

# Improving Robustness to Model Inversion Attacks via Sparse Coding Architectures

Sayanton V. Dibbo<sup>1,2</sup>, Adam Breuer<sup>1</sup>, Juston Moore<sup>2</sup>, and Michael Teti<sup>2</sup>

<sup>1</sup> Dartmouth College, Hanover, NH 03755, USA

<sup>2</sup> Los Alamos National Laboratory, Los Alamos, NM 87545, USA  
{f0048vh,adam.breuer}@dartmouth.edu  
{jmoore01,mteti}@lanl.gov

**Abstract.** Recent model inversion attack algorithms permit adversaries to reconstruct a neural network’s private and potentially sensitive training data by repeatedly querying the network. In this work, we develop a novel network architecture that leverages sparse-coding layers to obtain superior robustness to this class of attacks. Three decades of computer science research has studied sparse coding in the context of image denoising, object recognition, and adversarial misclassification settings, but to the best of our knowledge, its connection to state-of-the-art privacy vulnerabilities remains unstudied. In this work, we hypothesize that sparse coding architectures suggest an advantageous means to defend against model inversion attacks because they allow us to control the amount of irrelevant private information encoded by a network in a manner that is known to have little effect on classification accuracy. Specifically, compared to networks trained with a variety of state-of-the-art defenses, our sparse-coding architectures maintain comparable or higher classification accuracy while degrading state-of-the-art training data reconstructions by factors of 1.1 to 18.3 across a variety of reconstruction quality metrics (PSNR, SSIM, FID). This performance advantage holds across 5 datasets ranging from CelebA faces to medical images and CIFAR-10, and across various state-of-the-art SGD-based and GAN-based inversion attacks, including *Plug- $\mathcal{E}$ -Play* attacks. We provide a cluster-ready PyTorch codebase to promote research and standardize defense evaluations.

**Keywords:** Model Inversion Attack, Defense, Privacy Attacks, Sparse Coding.

## 1 Introduction

The popularization of machine learning has been accompanied by the widespread use of neural networks that were trained on private, sensitive, and proprietary datasets. This has given rise to a new generation of privacy attacks that seek to infer private information about the training dataset simply by inspecting the representation of the training data that remains encoded in the model’s parameters [13, 18, 19, 23, 24, 34, 40, 47, 52, 68, 72, 84, 93, 94, 98].

Of particular concern is a devastating stream of privacy attacks known as model inversion. Model inversion attacks leverage the network’s parameters or

classifications in order to reconstruct entire images or data that were used to train the network. Early work on model inversion focused on a white-box setting where the attacker has unfettered access to the model or auxiliary information about the training data [23, 33, 85, 87, 95]. However, recent work has shown that standard network architectures are vulnerable to model inversion attacks even when attackers have no knowledge of the model’s architecture or parameters, and only have access to the model’s classifications or its intermediate outputs, such as leaked outputs from a single hidden network layer [5, 25, 52, 53, 66, 92].

*Are different network architectures robust to model inversion attacks?*

Such attacks are feasible because each hidden layer of a standard network architecture captures a detailed representation of the training data. It is well-known that standard dense layers exhibit a tendency to memorize their inputs [11, 26, 63], so even a minimal leak of a network’s class distribution output or a leak of its intermediate outputs from a single layer is often sufficient to train an inverse mapping for data reconstruction. More concretely, state-of-the-art inversion attacks work by submitting externally obtained images to the model, observing leaked outputs, then using this data to train a new ‘inverted’ neural network that reconstructs (predicts) an input image given a leaked output. This can be accomplished either directly via SGD, or by optimizing a GAN, and we consider both approaches here. Such attacks on standard network architectures can reconstruct private training images that are clearly recognizable by humans familiar with the training data [4, 21, 25, 29, 33, 38, 72, 87, 92].

Recent work has pursued a diverse array of defense strategies to mitigate these attacks. For example, [25] augments the training dataset with GAN-generated fake samples designed to inject spurious features into the trained network that mislead the gradients computed during inversion attacks. In contrast, multiple recent defenses add regularization terms during training that attempt to penalize training data memorization [62, 82]. Other recent defenses noise the network weights to obfuscate memorized data [2, 54, 75], or noise and clip training gradients via DP-SGD [28]. All such approaches are costly: data augmentation-based defenses entail the computational burden of building a GAN and applying sophisticated parameter tuning techniques during training; regularization-based defenses explicitly trade away classification accuracy for less memorization, and noise-based defenses are also known to impose significant accuracy costs. Until very recently, there were no known provable guarantees for model inversion defenses, and the current best-known guarantees require a DP-SGD-based training algorithm that imposes a significant computational burden and accuracy loss to obtain privacy guarantees that are impractically weak for these attacks [28].

Very little is known about how a network’s architecture contributes to its robustness (or vulnerability). This is surprising since throughout three decades of research in other domains such as image denoising [3, 6, 10, 14, 22, 45, 58, 64], object recognition [27, 42, 59, 69], and adversarial misclassification [20, 43, 60, 73, 74], researchers seeking to control their model’s representations of the data have heavily studied sparse coding-based architectures that prune unnecessary details and preserve only the information that is essential to the model objective.

Specifically, sparse coding seeks to approximately represent an image (or layer) with only a small set of basis vectors selected from an overcomplete dictionary [10, 22, 58]. While it is well-known that computing a sparse representation using a standard objective function is NP-hard in general [17, 36, 56], we now benefit from fast approximation algorithms that efficiently compute high-quality sparse representations [8, 15, 36, 42, 45, 46, 55, 64]. Sparse coding architectures leverage this technique by inserting a sparse network layer after a dense layer, such that the sparse layer reduces the dense layer’s outputs to a sparse representation.

To our knowledge, sparse coding architectures have not been studied in the context of model inversion or privacy attacks. However, they suggest an advantageous means to prevent such attacks because they control the amount of irrelevant private information encoded in a model’s intermediate representations in a manner that is known to have little effect on its accuracy, that can be computed efficiently during training, and that adds little to the trained model’s overall parameter complexity. Put simply, sparsifying a network’s representations is a natural means to preclude memorization of detailed information about its inputs that is unnecessary to obtain high accuracy, so even an idealized ‘perfect attacker’ could only hope to recover a sparsified, un-detailed training image.

**Main contribution.** We begin by showing that an off-the-shelf sparse coding preprocessing step offers performance advantages compared to state-of-the-art data augmentation, regularization, and noise based defenses in terms of robustness to model inversion attacks. We then refine this idea into a network architecture that achieves superior performance. Our main result is a novel sparse-coding architecture, SCA, that is robust to state-of-the-art model inversion attacks.

SCA is defined by pairs of alternating sparse coded and dense layers that jettison unnecessary private information in the input image and ensure that downstream layers do not e.g., reconstruct this information. We show that SCA maintains comparable or higher classification accuracy while degrading state-of-the-art training data reconstructions 1.1 to 11.7 times more than 8 state-of-the-art data augmentation, regularization, and noise-based defenses in terms of PSNR and FID metrics and 1.1 to 720 times more in terms of SSIM. This performance advantage holds across 5 datasets ranging from CelebA faces to medical images and CIFAR-10, and across various state-of-the-art SGD-based and GAN-based inversion attacks, including *Plug-&-Play* attacks. SCA’s defense performance is also more stable than baselines across multiple runs. We emphasize that, unlike recent state-of-the-art defenses that require sophisticated parameter tuning to perform well, SCA obtains these results absent parameter tuning (i.e., using default sparsity parameters) because sparse coding naturally precludes networks from memorizing detailed representations of the training data.

