.

¹ \mathbf{p} is concatenated with 1 at the end to include bias as \mathbf{w}_i includes biases at the end.

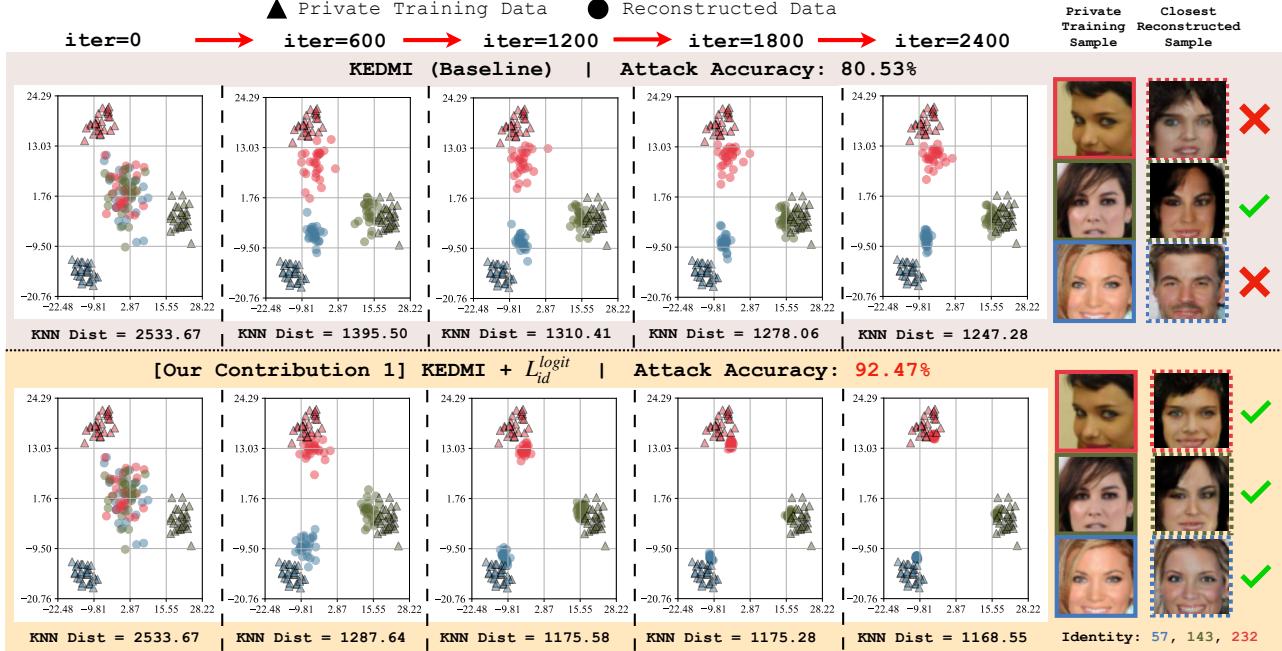


Figure 2. Visualization of the penultimate layer representations ($D_{priv} = \text{CelebA}$ [33], $D_{pub} = \text{CelebA}$ [33], Target Model = IR152 [7], Evaluation Model = face.evoLve [8], Inversion iterations = 2400) for private training data and reconstructed data using KEDMI [7]. Following the exact evaluation protocol in [7], we use face.evoLve [8] to extract representations. We show results for 3 randomly chosen identities. We include KNN distance (for different iterations) and final attack accuracy following the protocol in [7]. For each identity, we also include randomly selected private training data and the closest reconstructed sample at iteration=2400. **① Identity loss in SOTA MI methods [7, 47, 52] (Eqn. 2) is sub-optimal for MI (Top).** Using penultimate representations during inversion, we observe 2 instances (e.g. target identity 57 and 143) where KEDMI [7] (using Eqn. 2 for identity loss) is unable to reconstruct data close to private training data. Hence, private and reconstructed facial images are qualitatively different. **② Our proposed identity loss, L_{id}^{logit} (Eqn. 3), can effectively guide the reconstruction of data close to private training data (Bottom).** This can be clearly observed using both penultimate layer representations and KNN distances for all 3 target classes 57, 143 and 252. We show similar results using additional MI algorithms (GMI [52], VMI [47]) and target classifiers (face.evoLve, VGG16) in Supplementary Figures D.2, D.5 and D.8. Best viewed in color.

Logit Maximization as an improved MI identity loss. In light of our analysis / observations above, we propose to directly maximize the logit, $\mathbf{p}^T \mathbf{w}_k$, instead of maximizing the log likelihood of class k for MI. Our proposed identity loss objective is shown below:

$$L_{id}^{logit}(\mathbf{x}; y = k) = -\mathbf{p}^T \mathbf{w}_k + \lambda \|\mathbf{p} - \mathbf{p}_{reg}\|_2^2 \quad (3)$$

where $\lambda (> 0)$ is a hyper-parameter and \mathbf{p}_{reg} is used for regularizing \mathbf{p} . Particularly, if the regularization in Eqn. 3 is omitted and hence $\|\mathbf{p}\|$ is unbounded, a crude simplified way to solve Eqn. 3 is to maximize $\|\mathbf{p}\|$. Hence, we use \mathbf{p}_{reg} to regularize \mathbf{p} . Given that the attacker has no access to private training data, we estimate \mathbf{p}_{reg} by a simple method using *public* data (See Supplementary C.3). We remark that $\mathbf{p} = M^{\text{pen}}(\mathbf{x})$ where $\mathbf{x} = G(\mathbf{z})$ and $M^{\text{pen}}()$ operator returns the penultimate layer representations for a given input.

Our analysis shows that our proposed identity loss, L_{id}^{logit} (Eqn. 3), can significantly improve reconstruction of private training data compared to existing identity loss used in SOTA MI algorithms [7, 47, 52]. This can be clearly observed using both penultimate layer representations and

KNN distances in Fig. 2 (bottom row). Here KNN Dist refers to the shortest Euclidean feature distance from a reconstructed image to private training images for a given identity [7, 52]. Our proposed L_{id}^{logit} can be easily plugged in to all existing SOTA MI algorithms by replacing L_{id} with our proposed L_{id}^{logit} in Eqn. 1 (in the inversion step) with minimal computational overhead.

3.2. Overcoming MI Overfitting in SOTA methods

In this section, we discuss our second contribution. In particular, we formalize a concept of *MI overfitting*, observe its considerable impacts even in SOTA MI methods [7, 47, 52], and propose a new, simple solution to overcome this issue (Fig. 1 ③). To better discuss our MI overfitting concept, we first review the conventional concept of overfitting in machine learning: Given the fixed training dataset and our goal of learning a model, conventionally, overfitting is defined as instances which during model learning (training stage), the model fits too closely to the training data and adapts to the random variation and noise of training data, failing to adequately learn the semantics of the train-

ing data [1, 38, 44, 50, 54]. As the model lacks semantics of training data, it could be observed that the model performs poorly under unseen data (Fig. 1 ③ ④).

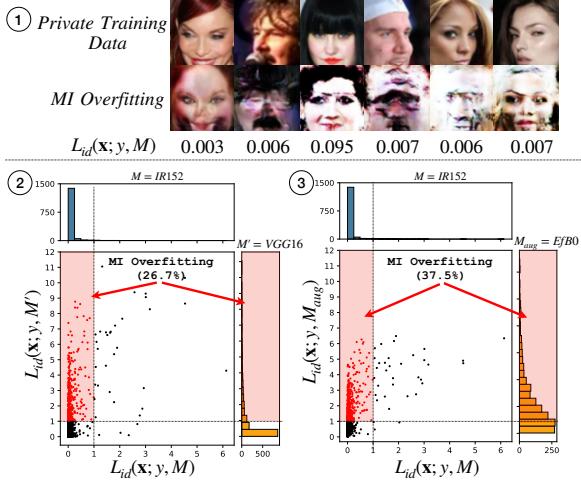


Figure 3. Qualitative / Quantitative studies to demonstrate MI overfitting in SOTA methods. We demonstrate this observation using KEDMI [7]. We use $\mathcal{D}_{priv} = \text{CelebA}$ [33], $\mathcal{D}_{pub} = \text{CelebA}$ [33] and $M = \text{IR152}$ [7]. **① We show qualitative results to illustrate MI overfitting.** We show 6 identities, top: private data, bottom: reconstructed data from M . The reconstructed samples have fit too closely to M during inversion resulting in samples with lack of identity semantics. Particularly, we remark that these samples have very low identity loss under the target model M . **② Quantitative results validating the prevalence of MI overfitting in SOTA MI methods.** We use an additional target classifier $M' = \text{VGG16}$ released by [7, 52] to quantitatively verify the presence of MI overfitting using identity loss. For 1,500 reconstructed samples from M , we visualize their identity loss w.r.t. M and M' in the scatter plot and respective histograms. Particularly, we find that there are 26.7% of samples with low identity loss under the target model M , but large identity loss under unseen VGG16 model M' , hinting that these samples might lack identity semantics. This result shows that MI overfitting is a considerable issue even in SOTA MI methods. Note that VGG16 is used here only for analysis and is not part of our solution, as private data is not available. **③ Model Augmentation to alleviate MI overfitting during inversion.** We repeat the above analysis, with $M' = \text{VGG16}$ replaced by $M_{aug} = \text{EfficientNet-B0}$. Importantly, M_{aug} is trained by *public data* using knowledge distillation [19]. We similarly observe samples with large identity loss under M_{aug} .

Overfitting in MI. We formalize the concept of overfitting in MI (Fig. 1 ③ ⑤). Given the fixed (target) model and our goal of learning reconstructed samples, we define MI overfitting as instances which during model inversion, the reconstructed samples fit too closely to the target model and adapt to the random variation and noise of the target model parameters, failing to adequately learn semantics of the identity. As these reconstructed samples lack identity semantics, it could be observed that they perform poorly under another unseen model.

Analysis. In what follows, we discuss our analysis to demonstrate MI overfitting and understand its impact in SOTA. See Fig. 3 for analysis setups and results. In particular, in Fig. 3 ①, we show some reconstructed samples which achieve low identity loss under the target model M , yet they lack identity semantics. In Fig. 3 ②, we show that for a considerable percentage of reconstructed samples from target model M with low identity loss under M , their identity loss under another unseen model M' is large as shown in the scatter plot and histograms, hinting that these samples might have suffered from MI overfitting and lack identity semantics. We note that the identity loss under M' is obtained by feeding the reconstructed sample into M' in a forward pass. We also note that SOTA KEDMI [7] is used in this analysis but the issue persists in [47, 52].