More broadly, our results show a deep connection between state-of-the-art ML privacy vulnerabilities and three decades of computer science research on sparse coding for other application domains. We provide a comprehensive cluster-ready PyTorch codebase to promote research and standardize defense evaluation.

## 2 Threat models

We consider three threat models that span the diverse range of powerful and well-informed attackers considered in recent work. We emphasize that a defense that performs well in all three settings provides strong evidence of its privacy protections under weaker, more realistic threat models with real-world attackers.

**1. Plug-&Play threat model [72].** Plug-&Play attacks are considered the most performant recent attacks. These attacks optimize the intermediate representation of StyleGAN’s input vectors so that generated images maximize the target network’s class prediction probability, which the attacker can query.

Separately, recent theoretical work on model inversion emphasizes that a strong threat model should capture ‘worst-case’ attackers with direct access to the information-rich, high-dimensional intermediate outputs of the target model that store private information about the training data, as well as ideal training data examples for training an inverted model [1, 28]. We consider two variants:

**2. End-to-end threat model.** We consider an attacker with access to all of the *last* hidden layer’s raw, high-dimensional outputs, as well as a large number of ideal training data examples drawn from the true training dataset [71, 83].

**3. Split network threat model (Federated Learning).** We also consider the split network threat model described by [75]. There has been much recent interest in Federated Learning architectures that split the network across multiple agents [7, 21, 29, 44, 51], particularly for privacy-fraught domains such as medicine where legal requirements limit data sharing [39, 77]. These architectures are known to be susceptible to model inversion attacks, [21, 29, 75], and defenses are urgently needed. This threat model also allows us to capture a different view of a ‘worst-case’ threat model: Model inversion attacks are known to be more effective when the attacker has access to outputs from earlier layers that may exhibit a more direct representation of the input images [29]. To capture this ‘worst-case’, we consider the setting where the attacker has access to raw intermediate outputs from the *first* linear network layer. As before, we also assume the attacker has access to ideal training data examples drawn from the actual training datasets. Appendix B provides additional details about the model inversion threat models.

## 3 SCA architecture

We now describe the SCA architecture, which is defined by alternating pairs of Sparse Coding Layers (SCL) and dense layers, followed by downstream linear and/or convolutional layers.

**Sparse Coding Layer (SCL).** Sparse coding converts raw inputs to sparse representations where only a few neurons whose features are useful in reconstructing the inputs are active. Our Sparse Coding Layer (SCL) performs sparse coding to obtain a sparse representation of a previous dense layer’s representation (if the SCL is not the first layer in the network) or of the inputs (if the SCL is the first layer in the network). Fig. 1 illustrates the working principle of SCL.

Formally, each SCL performs a reconstruction minimization problem to compute the sparse representation of its inputs (either a previous layer’s representation or of the inputs to the network). Suppose the input to a (2D convolutional) SCL is  $\mathcal{X} \in \mathbb{R}^{\mathcal{C} \times \mathcal{H} \times \mathcal{W}}$  with  $\mathcal{H}$  height,  $\mathcal{W}$  width, and  $\mathcal{C}$  channels/features. The goal is to find the sparse representation  $\mathcal{R}_x \in \mathbb{R}^{\mathcal{F} \times \lceil \mathcal{H}/S_h \rceil \times \lceil \mathcal{W}/S_w \rceil}$ , where  $\mathcal{R}_x$  has few active neurons and corresponds to a denoised version of the input  $\mathcal{X}$ , and  $S_w$  and  $S_h$  indicate convolutional strides across the width and height of the input, respectively.  $F$  is the number of convolutional features in the SCL layer’s dictionary,  $\Omega \in \mathbb{R}^{\mathcal{F} \times \mathcal{C} \times \mathcal{H}_f \times \mathcal{W}_f}$ , where  $\mathcal{H}_f$  and  $\mathcal{W}_f$  are the height and width of each convolutional feature, respectively. Per Figure 1, the sparse coding layer starts with its input,  $\mathcal{X}$ , and dictionary of features,  $\Omega$ , to produce  $\mathcal{R}_x$  by solving the following sparse reconstruction problem:

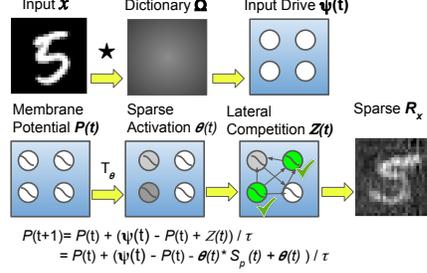
$$\min_{\mathcal{R}_x} \frac{1}{2} \|\mathcal{X} - \mathcal{R}_x \otimes \Omega\|_2^2 + \lambda \|\mathcal{R}_x\|_1 \quad (1)$$

where the first term represents how much information is preserved about  $\mathcal{X}$  by  $\mathcal{R}_x$  by measuring the difference between  $\mathcal{X}$  and its reconstruction,  $\mathcal{R}_x \otimes \Omega$ , computed with a transpose convolution,  $\otimes$ . The second term measures how sparse  $\mathcal{R}_x$  is, and  $\lambda$  is a constant which determines the trade-off between reconstruction fidelity and sparsity. Equation 1 is convex in  $\mathcal{R}_x$ , meaning we will always find the optimal  $\mathcal{R}_x$  that solves Equation 1.

Among different techniques to perform sparse coding, we leverage the commonly used Locally Competitive Algorithm (LCA) [64]. LCA implements a recurrent network of leaky integrate-and-fire neurons that incorporates the general principles of thresholding and feature-similarity-based competition between neurons to solve Equation 1. While Rozell introduced LCA in the non-convolutional setting, it can be readily adapted to the convolutional setting (see Appendix A.2 for details). Specifically, each LCA neuron has an internal membrane potential  $\mathcal{P}$  which evolves per the following differential equation:

$$\dot{\mathcal{P}}(t) = \frac{1}{\tau} [\Psi(t) - \mathcal{P}(t) - \mathcal{R}_x(t) * \mathcal{G}] \quad (2)$$

where  $\tau$  is a time constant,  $\Psi(t) = \mathcal{X} * \Omega$  is the neuron’s bottom-up drive from the input computed by taking the convolution,  $*$ , between the input,  $\mathcal{X}$ , and the dictionary,  $\Omega$ , and  $-\mathcal{P}(t)$  is the leak term [43, 74]. Lateral competition between neurons is performed via the term  $-\mathcal{R}_x(t) * \mathcal{G}$ , where  $\mathcal{G} = \Omega * \Omega - I$  is the similarity between each feature and the other  $\mathcal{F}$  features ( $-I$  prevents



**Fig. 1:** Pipeline of neuron (membrane potential) dynamics in Sparse Coding Layer (SCL) with lateral competitions.

self interactions).  $\mathcal{R}_x$  is computed by applying soft threshold activation  $T_\lambda(x) = \text{relu}(x - \lambda)$  to the neuron’s membrane potential, which produces nonnegative, sparse representations. Overall, this means that in LCA neurons will compete to determine which ones best represent the input, and thus will have non-zero activations in  $\mathcal{R}_x$ , the output of the SCL that is passed to the next layer.