Our proposed solution to MI overfitting. We propose a novel solution based on *model augmentation*. Our idea is inspired by the conventional issue of overfitting in model training and data augmentation as a solution to alleviate the issue. In particular, for conventional overfitting, augmenting the training dataset could alleviate the issue [29]. Therefore, we hypothesize that by augmenting the target model we can alleviate MI overfitting.

Specifically, we propose to apply knowledge distillation (KD) [19], with target model M_t as the teacher, to train augmented models $M_{aug}^{(i)}$. Importantly, as we do not have access to the private data, during KD, each $M_{aug}^{(i)}$ is trained on the *public dataset* to match its output to the output of M_t . We select different network architectures for $M_{aug}^{(i)}$ and they are different from M_t (Detailed discussion in the Supplementary B.1 and B.2). After performing KD, we apply $M_{aug}^{(i)}$ together with the target model M_t in the inversion step and compute the identity loss (with model augmentation):

$$\begin{aligned} L_{id}^{aug}(\mathbf{x}; y) &= \gamma_t \cdot L_{id}(\mathbf{x}; y, M_t) \\ &+ \gamma_{aug} \cdot \sum_{i=1}^{N_{aug}} L_{id}(\mathbf{x}; y, M_{aug}^{(i)}) \end{aligned} \quad (4)$$

Here, γ_t and γ_{aug} are two hyper-parameters. In particular, we use $\gamma_t = \gamma_{aug} = \frac{1}{N_{aug}+1}$, where N_{aug} is the number of augmented models. L_{id}^{aug} in Eqn. 4 is used to replace L_{id} in the inversion step in Eqn. 1. Furthermore, our proposed L_{logit} in Eqn. 3 can be used in Eqn. 4 to combine the improvements. See details in Supplementary C.1.

In Fig. 3 ③, we analyze the performance of $M_{aug}^{(i)}$. Similar to using the unseen model M' , we observe samples with large identity loss under $M_{aug}^{(i)}$, suggesting that samples with MI overfitting perform poorly under $M_{aug}^{(i)}$ as these samples lack identity semantic.

Table 2. We follow the exact experiment setups in [7] for GMI [52] and KEDMI [7]. For VMI [47], we follow the exact experiment setups in [47]. In total, we conduct 72 experiments spanning 18 setups to demonstrate the effectiveness of our proposed method.

Method	Private Dataset	Public Dataset	Target model	Evaluation Model	Model Augmentation
GMI [52] / KEDMI [7]	CelebA [33]	CelebA / FFHQ [24]	VGG16 [40] / IR152 [17] / face.evoLve [9]	face.evoLve	EfficientNet-B0 [43], EfficientNet-B1 [43], EfficientNet-B2 [43]
	CIFAR-10 [28]	CIFAR-10	VGG16	ResNet-18 [17]	
	MNIST [30]	MNIST	CNN(Conv3)	CNN(Conv5)	CNN(Conv2), CNN(Conv4)
VMI [47]	CelebA	CelebA	ResNet-34 [17]	IR-SE50 [12]	EfficientNet-B0, EfficientNet-B1, EfficientNet-B2
	MNIST	EMNIST [11]	ResNet-10	ResNet-10	CNN(Conv2), CNN(Conv4)

4. Experiments

In this section, we evaluate the performance of the proposed method in recovering a representative input from the target model, against current SOTA methods: GMI [52], VMI [47], and KEDMI [7]. More specifically, as our proposed method identifies two major limitations in current $L_{id}(\mathbf{x}; y)$ —used commonly in all SOTA MI approaches—we will evaluate the improvement brought by our improved identity loss L_{id}^{logit} , and model augmentation L_{id}^{aug} for all SOTA MI approaches.

4.1. Experimental Setup

In order to have a fair comparison, when evaluating our method against each SOTA MI approach, we follow the exactly same experimental setup of that approach. In what follows, we discuss the details of these setups.

Dataset. Following previous works, we evaluate the proposed method on different tasks: face recognition and digit classification is used for comparison with all three SOTA approaches, and image classification is used for comparison with GMI [52], and KEDMI [7]. For the face recognition task, we use CelebA dataset [33] that includes celebrity images, and the FFHQ dataset [24] which contains images with larger variation in terms of background, ethnicity, and age. The MNIST handwritten digits dataset [30] is used for digit classification. We utilize the CIFAR-10 dataset [28] for image classification.

Data Preparation Protocol. Following previous SOTA approaches [7, 47, 52], we split each dataset into two disjoint parts: one part is used as private dataset \mathcal{D}_{priv} for training target model, and another part is used as a public dataset \mathcal{D}_{pub} to extract the prior information. Most importantly, *throughout all experiments, public dataset \mathcal{D}_{pub} has no class intersection with private dataset \mathcal{D}_{priv} used for training target model*. Note that this is essential to make sure that adversary uses \mathcal{D}_{pub} only to gain prior knowledge about features that are general to that task (i.e., face recognition), and does not have access to information about class-specific and private information used for training target model.

Models. Following previous works, we implement several different models with varied complexities. As GMI [52] and KEDMI [7] use exactly similar model architecture in experiments, for comparison with these two algorithms, we use the same models. More specifically, for face recognition on CelebA and FFHQ, we use VGG16 [40], IR152 [17], and face.evoLve [9]—as SOTA face recognition model. For digit classification on MNIST, we use a CNN with 3 convolutional layers and 2 pooling layers. Finally, for image classification, following [7] we use VGG16 [40]. For a fair comparison with VMI, we follow its design in [47] and use ResNet-34 for face recognition CelebA, and ResNet-10 for digit classification on MNIST. The details of the target models, augmented models and datasets used in experiments are summarized in Table 2. We remark that when comparing our proposed method with each of the SOTA MI approaches, we use exactly the same target model and GAN for both SOTA and our approach.

Evaluation Metrics. To evaluate the performance of a MI attack, we need to assess whether the reconstructed image exposes private information about a target label/identity. In this work, following the literature, we conduct both qualitative evaluations by visual inspection, and quantitative evaluations using different metrics, including:

- **Attack Accuracy (Attack Acc).** Following [7, 47, 52], we use an *evaluation model* that predicts the label/identity of the reconstructed image. Similar to previous works, the evaluation model is different from the target model (different structure/ initialization seed), but it is trained on the same private dataset (see Table 2). Intuitively, considering a highly accurate evaluation model, it can be viewed as a proxy for human inspection [52]. Therefore, if the evaluation model infers high accuracy on reconstructed images, it means these images are exposing private information about the private dataset, i.e. high attack accuracy.
- **K-Nearest Neighbors Distance (KNN Dist).** KNN Dist indicates the distance between the reconstructed image for a specific label/id and corresponding images in the private training dataset. More specifically, it measures the

shortest feature distance from the reconstructed image to the real images in the private dataset, given a class/id. It is measured as l_2 distance between two images in the feature space, i.e., the penultimate layer of the evaluation model.

Private Training Data	KEDMI	Attack Acc. (\uparrow) KNN Dist (\downarrow)	
		Acc.	Dist
Existing SOTA		80.53%	1247.28
+ LOM (Ours)		92.47%	1168.55
+ MA (Ours)		84.73%	1220.23
+ LOMMA (Ours)		92.93%	1138.62

Figure 4. Qualitative / Quantitative (Top1 Attack Acc., KNN Dist) results to demonstrate the efficacy of our proposed method. We use KEDMI [7] (SOTA), $\mathcal{D}_{priv} = \text{CelebA}$ [33], $\mathcal{D}_{pub} = \text{CelebA}$ [33] and $M = \text{IR152}$ [17]. As one can observe, our proposed method achieves better reconstruction of private data both visually and quantitatively (validated by KNN results) resulting in a significant boost in attack performance.

4.2. Experimental Results

Comparison with previous state-of-the-art. We use GMI [52], KEDMI [7], and VMI [47] as SOTA MI baselines. We reproduce all baseline results using official public implementations. We report results for GMI and KEDMI for CelebA/ CelebA experiments in Table 3. We report VMI results for CelebA/ CelebA experiments in Table 4. For each baseline setup, we report results for 3 variants: • *LOM* (Logit Maximization, Sec. 3.1), • *MA* (Model Augmentation, Sec. 3.2), • *LOMMA* (Logit Maximization + Model Augmentation). The details are as follows:

1. + LOM (Ours): We replace existing identity loss, L_{id} with our improved identity loss L_{id}^{logit} (Sec. 3.1).
2. + MA (Ours): We replace existing identity loss, L_{id} with our proposed L_{id}^{aug} (Sec. 3.2).
3. + LOMMA (Ours): We combine both L_{id}^{logit} (Sec. 3.1) and L_{id}^{aug} (Sec. 3.2) for model inversion.

As one can clearly observe from Table 3 and Table 4, our proposed methods yield significant improvement in MI attack accuracy in *all experiment setups* showing the efficacy of our proposed methods. Further, by combining both our proposed methods, we significantly boost attack accuracy. The KNN results also clearly show that our proposed methods are able to reconstruct data close to the private training data compared to existing SOTA MI algorithms. Particularly, we improve the KEDMI baseline [7] attack accuracy by 12.4% under IR152 target classifier. We show private

Table 3. We report the results for KEDMI and GMI for IR152, face.evoLve and VGG16 target model. Following exact experiment setups in [7], here $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, evaluation model = face.evoLve. We report top 1 accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance. Top 5 attack accuracies are included in the Supplementary A.2. The best results are in **bold**. By alleviating both these major problems in MI algorithms, we achieve new SOTA MI performance (face.evoLve: 81.40% \rightarrow **93.20%**).