**SCA architecture.** The SCA architecture is defined by the use of multiple *pairs of sparse coding and dense (batch norm) layers* after the input image, which can then be followed by other (linear, convolutional) layers. Fig. 2 illustrates this design principle. The *key intuition* is that the first sparse layer jettisons unnecessary private information in the input image.

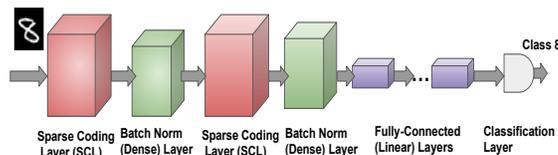


Fig. 2: Architecture of SCA.

Then, by alternating sparse-dense pairs of layers, we ensure that unnecessary information is also jettisoned from downstream layers. In this manner, downstream layers also do not convey unnecessary private information to the adversary, and they also do not e.g. learn to reconstruct private information jettisoned by the first sparse layer. In short, previous defenses work by trying to mislead attackers by pushing model features in a wrong direction, either randomly via noise or strategically via adversarial examples, or by disincentivizing memorization during training. In contrast, SCA directly removes the unnecessary private information. Training SCA is identical to training a standard network with one exception: after each backprop updates non-sparse layers, we perform a fast update on the sparse layers, except for the first sparse layer that sparse-codes the image input.<sup>3</sup>

**SCA training complexity & large scale applications.** While we focus on the neuron lateral competition approach to sparse coding as it is practically convenient and well-represented in recent work [74], we note that for large-scale machine learning applications, we now have practical parallel algorithms that learn the sparse coding dictionary near-optimally w.p. in parallel time (adaptivity) that is logarithmic in the size of the data [8, 15, 36]. Fast single-iteration heuristics are also available (see e.g. [88]). Thus, even for large-scale applications, computing sparse representations while training SCA adds little computational overhead compared to sophisticated optimization-based techniques necessary for recent defenses [25]. In practice, even our basic sparse coding research implementations (see Section 4 and Appendix A.3) are slightly faster than optimized Torch implementations of the best-performing recent defense [62].

<sup>3</sup> We can optionally also perform a backpropagation on sparse layers after updating them each iteration via sparse coding. We do this in our experiments.

## 4 Experiments

Our goal in this section is to show that SCA performs well compared to state-of-the-art defenses as well as practical defenses used in leading industry models in terms of both classification accuracy and various reconstruction quality metrics. To evaluate its performance comprehensively, we test *all combinations* of 5 diverse datasets, 3 threat models, 9 defense baselines, plus multiple runs-per-experiment and various sparsity parameters  $\lambda$ . See Appendix A.1-A.6.

**Five benchmark datasets.** We test our performance on *all 5* diverse datasets used to benchmark model inversion attacks across the recent literature: **CelebA** hi-res RGB faces, **Medical MNIST** medical images, **CIFAR-10** hi-res RGB objects, **MNIST** grayscale digits, and **Fashion MNIST** grayscale objects.

**Three attacks.** For each dataset, we conduct three sets of experiments corresponding to our three threat models. First, we test SCA and baselines’ defenses against a state-of-the-art *Plug-&Play attack* [72] that leverages *StyleGAN3* to obtain high-quality reconstructions. Second, we compare SCA networks to a variety of baselines in terms of their robustness to a state-of-the-art *end-to-end network attack that leverages leaked raw high-dimensional outputs from the networks’ last hidden layer, as well as held-out training data drawn from the true training dataset*. This allows us to assess SCA’s defenses in a realistic setting with a well-informed adversary. Our third set of experiments tests performance in a *split network setting of* [75] where the attacker has access to leaked raw outputs from the first linear network layer. Robustness in this setting is desirable because model inversion attacks are known to be more effective on earlier hidden layers [29], and also because an algorithm that is robust to such attacks would be an effective defense under novel security paradigms such as Federated Learning, which is vulnerable to inversion [75] (see Appendix A.5).

**Nine defense baselines.** We compare SCA to 9 baselines plus extra variants, including SOTA defenses and practical defenses used in leading industry models:

- **No-Defense.** The baseline target model with no added defenses.
- **Hayes et al.** [28]. We train a DP-SGD defense that noises and clips gradients during training. This is the only defense with provable guarantees.
- **Gong et al.** [25]. We train the very recent defense from [25] that uses sophisticated tuning and two types of GAN-generated images. We also try a ‘++’ version that adds extra Continual Learning accuracy optimizations.
- **Peng et al.** [62]. We train the Bilateral Dependency defense that adds a loss function for redundant input memorization during training.
- **Wang et al.** [82]. We train a Mutual Information Regularization defense that penalizes dependence between inputs and outputs during training.
- **Titcombe et al.** [75]. We train a state-of-the-art Laplace  $\mathcal{L}(\mu=0, b=0.5)$  noise defense as in [75]. We also try more noise—see Appendix D.
- **Sparse-Standard.** We train an off-the-shelf sparse coding architecture [74] with 1 sparse layer after the input image via lateral competition as in SCA.

- **GAN [common industry defense]**. We train a GAN for 25 epochs to generate fake samples, then train the target model with both original and GAN-generated samples. This defense is frequently used in industry.
- **Gaussian-Noise [common industry defense]**. We draw noises from  $\mathcal{N}(\mu=0, \sigma=0.5)$  and inject them into intermediate dense layers post-training.

**SCA without parameter tuning.** In all experiments, we consider the simplest case of SCA architecture that contains SCA’s alternating sparse-and-dense layer pairs followed by only linear layers. We note that adding downstream convolutional layers or more sophisticated downstream architectures is certainly possible, though we avoid this here in order to compare the essence of the SCA approach to the baselines. Appendix A.4 describes SCA details. In the split network setting, we are careful to use slightly shallower SCA architectures with fewer linear layers to match the split network experiments of [75].

Recent state-of-the-art defenses such as GAN-based defenses require sophisticated automatic parameter tuning techniques such as focal tuning and continual learning to obtain high performance [25]. To test whether SCA can be effective *absent* parameter tuning, we just run SCA with sparsity parameter  $\lambda$  set to **0.1**, **0.25**, or **0.5**—the default values from various sparse coding contexts.

**Performance metrics.** We evaluate attackers’ reconstructions using multiple standard metrics. Let  $X_{in}^*$  denote the reconstruction of training image  $X_{in}$ . Then:

- **Peak signal-to-noise ratio (PSNR)** [*lower=better*]. PSNR is the ratio of max squared pixel fluctuations from  $X_{in}$  to  $X_{in}^*$  over mean squared error.
- **Structural similarity (SSIM)** [86] [*lower=better*]. SSIM measures the product of luminance distortion, contrast distortion, & correlation loss:  

$$SSIM(X_{in}, X_{in}^*) = l_{dis}(X_{in}, X_{in}^*)c_{dis}(X_{in}, X_{in}^*)c_{loss}(X_{in}, X_{in}^*).$$
- **Fréchet inception distance (FID)** [31] [*higher=better*]. FID measures reconstruction quality as a distributional difference between  $X_{in}$  and  $X_{in}^*$ :  

$$FID^2(X_{in}, X_{in}^*) = \|\mu_{X_{in}} - \mu_{X_{in}^*}\|^2 + Tr(Co_{X_{in}} + Co_{X_{in}^*} - 2* \sqrt{Co_{X_{in}} \cdot Co_{X_{in}^*}})$$

**Target model.** We focus on privacy attacks on linear networks because they capture the essence of the privacy attack vulnerability [23,32], and because there is broad consensus that a principled understanding of their emerging privacy (and security) vulnerabilities<sup>4</sup> is urgently needed [30,48,67,89].