Method	Attack Acc \uparrow	Imp. \uparrow	KNN Dist \downarrow
CelebA/CelebA/IR152			
KEDMI	80.53 ± 3.86	-	1247.28
+ LOM (Ours)	92.47 ± 1.41	11.94	1168.55
+ MA (Ours)	84.73 ± 3.76	4.20	1220.23
+ LOMMA (Ours)	92.93 ± 1.15	12.40	1138.62
GMI	30.60 ± 6.54	-	1609.29
+ LOM (Ours)	78.53 ± 3.41	47.93	1289.62
+ MA (Ours)	61.20 ± 4.34	30.60	1389.99
+ LOMMA (Ours)	82.40 ± 4.37	51.80	1254.32
CelebA/CelebA/face.evoLve			
KEDMI	81.40 ± 3.25	-	1248.32
+ LOM (Ours)	92.53 ± 1.51	11.13	1183.76
+ MA (Ours)	85.07 ± 2.71	3.67	1222.02
+ LOMMA (Ours)	93.20 ± 0.85	11.80	1154.32
GMI	27.07 ± 6.72	-	1635.87
+ LOM (Ours)	61.67 ± 4.92	34.60	1405.35
+ MA (Ours)	74.13 ± 4.32	47.06	1352.25
+ LOMMA (Ours)	82.33 ± 3.51	55.26	1257.50
CelebA/CelebA/VGG16			
KEDMI	74.00 ± 3.10	-	1289.88
+ LOM (Ours)	89.07 ± 1.46	15.07	1218.46
+ MA (Ours)	82.00 ± 3.85	8.00	1248.33
+ LOMMA (Ours)	90.27 ± 1.36	16.27	1147.41
GMI	19.07 ± 4.47	-	1715.60
+ LOM (Ours)	69.67 ± 4.80	50.60	1363.81
+ MA (Ours)	51.73 ± 6.03	32.66	1467.68
+ LOMMA (Ours)	77.60 ± 4.64	58.53	1296.26

training data and reconstructed samples for KEDMI [7] under IR152 target model including all 3 variants in Fig. 4. We remark that in the standard CelebA benchmark, our method boosts attack accuracy significantly thereby achieving more than 90% attack accuracy (Table 3) for the first time in contemporary MI literature. We also include CIFAR-10, MNIST and additional results in Supplementary A.1.

Cross-dataset. Following [7], we conduct a series of experiments to study the effect of distribution shift between public and private data on attack performance and KNN distance. We use FFHQ [24] as the public dataset. In particular, we use FFHQ as public data for CelebA experiments. We train GAN models and three model augmentations using the public data. We remark that such setups closely replicate real-world MI attack scenario. We report top 1 accuracy and KNN distance for IR152, face.evoLve, and VGG16 target classifiers in Table 6. It is well known that baseline

Table 4. We follow exact the experiment setup of [47] for VMI experiments. Specifically, we use StyleGAN [23] and Flow model [26] to learn the distribution of z . The best results are in **bold**. Following exact experiment setups in [47], here \mathcal{D}_{priv} = CelebA, \mathcal{D}_{pub} = CelebA, target model = ResNet-34, evaluation model = IRSE50. We report top 1 attack accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance (KNN Dist). The top 5 attack accuracies are included in the Supplementary. The best results are in **bold**. By alleviating both these major problems in MI algorithms, we improve the attack accuracy by 14.94% (59.96% \rightarrow **74.90%**).

Method	Attack Acc \uparrow	Imp. \uparrow	KNN Dist \downarrow
CelebA/CelebA/ResNet-34			
VMI	59.96 \pm 0.27	-	1.144
+ LOM (Ours)	68.34 \pm 0.36	8.38	1.131
+ MA (Ours)	64.16 \pm 0.27	4.20	1.140
+ LOMMA (Ours)	74.90 \pm 0.34	14.94	1.109

Table 5. Results for SOTA defense model BiDO-HSIC [36]: Following exact experiment setups in BiDO-HSIC, \mathcal{D}_{priv} = CelebA, \mathcal{D}_{pub} = CelebA, evaluation model = face.evoLve, target model = BiDO-HSIC. We report top 1 attack accuracies (Attack Acc.), and KNN distance (KNN Dist).

Method	GMI		KEDMI	
	Attack Acc \uparrow	KNN Dist \downarrow	Attack Acc \uparrow	KNN Dist \downarrow
No Def.	19.07 \pm 4.47	1715.60	74.00 \pm 3.10	1289.88
Def. Model	5.20 \pm 2.75	1962.58	42.80 \pm 5.02	1469.75
+ LOM (Ours)	55.80 \pm 3.64	1397.05	64.33 \pm 1.82	1360.57
+ MA (Ours)	23.93 \pm 5.50	1634.84	49.27 \pm 4.02	1413.81
+ LOMMA (Ours)	62.13 \pm 4.04	1358.54	70.47 \pm 2.36	1293.25

attack performances will degrade due to distribution shift between public and private data [7]. But we remark that our proposed methods consistently improves the baseline SOTA attack performances. *i.e.* Our method boosts the attack accuracy of IR152 target model from 52.87% \rightarrow 77.27%.

MI under SOTA defense models. We further evaluate our method on SOTA MI defense models provided by BiDO-HSIC [36]. Specifically, we use the exact GAN and defense models provided by BiDO-HSIC which are trained on CelebA dataset. We then transfer knowledge from the defense model to $M_{aug} = \{\text{Efficientnet-B0}, \text{Efficientnet-B1}, \text{Efficientnet-B2}\}$ using D_{pub} . Results using GMI and KEDMI are shown in Table 5. We observe that SOTA defense BiDO-HSIC is rather ineffective for our proposed MI.

5. Discussion

Conclusion. We revisit SOTA MI and study two issues pertaining to all SOTA MI approaches. First, we analyze existing identity loss in SOTA and argue that it is sub-optimal for MI. We propose a new logit based identity loss that aligns better with the goal of MI. Second, we formalize the concept of MI overfitting and show that it has a considerable impact even in SOTA. Inspired by conventional data augmentation, we propose model augmentation to alleviate MI overfitting. Extensive experiments demonstrate that

Table 6. We report the results for KEDMI and GMI for IR152, face.evoLve and VGG16 target model. Here \mathcal{D}_{priv} = CelebA, \mathcal{D}_{pub} = FFHQ, evaluation model = face.evoLve. We report top 1 accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance. Top 5 attack accuracies are included in the Supplementary. The best results are in **bold**. By alleviating both these major problems in MI algorithms, we improve the attack accuracy 24.40% (IR152: 52.87% \rightarrow **77.27%**).

Method	Attack Acc \uparrow	Imp. \uparrow	KNN Dist \downarrow
CelebA/FFHQ/IR152			
KEDMI	52.87 \pm 4.96	-	1418.83
+ LOM (Ours)	67.73 \pm 2.29	14.86	1325.28
+ MA (Ours)	64.13 \pm 4.49	11.26	1373.42
+ LOMMA (Ours)	77.27 \pm 2.01	24.40	1292.80
GMI	17.20 \pm 5.31	-	1701.76
+ LOM (Ours)	56.00 \pm 5.20	38.80	1427.59
+ MA (Ours)	50.80 \pm 6.89	33.60	1462.92
+ LOMMA (Ours)	72.00 \pm 6.62	54.80	1338.35
CelebA/FFHQ/face.evoLve			
KEDMI	51.87 \pm 3.88	-	1440.19
+ LOM (Ours)	69.73 \pm 2.47	17.86	1379.73
+ MA (Ours)	65.73 \pm 3.51	13.86	1379.09
+ LOMMA (Ours)	73.20 \pm 2.24	21.33	1321.00
GMI	14.27 \pm 4.42	-	1744.47
+ LOM (Ours)	47.93 \pm 4.87	33.66	1498.19
+ MA (Ours)	46.07 \pm 4.88	31.80	1500.10
+ LOMMA (Ours)	64.33 \pm 4.69	50.06	1386.33
CelebA/FFHQ/VGG16			
KEDMI	41.27 \pm 3.50	-	1490.09
+ LOM (Ours)	55.07 \pm 1.88	13.80	1438.72
+ MA (Ours)	52.07 \pm 2.92	10.80	1428.77
+ LOMMA (Ours)	62.67 \pm 2.29	21.40	1366.94
GMI	10.93 \pm 3.47	-	1766.27
+ LOM (Ours)	44.40 \pm 5.96	33.47	1508.84
+ MA (Ours)	34.93 \pm 4.52	24.00	1547.93
+ LOMMA (Ours)	58.73 \pm 6.18	47.80	1415.06

our solutions can improve SOTA significantly, achieving for the first time over 90% attack accuracy under the standard benchmark. Our findings highlight rising threats based on MI and prompt serious consideration on privacy of machine learning.

Limitations and Ethical Concerns. We follow previous work in experimental setups. The scale of our experiments is comparable to previous works. Furthermore, extension of our methods for blackbox/ label-only attacks can be considered in future. While our improved MI methods could have negative societal impacts if it is used by malicious users, our work contributes to increased awareness about privacy attacks on DNNs.

Acknowledgements. This research is supported by the National Research Foundation, Singapore under its AI Singapore Programmes (AISG Award No.: AISG2-RP-2021-021; AISG Award No.: AISG2-TC-2022-007). This project is also supported by SUTD project PIE-SGP-AI-2018-01. We thank reviewers for their valuable comments. We also thank Loo Yi and Kelly Kuo for helpful discussion.

References

- [1] Milad Abdollahzadeh, Touba Malekzadeh, and Ngai-Man Cheung. Revisit multimodal meta-learning through the lens of multi-task learning. *Advances in Neural Information Processing Systems*, 34:14632–14644, 2021.
- [2] Brett K Beaulieu-Jones, Zhiwei Steven Wu, Chris Williams, Ran Lee, Sanjeev P Bhavnani, James Brian Byrd, and Casey S Greene. Privacy-preserving generative deep neural networks support clinical data sharing. *Circulation: Cardiovascular Quality and Outcomes*, 12(7):e005122, 2019.
- [3] Hervé Chabanne, Amaury De Wargny, Jonathan Milgram, Constance Morel, and Emmanuel Prouff. Privacy-preserving classification on deep neural network. *Cryptology ePrint Archive*, 2017.
- [4] Keshigeyan Chandrasegaran, Ngoc-Trung Tran, Alexander Binder, and Ngai-Man Cheung. Discovering Transferable Forensic Features for CNN-generated Images Detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, Oct 2022.
- [5] Keshigeyan Chandrasegaran, Ngoc-Trung Tran, and Ngai-Man Cheung. A Closer Look at Fourier Spectrum Discrepancies for CNN-Generated Images Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7200–7209, June 2021.
- [6] Keshigeyan Chandrasegaran, Ngoc-Trung Tran, Yunqing Zhao, and Ngai-Man Cheung. Revisiting label smoothing and knowledge distillation compatibility: What was missing? In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 2890–2916. PMLR, 17–23 Jul 2022.
- [7] Si Chen, Mostafa Kahla, Ruoxi Jia, and Guo-Jun Qi. Knowledge-enriched distributional model inversion attacks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16178–16187, 2021.
- [8] Yu Cheng, Jian Zhao, Zhecan Wang, Yan Xu, Karlekar Jayashree, Shengmei Shen, and Jiashi Feng. Know you at one glance: A compact vector representation for low-shot learning. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 1924–1932, 2017.
- [9] Yu Cheng, Jian Zhao, Zhecan Wang, Yan Xu, Karlekar Jayashree, Shengmei Shen, and Jiashi Feng. Know you at one glance: A compact vector representation for low-shot learning. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1924–1932, 2017.
- [10] Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-only membership inference attacks. In *International conference on machine learning*, pages 1964–1974. PMLR, 2021.
- [11] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 2921–2926. IEEE, 2017.
- [12] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019.
- [13] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pages 1322–1333, 2015.
- [14] Matthew Fredrikson, Eric Lantz, Somesh Jha, Simon Lin, David Page, and Thomas Ristenpart. Privacy in pharmacogenetics: An {End-to-End} case study of personalized warfarin dosing. In *23rd USENIX Security Symposium (USENIX Security 14)*, pages 17–32, 2014.
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [18] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.
- [20] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324, 2019.
- [21] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [22] Mostafa Kahla, Si Chen, Hoang Anh Just, and Ruoxi Jia. Label-only model inversion attacks via boundary repulsion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15045–15053, 2022.
- [23] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Proc. NeurIPS*, 2020.
- [24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [25] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.

- [26] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.
- [27] Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multimodal generation. *arXiv preprint arXiv:2301.13823*, 2023.
- [28] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). *URL http://www.cs.toronto.edu/kriz/cifar.html*, 5(4):1, 2010.
- [29] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- [30] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [31] Joon-Woo Lee, HyungChul Kang, Yongwoo Lee, Woosuk Choi, Jieun Eom, Maxim Deryabin, Eunsang Lee, Junghyun Lee, Donghoon Yoo, Young-Sik Kim, et al. Privacy-preserving machine learning with fully homomorphic encryption for deep neural network. *IEEE Access*, 10:30039–30054, 2022.
- [32] Swee Kiat Lim, Yi Loo, Ngoc Trung Tran, Ngai Man Cheung, Gemma Roig, and Yuval Elovici. Doping: Generative data augmentation for unsupervised anomaly detection with gan. In *18th IEEE International Conference on Data Mining, ICDM 2018*, pages 1122–1127. Institute of Electrical and Electronics Engineers Inc., 2018.
- [33] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015.
- [34] Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. Does label smoothing mitigate label noise? In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6448–6458. PMLR, 13–18 Jul 2020.
- [35] Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [36] Xiong Peng, Feng Liu, Jingfen Zhang, Long Lan, Junjie Ye, Tongliang Liu, and Bo Han. Bilateral dependency optimization: Defending against model-inversion attacks. *KDD*, 2022.
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [38] Claudio Filipi Gonçalves Dos Santos and João Paulo Papa. Avoiding overfitting: A survey on regularization methods for convolutional neural networks. *ACM Computing Surveys (CSUR)*, 54(10s):1–25, 2022.
- [39] Zhiqiang Shen, Zechun Liu, Dejia Xu, Zitian Chen, Kwang-Ting Cheng, and Marios Savvides. Is label smoothing truly incompatible with knowledge distillation: An empirical study. In *International Conference on Learning Representations*, 2021.
- [40] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [41] Warit Sirichotedumrong and Hitoshi Kiya. A gan-based image transformation scheme for privacy-preserving deep neural networks. In *2020 28th European Signal Processing Conference (EUSIPCO)*, pages 745–749. IEEE, 2021.
- [42] Nagesh Subbanna, Matthias Wilms, Anup Tuladhar, and Nils D Forkert. An analysis of the vulnerability of two common deep learning-based medical image segmentation techniques to model inversion attacks. *Sensors*, 21(11):3874, 2021.
- [43] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [44] Christopher TH Teo, Milad Abdollahzadeh, and Ngai-Man Cheung. Fair generative models via transfer learning. *arXiv preprint arXiv:2212.00926*, 2022.
- [45] Ngoc-Trung Tran, Viet-Hung Tran, Ngoc-Bao Nguyen, Trung-Kien Nguyen, and Ngai-Man Cheung. On data augmentation for gan training. *IEEE Transactions on Image Processing*, 30:1882–1897, 2021.
- [46] Ngoc-Trung Tran, Viet-Hung Tran, Ngoc-Bao Nguyen, Trung-Kien Nguyen, and Ngai-Man Cheung. On data augmentation for gan training. *IEEE Transactions on Image Processing*, 30:1882–1897, 2021.
- [47] Kuan-Chieh Wang, Yan Fu, Ke Li, Ashish Khisti, Richard Zemel, and Alireza Makhzani. Variational model inversion attacks. *Advances in Neural Information Processing Systems*, 34:9706–9719, 2021.
- [48] Ziqi Yang, Ee-Chien Chang, and Zhenkai Liang. Adversarial neural network inversion via auxiliary knowledge alignment. *arXiv preprint arXiv:1902.08552*, 2019.
- [49] Ziqi Yang, Jiyi Zhang, Ee-Chien Chang, and Zhenkai Liang. Neural network inversion in adversarial setting via background knowledge alignment. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, pages 225–240, 2019.
- [50] Chaojian Yu, Bo Han, Li Shen, Jun Yu, Chen Gong, Mingming Gong, and Tongliang Liu. Understanding robust overfitting of adversarial training and beyond. In *International Conference on Machine Learning*, pages 25595–25610. PMLR, 2022.
- [51] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 833–842, 2021.
- [52] Yuheng Zhang, Ruoxi Jia, Hengzhi Pei, Wenxiao Wang, Bo Li, and Dawn Song. The secret revealer: Generative model-inversion attacks against deep neural networks. In *Proceed-*

- ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 253–261, 2020.
- [53] Xuejun Zhao, Wencan Zhang, Xiaokui Xiao, and Brian Lim. Exploiting explanations for model inversion attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 682–692, 2021.
- [54] Yunqing Zhao, Keshigeyan Chandrasegaran, Milad Abdollahzadeh, and Ngai-Man Cheung. Few-shot image generation via adaptation-aware kernel modulation. In *Advances in Neural Information Processing Systems*, 2022.

Supplementary Materials

In this supplementary material, we provide additional experiments, analysis, ablation study, and reproducibility details to support our findings. We provide Pytorch code, demo and pre-trained models (target models/ evaluation models/ augmented models) at: https://ngocnguyen-0.github.io/re-thinking_model_inversion_attacks/.

Contents

A Additional experimental results	11
A.1 Experimental results on CIFAR-10 and MNIST	11
A.2 Experimental Results with Additional Metrics	13
B Ablation Study	13
B.1 Different number of augmented models M_{aug}	13
B.2 Different network architectures for M_{aug}	13
B.3 The effect of different sizes of public dataset	13
C Additional analysis and details on experimental setups	13
C.1 Details on combining L_{id}^{logit} and L_{id}^{aug}	13
C.2 Details on improving KEDMI baseline	14
C.3 Additional details on computing \mathbf{P}_{reg}	14
C.4 Details on regularization parameter λ	16
C.5 Computational overhead	16
C.6 Hyperparameters	16
C.7 Dataset	16
D Additional Visualizations	16
D.1 Additional Results for GMI	16
D.2 Penultimate layer visualization for GMI, KEDMI and VMI	17
D.3 Our reconstruction results	17
E Additional Related work	17

A. Additional experimental results

In this section, we provide additional experimental results that are not included in the main paper. More specifically, first, we evaluate the effect of the proposed method on improving SOTA approaches in new tasks including image classification and digit classification. Then, we use alternative metrics for evaluating SOTA MI approaches with and without proposed improvements on identity loss L_{id} . The additional experimental results in this section further support effectiveness of the proposed approach on improving MI attacks.

A.1. Experimental results on CIFAR-10 and MNIST

In Sec. 4.2 of the main paper, we mostly focus on the face recognition task (on the CelebA dataset) and show

that the proposed method significantly improves SOTA approaches by increasing **Attack Acc** (inference accuracy on reconstructed samples by an evaluation model; see Sec. 4.1. of the main paper) and decreasing **KNN Dist** (distance between the reconstructed samples of a specific class/id and corresponding data in the private dataset \mathcal{D}_{priv} ; see Sec. 4.1).

Table A.1. We report top 1 accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance for two experiment setups. Following exact experiment setups in [7]. For CIFAR-10 experiments, $\mathcal{D}_{priv} = \text{CIFAR-10}$, $\mathcal{D}_{pub} = \text{CIFAR-10}$, $M_t = \text{VGG16}$, evaluation model = ResNet-18. For MNIST experiments, $\mathcal{D}_{priv} = \text{MNIST}$, $\mathcal{D}_{pub} = \text{MNIST}$, $M_t = \text{CNN(Conv3)}$, evaluation model = CNN(Conv5). The best results are in **bold**.

Method	Attack Acc \uparrow	Imp. \uparrow	KNN Dist \downarrow
CIFAR-10/CIFAR-10/VGG16			
KEDMI	95.2 ± 7.96	-	78.24
+ LOM (Ours)	100 ± 0	4.80	52.12
+ MA (Ours)	100 ± 0	4.80	53.17
+ LOMMA (Ours)	100 ± 0	4.80	63.41
GMI	43.20 ± 19.80	-	96.11
+ LOM (Ours)	80.80 ± 14.65	37.60	70.47
+ MA (Ours)	80.00 ± 18.01	36.80	93.46
+ LOMMA (Ours)	95.20 ± 7.96	52.00	80.30
MNIST/MNIST/CNN(Conv3)			
KEDMI	46.40 ± 14.65	-	120.99
+ LOM (Ours)	55.20 ± 8.94	8.80	100.18
+ MA (Ours)	75.20 ± 6.57	28.80	72.38
+ LOMMA (Ours)	100.00 ± 0.00	53.60	58.81
GMI	8.00 ± 1.52	-	126.61
+ LOM (Ours)	15.20 ± 15.12	7.20	161.90
+ MA (Ours)	66.40 ± 19.86	58.40	78.38
+ LOMMA (Ours)	80.80 ± 17.38	72.80	83.56

Table A.2. We follow exact the experiment setup of [47] for the VMI experiments. Specifically, we use DCGAN and Flow model to learn the distribution of \mathbf{z} .