**PyTorch codebase, replicability, and evaluation standardization.** For all experiments, we consider the standard train test split of 70% and 30%. After training each defense model, we run attacks to reconstruct the entire training set and compare reconstruction performance. We run all the experiments on a standard industry production cluster with 4 nodes and DELL Tesla V100 GPUs with 40 cores. *We provide a cluster-ready PyTorch codebase on our project page at: <https://sayantondibbo.github.io/SCA>.*

<sup>4</sup> We also note that results on linear models may generalize better than results on more application-specific models, and linear models trained on private data remain ubiquitous among top industry products.

#### 4.1 Experimental results overview

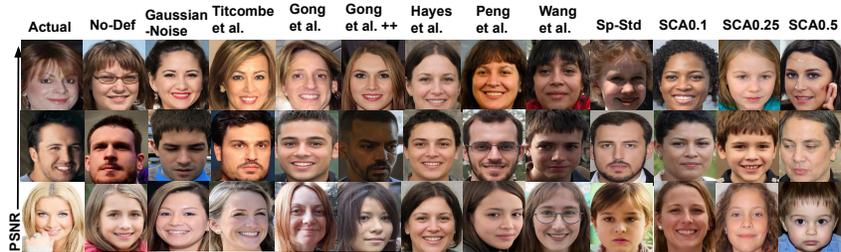
**Defense.** Across the 3 attack settings and 5 datasets, SCA maintains comparable or higher classification accuracy while degrading state-of-the-art training data reconstructions 1.1 to 11.7 times more than the 9 baselines in terms of PSNR & FID, and 1.1 to 720 times more in terms of SSIM. This performance gap exceeds the scale of improvements made by recent algorithms. SCA’s defense is also more stable than baselines across multiple runs. This is because unlike for baselines, even an idealized ‘perfect attacker’ can only hope to recover a sparsified, un-detailed training image from SCA. We show results here for the 2 most privacy-sensitive datasets of medical images and CelebA in 3 threat models, and defer the 9 other {dataset, attack} combinations to Appendix C.

**Accuracy.** Typically, obtaining greater defense means trading away accuracy. However, in 6 of the 15 experiments (MNIST + FashionMNIST)  $\times$  (Plug & Play + end-to-end + split networks), SCA *outperforms no-defense and all baselines’ accuracy*. SCA also outperforms all baselines’ accuracy on CelebA in Plug & Play, and a sparse approach is within 0.003 of the best accuracy on an 8th experiment. *No other single SOTA baseline wins on accuracy this consistently*. We emphasize that unlike baselines that do accuracy hyperparameter tuning, we obtain this result *absent* any such tuning. SCA drops accuracy on MedMNIST (which is the most imbalanced & has fewest training examples). However, tuning of SCA (kernel size 5  $\rightarrow$  7) improves SCA’s accuracy on MedMNIST in Table 3 from 94.6% to 97%—See Appendix G and Table 12.

##### Results of experiments set 1: Plug-&Play attack.

**Qualitative evaluations.** Fig. 3 shows hi-res CelebA reconstructions generated by Plug-&Play under various defenses. To avoid cherry-picking, Fig. 3 shows the 3 images with the highest (top row), median (middle), and lowest (bottom) PSNR reconstructions under No-Defense. Note that reconstructions under SCA totally differ from actual images (different race, hair, gender, child/adult), while those of other defenses are closer to actual images, indicating privacy leakage.

**Metric evaluations.** Table 1 reports reconstruction quality measures and accuracy for SCA and other baselines in the Plug-&Play attack [72] setting (*lower rows = better defense performance*). In terms of PSNR and SSIM, training data



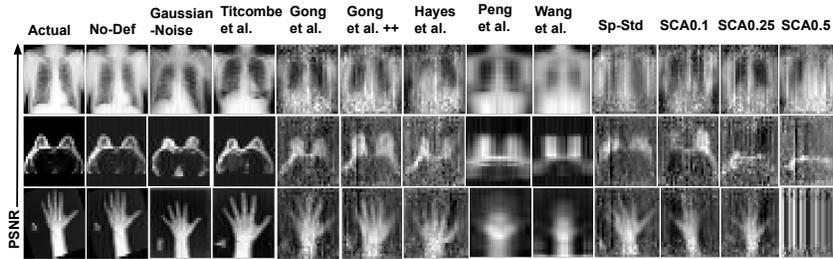
**Fig. 3:** Experiments set 1: Qualitative comparisons among actual and reconstructed images (Plug-&Play Attack [72]) under SCA & baselines on hi-res CelebA dataset.

**Table 1:** Experiments set 1: Performance comparison under Plug-&-Play attack [72] setting (*lower rows=better defense*) on hi-res CelebA faces and Medical MNIST images.

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
CelebA	NO-DEFENSE	11.35	0.718	256.4	0.779
	GAUSSIAN-NOISE	10.39	0.604	264.8	0.644
	GAN	10.15	0.613	289.7	0.635
	Titcombe et al. [75]	10.18	0.636	304.1	0.654
	Gong et al. [25]++	10.02	0.556	381.5	0.672
	Gong et al. [25]	10.11	0.595	350.5	0.614
	Peng et al. [62]	9.90	0.514	402.8	0.728
	Hayes et al. [28]	9.92	0.556	383.7	0.621
	Wang et al. [82]	9.85	0.527	402.8	0.742
	SPARSE-STANDARD	9.84	0.539	374.7	0.728
	<b>SCA0.1</b>	<b>9.79</b>	<b>0.451</b>	<b>391.9</b>	<b>0.726</b>
	<b>SCA0.25</b>	<b>9.45</b>	<b>0.442</b>	<b>411.0</b>	<b>0.739</b>
<b>SCA0.5</b>	<b>9.40</b>	<b>0.440</b>	<b>412.6</b>	<b>0.723</b>	
Medical MNIST	NO-DEFENSE	22.04	0.396	196.1	0.998
	GAUSSIAN-NOISE	21.83	0.382	209.4	0.862
	GAN	21.77	0.427	219.0	0.998
	Gong et al. [25]++	21.50	0.359	273.1	0.894
	Titcombe et al. [75]	21.68	0.360	286.3	0.899
	Gong et al. [25]	21.75	0.477	249.1	0.770
	Peng et al. [62]	21.82	0.381	268.3	0.927
	Hayes et al. [28]	21.72	0.337	259.7	0.823
	Wang et al. [82]	21.71	0.322	211.7	0.937
	SPARSE-STANDARD	20.97	0.086	239.3	0.907
	<b>SCA0.1</b>	<b>21.19</b>	<b>0.057</b>	<b>253.5</b>	<b>0.888</b>
	<b>SCA0.25</b>	<b>21.17</b>	<b>0.075</b>	<b>280.1</b>	<b>0.882</b>
<b>SCA0.5</b>	<b>20.06</b>	<b>0.055</b>	<b>288.8</b>	<b>0.881</b>	

reconstructions under the *least sparse* version SCA0.1 are degraded by factors of 1.01 to 6.7 and 1.01 to 5.6 compared to the regularization defenses of Peng et al. [62] and Wang et al. [82], respectively. Increasing SCA’s sparsity  $\lambda$  to 0.5 widens the performance gap, increasing these factors to 1.1 to 6.9 and 1.02 to 5.9, respectively, and making SCA outperform baselines’ FID. SCA0.1 also outperforms the noise-based approaches of Hayes et al. [28] and Titcombe et al. [75] on all metric by factors of 1.01 to 5.9 and 1.02 to 6.3. Increasing SCA’s sparsity  $\lambda$  to 0.5 increases these factors to 1.1 to 6.1 and 1.1 to 6.5. Finally, SCA0.1 outperforms the data augmentation defense of [25] by factors of 1.01 to 8.4, which widens to 1.1 to 8.7 for SCA0.5. All baselines also outperform common GAN and Noise-based industry defenses.