Method	Attack Acc \uparrow	Imp. \uparrow	KNN Dist \downarrow
MNIST/EMNIST/ResNet-10			
VMI	94.60 ± 0.13	-	68.53
+ LOM (Ours)	98.60 ± 0.09	4.00	88.13
+ MA (Ours)	98.98 ± 0.02	4.38	58.81
+ LOMMA (Ours)	100.00 ± 0.00	5.40	52.62

In this section, we provide results for other tasks. More specifically, as mentioned in Sec. 4.1, for GMI [52], and KEDMI [7], following their own setup, we use digit classification task MNIST dataset, and object classification task on the CIFAR-10 dataset. For each task, Table A.1 tabulates the performance of the SOTA approach together with three variants of our proposed approach:

1. + LOM (Ours): We replace existing identity loss, L_{id}

with our improved identity loss L_{id}^{logit} (Sec. 3.1).

2. + MA (Ours): We replace existing identity loss, L_{id} with our proposed L_{id}^{aug} (Sec. 3.2).
3. + LOMMA (Ours): We combine both L_{id}^{logit} and L_{id}^{aug} for model inversion.

As one can see, on average each of the proposed solutions drastically improves the SOTA approaches, and combining these two solutions works even better.

Table A.3. We report the results for KEDMI and GMI for IR152, face.evoLve and VGG16 target model. Following exact experiment setups in [7], here $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, evaluation model = face.evoLve. We report top-5 accuracies, the improvement compared to the SOTA MI (Imp.), and FID scores.

Method	Top-5 Attack Acc \uparrow	Imp. \uparrow	FID \downarrow
CelebA/CelebA/IR152			
KEDMI	98.00 ± 1.96	-	28.06
+ LOM (Ours)	98.67 ± 0.00	0.67	39.03
+ MA (Ours)	98.33 ± 1.19	0.33	28.38
+ LOMMA (Ours)	98.67 ± 0.37	0.67	36.78
GMI	55.67 ± 7.14	-	57.11
+ LOM (Ours)	93.00 ± 3.41	37.33	48.87
+ MA (Ours)	89.00 ± 4.10	33.33	45.24
+ LOMMA (Ours)	97.67 ± 2.41	42.00	45.02
CelebA/CelebA/face.evoLve			
KEDMI	97.33 ± 1.73	-	31.26
+ LOM (Ours)	99.33 ± 0.18	2.00	42.45
+ MA (Ours)	98.00 ± 0.94	0.67	32.08
+ LOMMA (Ours)	99.33 ± 0.33	2.00	38.69
GMI	45.33 ± 8.05	-	59.76
+ LOM (Ours)	84.33 ± 4.49	39.00	44.27
+ MA (Ours)	92.00 ± 2.25	46.67	51.15
+ LOMMA (Ours)	93.67 ± 2.42	48.33	44.07
CelebA/CelebA/VGG16			
KEDMI	93.33 ± 3.36	-	25.46
+ LOM (Ours)	99.00 ± 0.18	5.67	34.45
+ MA (Ours)	95.33 ± 1.60	2.00	24.65
+ LOMMA (Ours)	98.00 ± 0.61	4.67	33.91
GMI	40.33 ± 4.74	-	58.03
+ LOM (Ours)	89.33 ± 2.73	49.00	46.40
+ MA (Ours)	81.33 ± 5.88	41.00	44.90
+ LOMMA (Ours)	95.67 ± 2.16	55.34	43.21

Additionally, as mentioned in Sec. 4.1, for VMI [47], following the setup in [47], we evaluate its performance for digit classification on MNIST, and improvement brought by the proposed method. Note that for a fair comparison, following VMI implementation in [47], in this experiment we use EMNIST [11] as public dataset \mathcal{D}_{pub} to acquire prior knowledge. Results are shown in Table A.2 for three variants of our proposed method, which indicates better performance in terms of both attack accuracy (reaching 100% attack accuracy) and decreasing KNN Distance.

A.2. Experimental Results with Additional Metrics

As mentioned in Sec. 4.1 of the main paper, Attack Acc and KNN Dist are common metrics used in literature to evaluate the MI attacks. In this section, we include results on two additional metrics namely: Top-5 Attack Acc and FID [18]. Results in Table A.3, Table A.4, and Table A.5 show that the proposed method achieves better performance in terms of Top-5 Attack Acc, and FID value.

Table A.4. We report the results for VMI. Following exact experiment setups in [47], here $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, $M_t = \text{ResNet-34}$, evaluation model = IR-SE50. We report top-5 accuracies, the improvement compared to the SOTA MI (Imp.), and FID scores.

Method	Top-5 Attack Acc \uparrow	Imp. \uparrow	FID \downarrow
CelebA/CelebA/ResNet-34			
VMI	82.32 ± 0.21	-	16.82
+ LOM (Ours)	86.56 ± 0.27	4.24	25.42
+ MA (Ours)	86.16 ± 0.19	3.84	17.60
+ LOMMA (Ours)	91.02 ± 0.22	8.70	23.56

B. Ablation Study

B.1. Different number of augmented models M_{aug}

In Sec 3.2, we propose a model augmentation idea with augmented models M_{aug} . Here, we experiment using a different number of networks for M_{aug} . Table B.1 show that increasing the number of the augmented models will improve attack accuracy. We use 3 augmented models in our main result as this configuration achieves a good tradeoff in accuracy and computation.

B.2. Different network architectures for M_{aug}

In this section, we provide additional results by using different structures for augmenting the target model using M_{aug} in the MI process. Note that the architecture of all these models is different from the one used for target model M_t .

More specifically, we use three different combinations for M_{aug} , each of which contains three models: (i) {EfficientNet-B0, EfficientNet-B1, EfficientNet-B2}, and (ii) {DenseNet121, DenseNet161, DenseNet169}, and (iii) {EfficientNet-B0, DenseNet121, MobileNetV3}. Results in Table B.2 shows that +MA (Ours) consistently improves the attack accuracy and KNN distance with different network architectures.

B.3. The effect of different sizes of public dataset

We conduct additional experiments using different sizes of \mathcal{D}_{pub} (10%, 50%) to emulate the different quality of prior

Table A.5. We report the results for KEDMI and GMI for IR152, face.evoLve and VGG16 target model. Following exact experiment setups in [7], here $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{FFHQ}$, evaluation model = face.evoLve. We report top-5 accuracies, the improvement compared to the SOTA MI (Imp.), and FID scores.

Method	Top-5 Attack Acc \uparrow	Imp. \uparrow	FID \downarrow
CelebA/FFHQ/IR152			
KEDMI	85.33 ± 4.01	-	41.71
+ LOM (Ours)	88.67 ± 1.18	3.33	50.84
+ MA (Ours)	87.67 ± 2.28	2.33	39.88
+ LOMMA (Ours)	92.00 ± 0.57	6.60	45.67
GMI	36.33 ± 3.98	-	47.72
+ LOM (Ours)	80.33 ± 4.21	44.00	40.18
+ MA (Ours)	84.00 ± 5.35	47.67	35.41
+ LOMMA (Ours)	90.33 ± 3.16	54.00	37.58
CelebA/FFHQ/face.evoLve			
KEDMI	80.67 ± 2.83	-	38.09
+ LOM (Ours)	91.33 ± 0.47	10.67	47.30
+ MA (Ours)	88.67 ± 2.44	8.00	35.94
+ LOMMA (Ours)	94.00 ± 0.68	13.33	47.51
GMI	33.33 ± 6.18	-	52.84
+ LOM (Ours)	74.67 ± 4.78	41.33	44.01
+ MA (Ours)	72.00 ± 4.64	38.67	35.58
+ LOMMA (Ours)	89.00 ± 2.73	55.67	40.03
CelebA/FFHQ/VGG16			
KEDMI	74.00 ± 4.05	-	36.18
+ LOM (Ours)	81.67 ± 1.19	7.67	43.76
+ MA (Ours)	80.33 ± 3.27	6.33	35.02
+ LOMMA (Ours)	85.33 ± 1.98	11.33	40.26
GMI	25.67 ± 5.13	-	53.17
+ LOM (Ours)	70.67 ± 3.92	45.00	42.60
+ MA (Ours)	62.33 ± 5.36	36.67	36.04
+ LOMMA (Ours)	86.33 ± 5.17	60.67	35.59

information. The results for KEDMI [7] are shown in Table B.3. The key observations are:

- Baseline attack accuracies are poorer under limited \mathcal{D}_{pub} , i.e. $\mathcal{D}_{pub} = 10\%$.
- Our proposed method can outperform existing SOTA under varying degrees of prior information although the improvement obtained by KD is marginal under $\mathcal{D}_{pub} = 10\%$.

C. Additional analysis and details on experimental setups

C.1. Details on combining L_{id}^{logit} and L_{id}^{aug}

We provide details of combining L_{id}^{logit} and L_{id}^{aug} . We substitute L_{id}^{logit} (Eqn. 3 of main paper) into L_{id}^{aug} (Eqn. 4 of main paper) for an inversion targeting class k of the target model M_t , using augmented model $M_{aug}^{(i)}$. In particular,

Table B.1. We report top-1 attack accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance for using different numbers N_{aug} of network M_{aug} . Following exact experiment setups in [7], here method = KEDMI, $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, $M_t = \text{IR152}$, evaluation model = face.evoLve. We select M_{aug} from the set of 4 networks including EfficientNet-B0, EfficientNet-B1, EfficientNet-B2, EfficientNet-B3. The number of network M_{aug} increases from 0 (Baseline KEDMI) to 4. It shows that using more M_{aug} improves the attack accuracy and KNN distance.

Method	N_{aug}	M_{aug}	Attack Acc \uparrow	Imp. \uparrow	KNN dist \downarrow
CelebA/CelebA/IR152					
KEDMI	-	-	80.53 ± 3.86	-	1247.28
+ MA	1	EfficientNet-B0	81.20 ± 3.75	0.67	1234.16
+ MA	2	EfficientNet-B0, EfficientNet-B1	84.47 ± 2.99	3.94	1223.56
+ MA	3	EfficientNet-B0, EfficientNet-B1, EfficientNet-B2	84.73 ± 3.76	4.20	1220.23
+ MA	4	EfficientNet-B0, EfficientNet-B1, EfficientNet-B2, EfficientNet-B3	85.87 ± 2.63	5.34	1217.15

Table B.2. We report top-1 attack accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance for different structures of network M_{aug} . Following exact experiment setups in [7], here method = KEDMI, $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, $M_t = \text{IR152}$, evaluation model = face.evoLve. We select different network architectures for our experiment. Specifically, we use Ours-1 = {EfficientNet-B0 [43], EfficientNet-B1 [43], EfficientNet-B2 [43]}, Ours-2 = {DenseNet121 [21], DenseNet161 [21], DenseNet169 [21]}, Ours-3 = {EfficientNet-B0, DenseNet121 [21], MobileNetV3-large [20]}. It shows that using different network architectures M_{aug} consistently improves the attack accuracy and KNN distance.

Method	M_{aug}	Attack Acc \uparrow	Imp. \uparrow	KNN dist \downarrow
CelebA/CelebA/IR152				
KEDMI	-	80.53 ± 3.86	-	1247.28
+ MA (Ours-1)	EfficientNet-B0, EfficientNet-B1, EfficientNet-B2	84.73 ± 3.76	4.20	1220.23
+ MA (Ours-2)	DenseNet121, DenseNet161, DenseNet169	89.07 ± 3.32	8.54	1211.73
+ MA (Ours-3)	EfficientNet-B0, DenseNet121, MobileNetV3-large	86.53 ± 1.98	6.00	1204.94

starting from Eqn. 4 of the main paper:

$$\begin{aligned}
L_{id}^{aug}(\mathbf{x}; y) &= \gamma_t \cdot L_{id}(\mathbf{x}; y, M_t) \\
&\quad + \gamma_{aug} \cdot \sum_{i=1}^{N_{aug}} L_{id}(\mathbf{x}; y, M_{aug}^{(i)}) \\
&= \gamma_t \cdot L_{id}^{logit}(\mathbf{x}; y, M_t) \\
&\quad + \gamma_{aug} \cdot \sum_{i=1}^{N_{aug}} L_{id}^{logit}(\mathbf{x}; y, M_{aug}^{(i)}) \\
&= \gamma_t \cdot (-\mathbf{p}_t^T \mathbf{w}_{t,k} + \lambda \|\mathbf{p}_t - \mathbf{p}_{reg}\|_2^2) \\
&\quad + \gamma_{aug} \cdot \sum_{i=1}^{N_{aug}} (-(\mathbf{p}_{aug}^{(i)})^T (\mathbf{w}_{aug,k}^{(i)})) \\
&\quad + \lambda \|\mathbf{p}_{aug}^{(i)} - \mathbf{p}_{reg}\|_2^2 \\
&\approx \gamma_t \cdot (-\mathbf{p}_t^T \mathbf{w}_{t,k}) \\
&\quad + \gamma_{aug} \cdot \sum_{i=1}^{N_{aug}} (-(\mathbf{p}_{aug}^{(i)})^T (\mathbf{w}_{aug,k}^{(i)})) \\
&\quad + \lambda' \|\mathbf{p}_t - \mathbf{p}_{reg}\|_2^2 \tag{5}
\end{aligned}$$

Here, $\mathbf{p}_t, \mathbf{w}_{t,k}$ are penultimate layer activation and last layer weight for the target model M_t ; $\mathbf{p}_{aug}^{(i)}, \mathbf{w}_{aug,k}^{(i)}$ are penultimate layer activation and last layer weight for the aug-

mented model $M_{aug}^{(i)}$. Note that one regularization is sufficient as shown in the last step. Eqn. 5 above is used in Eqn. 1 of the main paper in the inversion step using the proposed method.

C.2. Details on improving KEDMI baseline

We apply a simple technique that is introduced by GMI [52] to get better results for KEDMI [7]. Specifically, after model inversion, and sampling \mathbf{z} from the learned distribution, we clip all elements of \mathbf{z} into $[-1, 1]$, which is shown to be beneficial in [52]. In Table C.1, we observe that clipping \mathbf{z} help to boost the attack accuracy of KEDMI and the reconstructed images are more similar to the private dataset as KNN distances are reduced. Therefore, for all the experiments with KEDMI in the main paper and Supp, we clip \mathbf{z} to get better results and we compare with this better version of KEDMI.

C.3. Additional details on computing \mathbf{p}_{reg}

In Sec. 3.1, we propose an improved formulation for identity loss L_{id}^{logit} which includes a regularization term $\|\mathbf{p} - \mathbf{p}_{reg}\|_2^2$ to prevent unbound growth of norm during optimization. Here we provide additional details on computing \mathbf{p}_{reg} .

Given that the attacker has no access to private training data, we estimate \mathbf{p}_{reg} by a simple method using *public*

Table B.3. Sensitivity of the proposed method to prior information, \mathcal{D}_{pub} : We use $\mathcal{D}_{priv}/\mathcal{D}_{pub} = \text{CelebA}$, $M_t = \text{face.evoLve}$, evaluation = face.evoLve and KEDMI [7]. We report top 1 MI attack accuracy and KNN distance using 10%, 50% and 100% of D_{pub} . As GAN is trained on D_{pub} , it affects the baseline KEDMI and our proposed method. The results show that + LOM and + MA consistently improve upon the baseline.

	$\mathcal{D}_{pub} = 10\%$		$\mathcal{D}_{pub} = 50\%$		$\mathcal{D}_{pub} = 100\%$	
	Attack Acc \uparrow	KNN Dist \downarrow	Attack Acc \uparrow	KNN Dist \downarrow	Attack Acc \uparrow	KNN Dist \downarrow
KEDMI	58.33 ± 5.25	1450.06	79.07 ± 3.76	1265.37	81.40 ± 3.25	1248.32
+ LOM (Ours)	67.27 ± 1.83	1395.38	89.27 ± 0.96	1202.45	92.53 ± 1.51	1183.76
+ MA (Ours)	61.80 ± 3.03	1421.83	82.20 ± 2.77	1244.21	85.07 ± 2.71	1222.02
+ LOMMA (Ours)	74.40 ± 2.21	1328.79	89.67 ± 0.76	1170.37	93.20 ± 0.85	1154.32

Table C.1. We apply a simple technique that is introduced by GMI [52] to get better baseline results for KEDMI [7]. We report the results for KEDMI with and without \mathbf{z} clipping for IR152, face.evoLve, and VGG16 target model. Following exact experiment setups in [7], here $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, evaluation model = face.evoLve. We report top-1 attack accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance. The improvement using \mathbf{z} clipping is clear.

Method	Attack Acc \uparrow	Imp. \uparrow	KNN dist \downarrow
CelebA/CelebA/IR152			
KEDMI w/o \mathbf{z} clipping	78.53 ± 3.45	-	1270.87
KEDMI with \mathbf{z} clipping	80.53 ± 3.86	2.00	1247.28
CelebA/CelebA/face.evoLve			
KEDMI w/o \mathbf{z} clipping	78.00 ± 4.09	-	1290.62
KEDMI with \mathbf{z} clipping	81.40 ± 3.25	3.40	1248.32
CelebA/CelebA/VGG16			
KEDMI w/o \mathbf{z} clipping	67.93 ± 4.24	-	1345.03
KEDMI with \mathbf{z} clipping	74.00 ± 3.10	6.07	1289.88

data. We firstly construct the set of penultimate layer features of public data using the target model and estimate the mean μ_{pen} and variance σ_{pen}^2 :

$$\mu_{pen} = \frac{1}{N} \sum_{i=1}^N M^{pen}(\mathbf{x}_i) \quad (6)$$

$$\sigma_{pen}^2 = \frac{1}{N} \sum_{i=1}^N (M^{pen}(\mathbf{x}_i) - \mu_{pen})^2 \quad (7)$$

where \mathbf{x}_i is a sample from public dataset \mathcal{D}_{pub} , and $M^{pen}()$ operator returns the penultimate layer representations of the target model M_t for a given input \mathbf{x} . We analyze two ways to estimate \mathbf{p}_{reg} as follow:

- Fixed \mathbf{p}_{reg} where $\mathbf{p}_{reg} = \mu_{pen}$.
- \mathbf{p}_{reg} is sampled using the prior distribution $\mathcal{N}(\mu_{pen}, \sigma_{pen})$.

Empirically, we use $N = 5,000$ images from the public dataset \mathcal{D}_{pub} to estimate μ_{pen} and σ_{pen} . The results show that using \mathbf{p}_{reg} which is sampled from $\mathcal{N}(\mu_{pen}, \sigma_{pen})$ gives better performance than using fixed $\mathbf{p}_{reg} = \mu_{pen}$ (see Table C.2). Therefore, all the results reported in the main paper

Table C.2. We report the results for KEDMI using a fixed \mathbf{p}_{reg} or sampling from a distribution approximated for \mathbf{p}_{reg} . We use three different target models: IR152, face.evoLve, and VGG16. Following exact experiment setups in [7], here $\mathcal{D}_{priv} = \text{CelebA}$, $\mathcal{D}_{pub} = \text{CelebA}$, evaluation model = face.evoLve. We report top-1 attack accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance.

Method	Attack Acc \uparrow	KNN dist \downarrow
CelebA/CelebA/IR152		
+ LOM (Fixed p_{reg})	92.27 ± 1.37	1155.92
+ LOM (Ours)	92.47 ± 1.41	1168.55
CelebA/CelebA/face.evoLve		
+ LOM (Fixed p_{reg})	90.40 ± 1.68	1257.95
+ LOM (Ours)	92.53 ± 1.51	1183.76
CelebA/CelebA/VGG16		
+ LOM (Fixed p_{reg})	85.60 ± 1.79	1259.60
+ LOM (Ours)	89.07 ± 1.46	1218.46

use the $\mathbf{p}_{reg} \sim \mathcal{N}(\mu_{pen}, \sigma_{pen})$. We remark again that \mathbf{p}_{reg} is estimated from *public* dataset.

Table C.3. We report the results for KEDMI with different λ values using IR152 as target model. Following exact experimental setups in [7], here \mathcal{D}_{priv} = CelebA, \mathcal{D}_{pub} = CelebA, evaluation model = face.evoLve. We report top-1 attack accuracies, the improvement compared to the SOTA MI (Imp.), and KNN distance.

Method	λ	Attack Acc \uparrow	Imp. \uparrow	KNN dist \downarrow
CelebA/CelebA/IR152				
KEDMI	-	80.53 ± 3.86	-	1247.28
+ LOM	0	90.33 ± 1.64	9.80	1198.39
+ LOM	0.5	89.53 ± 1.21	9.00	1175.35
+ LOM	1.0	92.47 ± 1.41	11.94	1168.55
+ LOM	2.0	91.87 ± 1.09	11.34	1125.54
+ LOM	10.0	85.80 ± 1.24	5.27	1110.80

C.4. Details on regularization parameter λ

In Sec 3.1 of the main paper, the regularization term $\|\mathbf{p} - \mathbf{p}_{reg}\|_2^2$ includes a parameter λ which controls the effect of this term. In this section, we evaluate the effect of this parameter by examining different values of λ on model inversion performance. Results in Table C.3 show that attack accuracy is improved over SOTA KEDMI with our proposed logit loss even without the regularization term ($\lambda = 0$). However, we get better results if the regularization is added *e.g.* $\lambda = 1.0$. Due to its better performance, we use $\lambda = 1.0$ in all experiments with the proposed method.