**Basic SPARSE-STANDARD outperforms SOTA baselines.** Our basic SPARSE-STANDARD baseline outperforms the best baselines’ PSNR on both CelebA and Medical MNIST, and also outperforms baselines’ SSIM on the latter. SCA0.5 then outperforms SPARSE-STANDARD on *all metrics*, obtaining SSIM and FID that are better by factors of 1.2 to 1.6 and 1.1 to 1.2, respectively, and slightly better PSNR. Thus, while SPARSE-STANDARD offers an inferior defense vs. SCA, it still offers a fast practical defense for less privacy-critical domains.



**Fig. 4:** Experiments set 2: Qualitative comparisons among actual & reconstructed images (*end-to-end* setting) under SCA & baselines on the Medical MNIST dataset.

**Table 2:** Experiments set 2: Performance comparison in *end-to-end* setting (*lower rows=better defense*) on hi-res CelebA faces and Medical MNIST images.

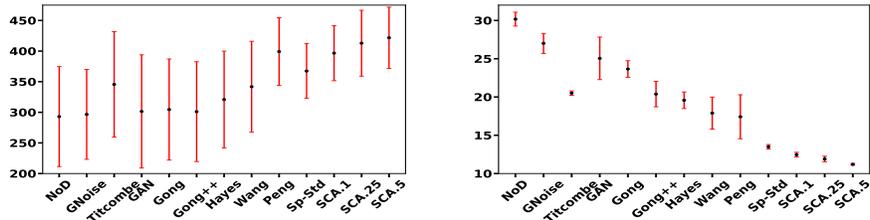
Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
CelebA	NO-DEFENSE	16.26	0.262	201.8	0.773
	GAUSSIAN-NOISE	16.08	0.262	220.4	0.638
	GAN	13.55	0.133	199.6	0.668
	Titcombe et al. [75]	15.13	0.191	197.7	0.695
	Gong et al. [25]++	13.10	0.032	204.8	0.704
	Gong et al. [25]	13.15	0.119	199.6	0.682
	Peng et al. [62]	13.78	0.141	218.8	0.716
	Hayes et al. [28]	14.10	0.004	199.0	0.664
	Wang et al. [82]	13.63	0.0011	203.2	0.744
	SPARSE-STANDARD	13.09	0.002	222.1	0.749
	<b>SCA0.1</b>	<b>12.89</b>	<b>0.004</b>	<b>228.5</b>	<b>0.748</b>
<b>SCA0.25</b>	<b>12.73</b>	<b>0.004</b>	<b>218.8</b>	<b>0.737</b>	
<b>SCA0.5</b>	<b>12.42</b>	<b>0.001</b>	<b>231.9</b>	<b>0.741</b>	
Medical MNIST	NO-DEFENSE	31.48	0.935	10.66	0.998
	GAUSSIAN-NOISE	30.46	0.920	12.23	0.862
	GAN	27.34	0.480	33.77	0.998
	Gong et al. [25]++	18.37	0.353	81.52	0.894
	Titcombe et al. [75]	21.33	0.431	30.60	0.899
	Gong et al. [25]	21.52	0.436	64.88	0.770
	Peng et al. [62]	19.05	0.420	107.9	0.908
	Hayes et al. [28]	18.48	0.007	150.9	0.824
	Wang et al. [82]	20.48	0.549	30.01	0.946
	SPARSE-STANDARD	14.79	0.119	250.6	0.907
	<b>SCA0.1</b>	<b>13.43</b>	<b>0.004</b>	<b>352.1</b>	<b>0.888</b>
<b>SCA0.25</b>	<b>12.32</b>	<b>0.004</b>	<b>375.9</b>	<b>0.882</b>	
<b>SCA0.5</b>	<b>12.04</b>	<b>0.003</b>	<b>369.9</b>	<b>0.881</b>	

### Results of experiments set 2: End-to-end networks.

*Qualitative evaluations.* Fig. 4, shows Medical MNIST reconstructions in the end-to-end threat model. SCA’s reconstructions are visually destroyed (esp. SCA0.5), whereas baselines admit noisy-but-recognizable reconstructions.

*Metric evaluations.* Table 2 reports performance in the end-to-end network setting. In this setting, SCA’s performance advantage *widens*. In terms of PSNR & FID, training data reconstructions under SCA0.1 are degraded by factors of 1.1 to 3.3 and 1.1 to 11.7 compared to Peng et al. [62] and Wang et al. [82], respectively (but slightly worse SSIM vs. Wang et al. on CelebA). On all metrics, SCA0.1 outperforms Hayes et al. [28] and Titcombe et al. [75] by factors of 1.1 to 2.4 and 1.1 to 107.7, respectively. SCA0.1 also outperforms Gong et al. [25] by factors of 1.01 to 109. Increasing SCA’s  $\lambda$  to 0.5 causes it to outperform the same baselines *on all metrics* by factors of 1.2 to 141 [62], 1.2 to 183 [82], 1.2 to 4.0 [28], 1.2 to 191 [75], and 1.1 to 145.3 [25], respectively.

*Stability of SCA’s defense vs. baselines.* SCA’s performance is also as stable or more stable than baselines’ performance over multiple runs. For example, Fig. 5a plots the means & std. devs. of SCA & baselines’ per-run FID (the standard metric for faces) over multiple runs for CelebA faces in the Plug-&-Play setting, and Fig. 5b plots this for PSNR (the standard metric for grayscale images) on Medical MNIST in the end-to-end setting. SCA obtains better performance while also exhibiting stability of defense performance on par with the best baseline (and better stability than most baselines). See Appendix E.



(a) FID↑↑ for CelebA faces under Plug-&-Play (b) PSNR↓↓ for Med.MNIST, end-to-end setting

**Fig. 5:** Stability of SCA & baselines’ defense performance (mean  $\pm$  std. dev.) of PSNR and FID across multiple runs on CelebA and Medical MNIST.