C.5. Computational overhead

In order to investigate the computational overhead introduced by our proposed method, in this section, we report the running time for reconstructing images of 300 identities on CelebA/CelebA/IR152 setup for KEDMI and GMI, and 100 identities on CelebA/CelebA/ResNet-34 for VMI. All the experiments of KEDMI and GMI are performed on an NVIDIA GeForce RTX 3090 GPU, and the experiments of VMI are performed on an NVIDIA RTX A5000 GPU. The results in Table C.4 show that + LOM does not affect the training time compared to the baseline. However, + MA adds some computational overhead as it uses additional networks M_{aug} during the inversion.

C.6. Hyperparameters

In the experiments of GMI and KEDMI, we do the inversion using SGD optimizer with the learning rate $lr = 0.02$ in 2400 iterations which are used from the released code of KEDMI². We set $\gamma_t = \gamma_{aug} = 100/(N_{aug} + 1)$ and $\lambda = 100$, where N_{aug} is the number of models used for augmented model M_{aug} . We estimate \mathbf{p}_{reg} for each classifier by using $N = 5,000$ images from the public dataset \mathcal{D}_{pub} . In the experiments of VMI, we use 20 epochs (equal to 3120 iterations) to learn the distribution of each identity.

²<https://github.com/SCccc21/Knowledge-Enriched-DMI>

Table C.4. Computational complexity of different algorithms in terms of average running time (GPU hours) using single GPU. We use KEDMI, GMI, and VMI approaches as the baseline. We have also included the running time Ratio when compared to the corresponding baseline.

Method	RunTime (hrs)	Ratio
KEDMI	0.35	1.00
+ LOM (Ours)	0.35	1.00
+ MA (Ours)	0.60	1.71
+ LOMMA (Ours)	0.60	1.72
GMI	1.61	1.00
+ LOM (Ours)	1.61	1.00
+ MA (Ours)	2.83	1.76
+ LOMMA (Ours)	2.85	1.77
VMI	364.67	1.00
+ LOM (Ours)	368.24	1.01
+ MA (Ours)	368.69	1.01
+ LOMMA (Ours)	379.41	1.04

C.7. Dataset

Experiments of KEDMI and GMI. We follow exact experimental setups in [7]. For the CelebA task, we use the dataset divided by [7] for all of the experiments. In particular, the private dataset has 30,027 images of 1000 identities and the public dataset has 30000 images that are non-overlapping identities with the private dataset. In the experiments in Table 5 (main paper), we use FFHQ [24] as the public dataset to train GAN and distill knowledge to augmented models. For MNIST and CIFAR-10 tasks, the private dataset contains images with labels from 0 to 4 and the public dataset includes the rest of the dataset with labels from 5 to 9.

Experiments of VMI. We follow exact experimental setup in [47]. We use the CelebA dataset and MNIST dataset for VMI experiments. For CelebA, we follow [47] to divide the dataset into two parts. The first part contains images of 1000 most frequent identities which uses as private dataset. The rest of dataset is used as public dataset. For the experiments on MNIST dataset, we use EMNIST [11] as public dataset to train GAN and M_{aug} .

D. Additional Visualizations

D.1. Additional Results for GMI

Similar to results reported for KEDMI (Figure 4, main paper), in this section, we show results for GMI [52] under IR152 target classifier to show the efficacy of our proposed methods. The result is shown in Figure D.1.

	GMI	Attack Acc. (\uparrow)	KNN (\downarrow)
Private Training Data			
Existing SOTA		30.60%	1609.29
+ LOM (Ours)		78.53%	1289.62
+ MA (Ours)		61.20%	1389.99
+ LOMMA (Ours)		82.40%	1254.32

Figure D.1. Qualitative / Quantitative (Top1 Attack Acc., KNN Dist) results to demonstrate the efficacy of our proposed method. We use GMI [52], \mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = CelebA [33] and M = IR152 [17]. As one can observe, our proposed method achieves better reconstruction of private data both visually and quantitatively (validated by KNN results) resulting in a significant boost in attack performance.

D.2. Penultimate layer visualization for GMI, KEDMI and VMI

In this section, we show additional penultimate layer visualizations to support our formulation of L_{id}^{logit} as an improved MI Identity Loss. We show visualizations for GMI [52] and VMI [47] in Figures D.2 and D.5 respectively. Further, we show penultimate layer visualization for an additional target classifier, face.evoLve using KEDMI [7] in Figure D.8 to validate our findings.

D.3. Our reconstruction results

Given that the goal of MI is to reconstruct private training data, in this section, we show reconstructed samples for 5 additional setups using our proposed method. We show reconstruction results using GMI [52] and VMI [47] in Figures D.7 and D.9 respectively. Further, we show additional reconstruction results for GMI and KEDMI using a different target classifier (face.evoLve) in Figures D.3 and D.4 to validate the efficacy of our proposed method. Finally, we show reconstruction results for Cross-dataset MI in Figure D.6. We remark that cross-dataset MI is a challenging attack setup due to large distribution shift between private and public data. Following [47], we use FFHQ [24] as the public dataset. To conclude, we remark that the samples reconstructed using our proposed method closely resembles the private training data in many instances, and this is quantitatively validated using MI attack accuracy.

E. Additional Related work

Given a trained model, Model Inversion (MI) aims to extract information about training data. Fredrikson et al.

[14] propose one of the first methods for MI. The authors found that attackers can extract genomic and demographic information about patients using the ML model. In [13], Fredrikson et al. extended the problem to the facial recognition setup where the authors can recover the face images. In [49], Yang et al. proposed adversarial model inversion which uses the target classifier as an encoder to produce a prediction vector. A second network takes the prediction vector as the input to reconstruct the data.

Instead of performing MI attacks directly on high-dimensional space (e.g. image space), recent works have proposed to reduce the search space to latent space by training a deep generator [7, 47, 48, 52]. In particular, a generator is trained on an auxiliary dataset that has a similar structure to the target image space. In [52], the authors proposed GMI which uses a pretrained GAN to learn the image structure of the auxiliary dataset and finds the inversion images through the latent vector of the generator. Chen et al. [7] extend GMI by training discriminator to distinguish the real and fake samples and to be able to predict the label as the target model. Furthermore, the authors proposed modeling the latent distribution to reduce the inversion time and improve the quality of reconstructed samples. VMI [47] provides a probabilistic interpretation for MI and proposes a variational objective to approximate the latent space of target data.

Zhao et al. [53] propose to embed the information of model explanations for model inversion. A model explanation is trained to analyze and constrain the inversion model to learn useful activations. Another MI attack type is called label-only MI attacks which attackers only access the predicted label without a confidence probability [10, 22]. Recently, Kahla et al. [22] propose to estimate the direction to reach the target class’s centroid for an MI attack. In this work, we instead focus on a different problem and propose two improvements to the identity loss which is common among all SOTA MI approaches. In future work, we hope to explore model inversion for different tasks including multimodal learning and data-centric applications [4, 5, 16, 25, 27, 32, 37, 46, 51].

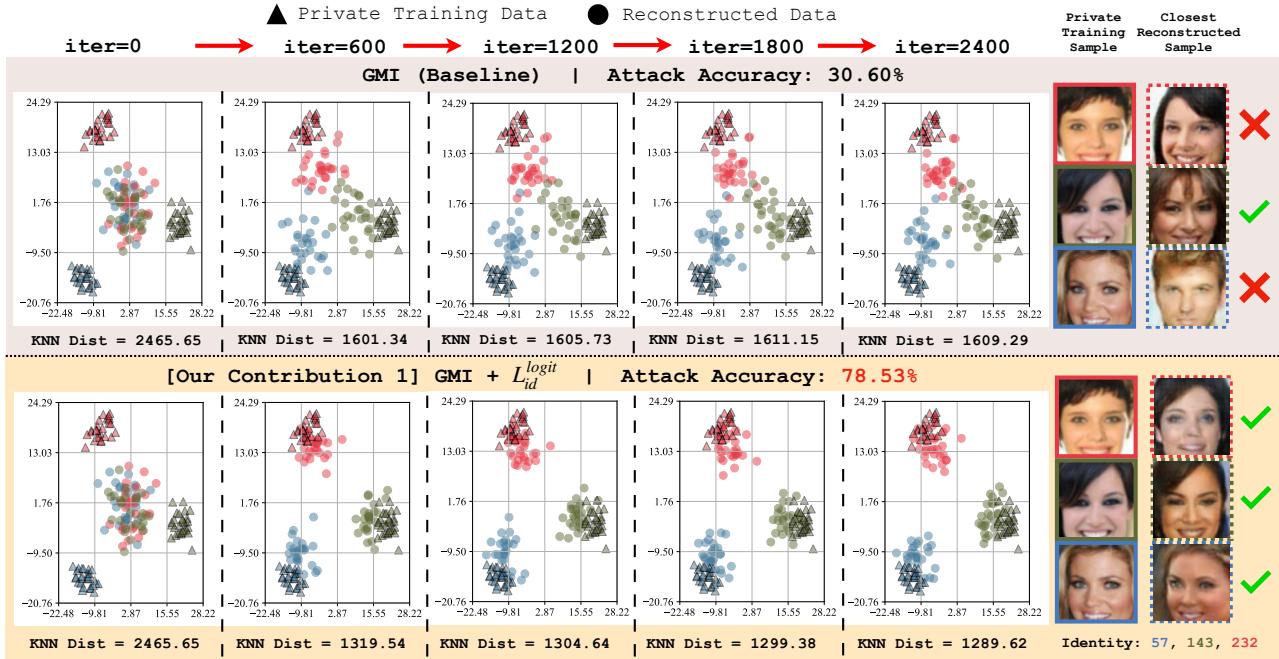


Figure D.2. Visualization of the penultimate layer representations (\mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = CelebA [33], M_t = IR152 [17], Evaluation Model = face.evoLve [8], Inversion iterations = 2400) for private training data and reconstructed data using GMI [52]. Following exact evaluation protocol in [7], we use face.evoLve [8] to extract representations. We show results for 3 randomly chosen identity. We include KNN distance (for different iterations) and final attack accuracy following the protocol in [7]. For each identity, we also include a randomly selected private training data and the closest reconstructed sample at iteration=2400. **① Identity loss in SOTA MI methods [7, 47, 52] (Eqn. 2, main paper) is sub-optimal for MI (Top).** Using penultimate representations during inversion, we observe 2 instances (e.g. target identity 57 and 232) where GMI [52] (using Eqn. 2, main paper for identity loss) is unable to reconstruct data close to private training data. Hence, private and reconstructed facial images are qualitatively different. **② Our proposed identity loss, L_{id}^{logit} (Eqn. 3, main paper), can effectively guide reconstruction of data close to private training data (Bottom).** This can be clearly observed using both penultimate layer representations and KNN distances for all 3 target classes 57, 143 and 232. Best viewed in color.