*SCA’s sparsity vs. performance.* We also try varying SCA’s and SPARSE-STANDARD’s sparsity parameters  $\lambda$  and recompute PSNR, SSIM, FID, and accuracy. Appendix A.6 shows that for each  $\lambda$  and defense metric, SCA significantly outperforms the off-the-shelf SPARSE-STANDARD architecture for a small accuracy cost. Thus, for a given  $\lambda$  with SPARSE-STANDARD, we can use a (smaller)  $\lambda$  with SCA to obtain better reconstruction *and* higher or equal (within 0.0017) accuracy. SCA is also amenable to more sophisticated tuning (and performance improvements) by tuning different  $\lambda$  per sparse layer (e.g., by having a sparser representation of inputs but less sparse reductions of downstream layers). We *avoid* such tuning here as it is unnecessary for good performance.

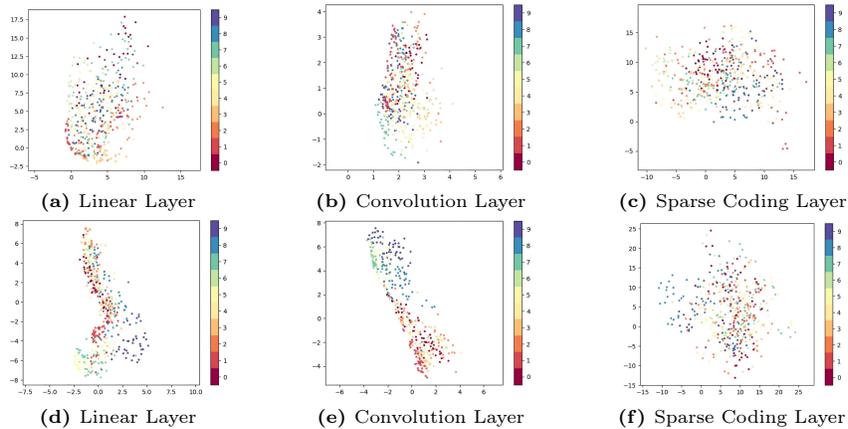
**Table 3:** Experiments set 3: Performance comparison in *split network* setting (*lower rows=better defense*) on hi-res CelebA faces and sensitive Medical MNIST images.

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
CelebA	No-DEFENSE	16.49	0.302	185.8	0.766
	GAUSSIAN-NOISE	15.44	0.227	191.1	0.753
	GAN	15.57	0.253	176.7	0.646
	Titcombe et al. [75]	14.99	0.144	194.2	0.725
	Gong et al. [25] <sup>++</sup>	15.06	0.038	190.5	0.756
	Gong et al. [25]	15.65	0.044	185.8	0.653
	Peng et al. [62]	16.23	0.211	198.6	0.717
	Hayes et al. [28]	15.06	0.005	178.8	0.672
	Wang et al. [82]	14.82	0.173	189.6	0.652
	SPARSE-STANDARD	15.39	0.009	187.0	0.746
	<b>SCA0.1</b>	<b>15.05</b>	<b>0.005</b>	<b>178.7</b>	<b>0.745</b>
	<b>SCA0.25</b>	<b>14.76</b>	<b>0.003</b>	<b>191.1</b>	<b>0.743</b>
<b>SCA0.5</b>	<b>14.71</b>	<b>0.003</b>	<b>206.1</b>	<b>0.739</b>	
Medical MNIST	No-DEFENSE	23.47	0.776	45.57	0.993
	GAUSSIAN-NOISE	21.93	0.722	44.72	0.811
	GAN	21.67	0.719	48.49	0.912
	Gong et al. [25] <sup>++</sup>	21.07	0.573	67.53	0.931
	Titcombe et al. [75]	21.35	0.704	48.82	0.961
	Gong et al. [25]	21.33	0.720	41.74	0.925
	Peng et al. [62]	18.98	0.426	124.8	0.914
	Hayes et al. [28]	21.46	0.442	137.4	0.850
	Wang et al. [82]	20.03	0.538	65.17	0.986
	SPARSE-STANDARD	15.33	0.149	142.4	0.955
	<b>SCA0.1</b>	<b>13.95</b>	<b>0.008</b>	<b>244.9</b>	<b>0.946</b>
	<b>SCA0.25</b>	<b>12.31</b>	<b>0.008</b>	<b>255.3</b>	<b>0.928</b>
<b>SCA0.5</b>	<b>12.27</b>	<b>0.001</b>	<b>285.3</b>	<b>0.909</b>	

**Results of experiments set 3: Split networks.** Table 3 reports performance in the split network setting. As expected, all baselines and SCA perform slightly worse in this threat model compared to the end-to-end model (aside from Gaussian and GAN heuristics). On all metrics, SCA’s performance advantage remains consistent: SCA0.5 outperforms all baselines by factors of 1.1 to 720.

## 5 Empirical analysis of sparse coding robustness to attack

Sparse-coding layers’ robustness to privacy attacks can be observed empirically. Consider that the attacker trains the attack to map leaked raw hidden layer outputs back to input images. Attacks are thus highly dependent on these outputs’ distributions. Recall that UMAP projections compute a 2D visualization of the global structure of distances between different training images’ features according to a particular layer [50]. Fig. 6 plots UMAP 2D projections of linear layer feature distributions of training inputs *after* either two linear layers (Figs. 6a & 6d), two convolutional layers (Figs. 6b & 6e), or two sparse coding layers (with



**Fig. 6:** UMap 2D projections of input images’ features by class after 2 linear layers, 2 conv. layers, or 2 sparse-coded layers on MNIST (top) & Fashion MNIST (bottom).

interspersed dense layers – Figs. 6c & 6f). Importantly, observe that after two linear or two convolutional layers, points are clustered by color, i.e., input images’ features are highly clustered by label. This class-clustered property leaves such layers vulnerable to model inversion attacks, as an attacker can ‘home in on’ examples from a specific class. In contrast, the goal in sparse coding is not to optimize the classification objective by separating classes, but rather to jettison unnecessary information. Here, this means that unnecessary information is jettisoned both from the input image and also the downstream dense layer. Per Figs. 6c & 6f, this tends to ‘uncluster’ remaining non-sparsified features of training examples from the same class, making it much harder for an attacker to compute informative gradients to home in on a training example.

## 6 Discussion & Conclusion

In this paper, we have provided the first study of sparse coding-based neural network architectures that are robust to model inversion attacks. Specifically, we have shown that the natural properties of sparse coded layers can control the extraneous private information about the training data that is encoded in a network without resorting to complex and computationally intensive parameter tuning techniques. Our work reveals a deep connection between state-of-the-art privacy vulnerabilities and three decades of computer science research on sparse coding for other application domains. Currently, our basic research implementation of SCA achieves compute times that in the worst-case are no better than some SOTA baselines (see Appendix F). However, given the rich theoretic body of work on fast algorithms and provable guarantees for sparse coding, we believe these aspects are opportune areas for future improvements.

## Acknowledgements

We are grateful for generous support from OpenAI, as well as the Dartmouth College Cybersecurity Cluster Research. This work was partially funded by the Center for Nonlinear Studies and the Information Science and Technology Institute’s Cyber Security Summer School at Los Alamos National Laboratory, as well as an award from the Department of Energy’s Advanced Scientific Computing Research program (#77902).

## References

1. Abadi, M., Chu, A., Goodfellow, I., McMahan, H.B., Mironov, I., Talwar, K., Zhang, L.: Deep learning with differential privacy. In: Proceedings of the 2016 ACM SIGSAC conference on computer and communications security. pp. 308–318 (2016) [4](#)
2. Abuadba, S., Kim, K., Kim, M., Thapa, C., Camtepe, S.A., Gao, Y., Kim, H., Nepal, S.: Can we use split learning on 1d cnn models for privacy preserving training? In: Proceedings of the 15th ACM Asia Conference on Computer and Communications Security. pp. 305–318 (2020) [2](#)
3. Ahmad, S., Scheinkman, L.: How can we be so dense? the benefits of using highly sparse representations. arXiv preprint arXiv:1903.11257 (2019) [2](#)
4. Aïvodji, U., Gambs, S., Ther, T.: Gamin: An adversarial approach to black-box model inversion. arXiv preprint arXiv:1909.11835 (2019) [2](#), [23](#)
5. An, S., Tao, G., Xu, Q., Liu, Y., Shen, G., Yao, Y., Xu, J., Zhang, X.: Mirror: Model inversion for deep learning network with high fidelity. In: Proceedings of the 29th Network and Distributed System Security Symposium (2022) [2](#)
6. Barlow, H.B.: The coding of sensory messages. *Current problems in animal behavior* (1961) [2](#)
7. Bonawitz, K., Eichner, H., Grieskamp, W., Huba, D., Ingerman, A., Ivanov, V., Kiddon, C., Konečný, J., Mazzocchi, S., McMahan, B., et al.: Towards federated learning at scale: System design. *Proceedings of machine learning and systems* **1**, 374–388 (2019) [4](#)
8. Breuer, A., Balkanski, E., Singer, Y.: The fast algorithm for submodular maximization. In: International Conference on Machine Learning. pp. 1134–1143. PMLR (2020) [3](#), [6](#)
9. Breuer, A., Khosravani, N., Tingley, M., Cattel, B.: Preemptive detection of fake accounts on social networks via multi-class preferential attachment classifiers. In: Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. pp. 105–116 (2023) [24](#)
10. Candès, E.J., Donoho, D.L.: New tight frames of curvelets and optimal representations of objects with piecewise  $c^2$  singularities. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences* **57**(2), 219–266 (2004) [2](#), [3](#)
11. Carlini, N., Jagielski, M., Zhang, C., Papernot, N., Terzis, A., Tramer, F.: The privacy onion effect: Memorization is relative. *Advances in Neural Information Processing Systems* **35**, 13263–13276 (2022) [2](#)
12. Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., et al.: Extracting training data from large language models. In: 30th USENIX Security Symposium (USENIX Security 21). pp. 2633–2650 (2021) [23](#)

13. Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwag, V., Tramer, F., Balle, B., Ippolito, D., Wallace, E.: Extracting training data from diffusion models. In: 32nd USENIX Security Symposium (USENIX Security 23). pp. 5253–5270 (2023) [1](#), [23](#)
14. Chen, S.S., Donoho, D.L., Saunders, M.A.: Atomic decomposition by basis pursuit. *SIAM review* **43**(1), 129–159 (2001) [2](#)
15. Chen, Y., Dey, T., Kuhnle, A.: Best of both worlds: Practical and theoretically optimal submodular maximization in parallel. *Advances in Neural Information Processing Systems* **34**, 25528–25539 (2021) [3](#), [6](#)
16. Choquette-Choo, C.A., Tramer, F., Carlini, N., Papernot, N.: Label-only membership inference attacks. In: International conference on machine learning. pp. 1964–1974. PMLR (2021) [23](#), [24](#)
17. Davis, G., Mallat, S., Avellaneda, M.: Adaptive greedy approximations. *Constructive approximation* **13**, 57–98 (1997) [3](#)
18. Dibbo, S.V.: Sok: Model inversion attack landscape: Taxonomy, challenges, and future roadmap. In: IEEE 36th Computer Security Foundations Symposium. pp. 408–425. IEEE Computer Society (2023) [1](#), [23](#), [24](#)
19. Dibbo, S.V., Chung, D.L., Mehnaz, S.: Model inversion attack with least information and an in-depth analysis of its disparate vulnerability. In: 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML). pp. 119–135. IEEE (2023) [1](#), [23](#)
20. Dibbo, S.V., Moore, J.S., Kenyon, G.T., Teti, M.A.: Lcanets++: Robust audio classification using multi-layer neural networks with lateral competition. arXiv preprint arXiv:2308.12882 (2023) [2](#)
21. Fang, H., Chen, B., Wang, X., Wang, Z., Xia, S.T.: Gifd: A generative gradient inversion method with feature domain optimization. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4967–4976 (2023) [2](#), [4](#), [24](#)
22. Field, D.J.: What is the goal of sensory coding? *Neural computation* **6**(4), 559–601 (1994) [2](#), [3](#)
23. Fredrikson, M., Jha, S., Ristenpart, T.: Model inversion attacks that exploit confidence information and basic countermeasures. In: Proceedings of the 22nd ACM SIGSAC conference on computer and communications security. pp. 1322–1333 (2015) [1](#), [2](#), [8](#), [24](#)
24. Gong, N.Z., Liu, B.: You are who you know and how you behave: Attribute inference attacks via users’ social friends and behaviors. In: 25th USENIX Security Symposium (USENIX Security 16). pp. 979–995 (2016) [1](#), [23](#)
25. Gong, X., Wang, Z., Li, S., Chen, Y., Wang, Q.: A gan-based defense framework against model inversion attacks. *IEEE Transactions on Information Forensics and Security* (2023) [2](#), [6](#), [7](#), [8](#), [10](#), [11](#), [12](#), [13](#), [22](#), [24](#), [26](#), [27](#), [28](#), [31](#), [32](#)
26. Haim, N., Vardi, G., Yehudai, G., Shamir, O., Irani, M.: Reconstructing training data from trained neural networks. *Advances in Neural Information Processing Systems* **35**, 22911–22924 (2022) [2](#)
27. Hannan, D., Nesbit, S.C., Wen, X., Smith, G., Zhang, Q., Goffi, A., Chan, V., Morris, M.J., Hunninghake, J.C., Villalobos, N.E., et al.: Mobileptx: sparse coding for pneumothorax detection given limited training examples. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 15675–15681 (2023) [2](#)
28. Hayes, J., Mahloujifar, S., Balle, B.: Bounding training data reconstruction in dp-sgd. arXiv preprint arXiv:2302.07225 (2023) [2](#), [4](#), [7](#), [10](#), [11](#), [12](#), [13](#), [26](#), [27](#), [28](#), [31](#), [32](#)