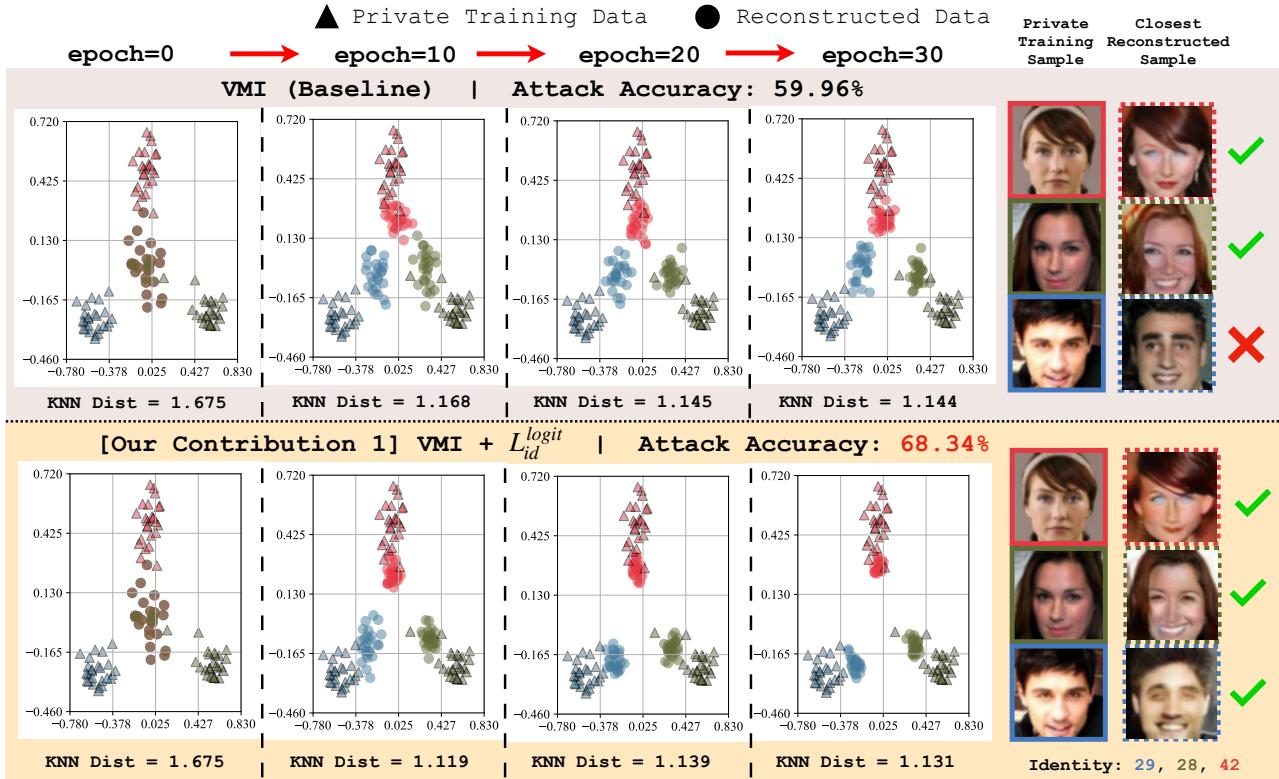


Figure D.5. Visualization of the penultimate layer representations (\mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = CelebA [33], M_t = ResNet34 [47], Evaluation Model = IR-SE50 [47], Inversion epochs = 30) for private training data and reconstructed data using VMI [47]. Following exact evaluation protocol in [47], we use IR-SE50 to extract representations. We show results for 3 randomly chosen identity. We include KNN distance and final attack accuracy. Given that we strictly follow [47], we remark that due to the use of IR-SE50 evaluation classifier to extract penultimate layer representations, the features have different scales resulting in lower KNN distances (compared to KEDMI and GMI results). For each identity, we include a randomly selected private training data and the closest reconstructed sample (epoch = 30). **① Identity loss in SOTA MI methods [7,47,52] (Eqn. 2, main paper) is sub-optimal for MI (Top).** Using penultimate representations during inversion, we observe an instance (e.g. target identity 29) where VMI [47] (using Eqn. 2, main paper for identity loss) is unable to reconstruct data close to private training data. Hence, private and reconstructed facial images are qualitatively different. **② Our proposed identity loss, L_{id}^{logit} (Eqn. 3, main paper), can effectively guide reconstruction of data close to private training data (Bottom).** This can be observed using penultimate layer representations and KNN distances for all 3 target classes 29, 28 and 42. Best viewed in color.



Figure D.6. *Cross-dataset MI results*. We show private data (top), *our* reconstruction results (bottom) and Attack accuracy (\mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = FFHQ [24], M_t = IR152 [17], Evaluation Model = face.evoLve [8], Inversion iterations = 2400) using KEDMI [7]. Cross-dataset MI is a challenging setup due to the large distribution shift between private and public data. We remark that these results are obtained by combining L_{id}^{logit} and L_{id}^{aug} (referred to as + LOMMA throughout the paper).



Figure D.7. We show private data (top), *our* reconstruction results (bottom) and Attack accuracy (\mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = CelebA [33], M_t = IR152 [17], Evaluation Model = face.evoLve [8], Inversion iterations = 2400) using GMI [52]. We remark that these results are obtained by combining L_{id}^{logit} and L_{id}^{aug} (referred to as + LOMMA throughout the paper).

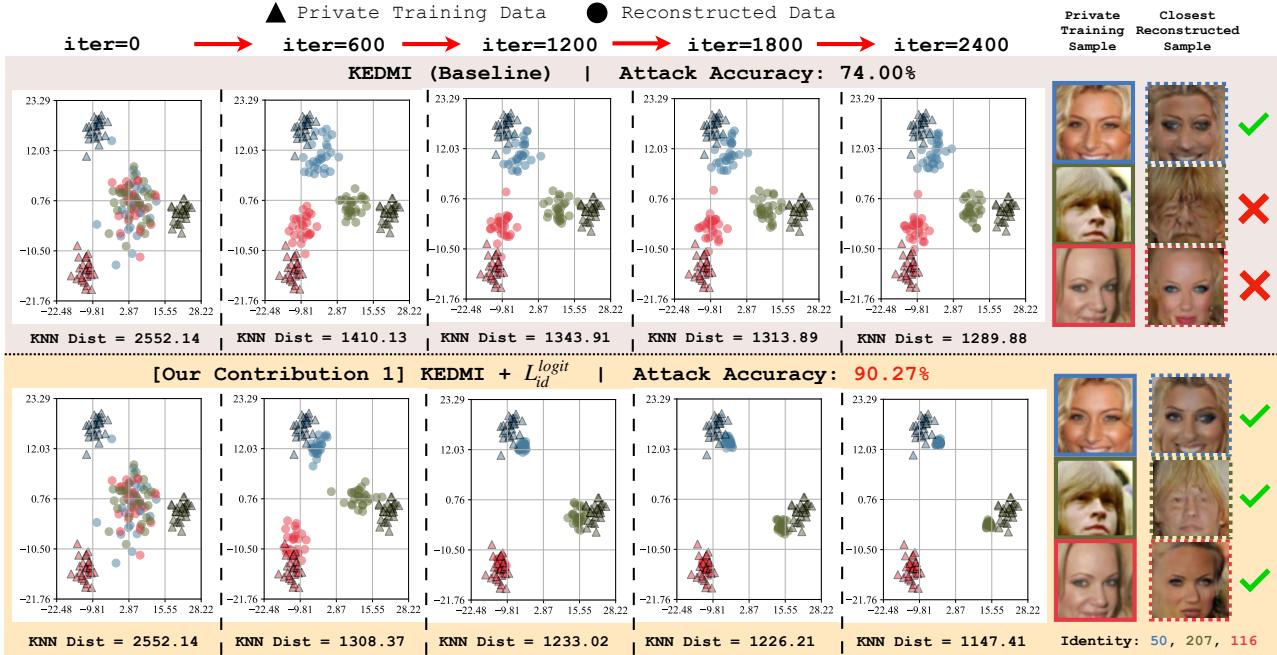


Figure D.8. Visualization of the penultimate layer representations (\mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = CelebA [33], M_t = VGG16 [40], Evaluation Model = face.evoLve [8], Inversion iterations = 2400) for private training data and reconstructed data using KEDMI [7]. Following exact evaluation protocol in [7], we use face.evoLve [8] to extract representations. We show results for 3 randomly chosen identity. We include KNN distance (for different iterations) and final attack accuracy following the protocol in [7]. For each identity, we also include a randomly selected private training data and the closest reconstructed sample at iteration=2400. ① **Identity loss in SOTA MI methods [7, 47, 52] (Eqn. 2, main paper) is sub-optimal for MI (Top).** Using penultimate representations during inversion, we observe 2 instances (e.g. target identity 207 and 116) where KEDMI [7] (using Eqn. 2, main paper for identity loss) is unable to reconstruct data close to private training data. Hence, private and reconstructed facial images are qualitatively different. ② **Our proposed identity loss, L_{id}^{logit} (Eqn. 3, main paper), can effectively guide reconstruction of data close to private training data (Bottom).** This can be clearly observed using both penultimate layer representations and KNN distances for all 3 target classes 50, 207 and 116. Best viewed in color.



Figure D.9. We show private data (top), *our* reconstruction results (bottom) and Attack accuracy (\mathcal{D}_{priv} = CelebA [33], \mathcal{D}_{pub} = CelebA [33], M_t = ResNet34 [47], Evaluation Model = IR-SE50 [12], Inversion epochs = 30) using VMI [47]. We remark that these results are obtained by combining L_{id}^{logit} and L_{id}^{aug} (referred to as + LOMMA throughout the paper).