29. He, Z., Zhang, T., Lee, R.B.: Model inversion attacks against collaborative inference. In: Proceedings of the 35th Annual Computer Security Applications Conference. pp. 148–162 (2019) [2](#), [4](#), [7](#), [24](#)
30. Heredia, L.G., Negrevergne, B., Chevaleyre, Y.: Adversarial attacks for mixtures of classifiers. arXiv preprint arXiv:2307.10788 (2023) [8](#)
31. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* **30** (2017) [8](#)
32. Hidano, S., Murakami, T., Katsumata, S., Kiyomoto, S., Hanaoka, G.: Model inversion attacks for prediction systems: Without knowledge of non-sensitive attributes. In: 2017 15th Annual Conference on Privacy, Security and Trust (PST). pp. 115–11509. IEEE (2017) [8](#)
33. Hitaj, B., Ateniese, G., Perez-Cruz, F.: Deep models under the gan: information leakage from collaborative deep learning. In: Proceedings of the 2017 ACM SIGSAC conference on computer and communications security. pp. 603–618 (2017) [2](#)
34. Hu, H., Salcic, Z., Sun, L., Dobbie, G., Yu, P.S., Zhang, X.: Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)* **54**(11s), 1–37 (2022) [1](#), [24](#)
35. Jia, J., Gong, N.Z.: {AttriGuard}: A practical defense against attribute inference attacks via adversarial machine learning. In: 27th USENIX Security Symposium (USENIX Security 18). pp. 513–529 (2018) [24](#)
36. Jiang, Z., Zhang, G., Davis, L.S.: Submodular dictionary learning for sparse coding. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 3418–3425. IEEE (2012) [3](#), [6](#)
37. Juuti, M., Szyller, S., Marchal, S., Asokan, N.: Prada: protecting against dnn model stealing attacks. In: 2019 IEEE European Symposium on Security and Privacy (EuroS&P). pp. 512–527. IEEE (2019) [24](#)
38. Kahla, M., Chen, S., Just, H.A., Jia, R.: Label-only model inversion attacks via boundary repulsion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15045–15053 (2022) [2](#)
39. Kaissis, G.A., Makowski, M.R., Rückert, D., Braren, R.F.: Secure, privacy-preserving and federated machine learning in medical imaging. *Nature Machine Intelligence* **2**(6), 305–311 (2020) [4](#)
40. Kariyappa, S., Prakash, A., Qureshi, M.K.: Maze: Data-free model stealing attack using zeroth-order gradient estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13814–13823 (2021) [1](#)
41. Karras, T., Aittala, M., Laine, S., Härkönen, E., Hellsten, J., Lehtinen, J., Aila, T.: Alias-free generative adversarial networks. In: Proc. NeurIPS (2021) [23](#)
42. Kavukcuoglu, K., Ranzato, M., LeCun, Y.: Fast inference in sparse coding algorithms with applications to object recognition. arXiv preprint arXiv:1010.3467 (2010) [2](#), [3](#)
43. Kim, E., Rego, J., Watkins, Y., Kenyon, G.T.: Modeling biological immunity to adversarial examples. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4666–4675 (2020) [2](#), [5](#), [21](#)
44. Konečný, J., McMahan, H.B., Yu, F.X., Richtárik, P., Suresh, A.T., Bacon, D.: Federated learning: Strategies for improving communication efficiency. arXiv preprint arXiv:1610.05492 (2016) [4](#)
45. Krause, A., Cevher, V.: Submodular dictionary selection for sparse representation. In: International Conference on Machine Learning (ICML) (2010) [2](#), [3](#)
46. Lee, H., Battle, A., Raina, R., Ng, A.: Efficient sparse coding algorithms. *Advances in neural information processing systems* **19** (2006) [3](#)

47. Li, L., Xie, T., Li, B.: Sok: Certified robustness for deep neural networks. In: 2023 IEEE Symposium on Security and Privacy (SP). pp. 1289–1310. IEEE (2023) [1](#), [23](#)
48. Liu, G., Wang, C., Peng, K., Huang, H., Li, Y., Cheng, W.: Socinf: Membership inference attacks on social media health data with machine learning. *IEEE Transactions on Computational Social Systems* **6**(5), 907–921 (2019) [8](#)
49. Liu, Y., Wen, R., He, X., Salem, A., Zhang, Z., Backes, M., De Cristofaro, E., Fritz, M., Zhang, Y.: {ML-Doctor}: Holistic risk assessment of inference attacks against machine learning models. In: 31st USENIX Security Symposium (USENIX Security 22). pp. 4525–4542 (2022) [23](#)
50. McInnes, L., Healy, J., Melville, J.: Umap: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426 (2018) [13](#)
51. McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communication-efficient learning of deep networks from decentralized data. In: Artificial intelligence and statistics. pp. 1273–1282. PMLR (2017) [4](#)
52. Mehnaz, S., Dibbo, S.V., Kabir, E., Li, N., Bertino, E.: Are your sensitive attributes private? novel model inversion attribute inference attacks on classification models. In: 31st USENIX Security Symposium (USENIX Security 22). pp. 4579–4596. USENIX Association, Boston, MA (Aug 2022) [1](#), [2](#), [24](#)
53. Melis, L., Song, C., De Cristofaro, E., Shmatikov, V.: Exploiting unintended feature leakage in collaborative learning. In: 2019 IEEE symposium on security and privacy (SP). pp. 691–706. IEEE (2019) [2](#)
54. Mireshghallah, F., Taram, M., Ramrakhiani, P., Jalali, A., Tullsen, D., Esmailzadeh, H.: Shredder: Learning noise distributions to protect inference privacy. In: Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems. pp. 3–18 (2020) [2](#)
55. Mirzasoleiman, B., Badanidiyuru, A., Karbasi, A., Vondrák, J., Krause, A.: Lazier than lazy greedy. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 29 (2015) [3](#)
56. Natarajan, B.K.: Sparse approximate solutions to linear systems. *SIAM journal on computing* **24**(2), 227–234 (1995) [3](#)
57. Naveed, M., Ayday, E., Clayton, E.W., Fellay, J., Gunter, C.A., Hubaux, J.P., Malin, B.A., Wang, X.: Privacy in the genomic era. *ACM Computing Surveys (CSUR)* **48**(1), 1–44 (2015) [24](#)
58. Olshausen, B.A., Field, D.J.: Sparse coding of sensory inputs. *Current opinion in neurobiology* **14**(4), 481–487 (2004) [2](#), [3](#)
59. Olshausen, B.A., Field, D.J., et al.: Sparse coding of natural images produces localized, oriented, bandpass receptive fields. Submitted to Nature. Available electronically as ftp://redwood.psych.cornell.edu/pub/papers/sparse-coding.ps (1995) [2](#)
60. Paiton, D.M., Frye, C.G., Lundquist, S.Y., Bowen, J.D., Zarccone, R., Olshausen, B.A.: Selectivity and robustness of sparse coding networks. *Journal of vision* **20**(12), 10–10 (2020) [2](#)
61. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019) [21](#)
62. Peng, X., Liu, F., Zhang, J., Lan, L., Ye, J., Liu, T., Han, B.: Bilateral dependency optimization: Defending against model-inversion attacks. In: Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. pp. 1358–1367 (2022) [2](#), [6](#), [7](#), [10](#), [11](#), [12](#), [13](#), [26](#), [27](#), [28](#), [31](#), [32](#)

63. Rigaki, M., Garcia, S.: A survey of privacy attacks in machine learning. *ACM Computing Surveys* (2020) [2](#)
64. Rozell, C.J., Johnson, D.H., Baraniuk, R.G., Olshausen, B.A.: Sparse coding via thresholding and local competition in neural circuits. *Neural computation* **20**(10), 2526–2563 (2008) [2](#), [3](#), [5](#), [21](#)
65. Sablayrolles, A., Douze, M., Schmid, C., Ollivier, Y., Jégou, H.: White-box vs black-box: Bayes optimal strategies for membership inference. In: *International Conference on Machine Learning*. pp. 5558–5567. PMLR (2019) [23](#)
66. Salem, A., Bhattacharya, A., Backes, M., Fritz, M., Zhang, Y.: {Updates-Leak}: Data set inference and reconstruction attacks in online learning. In: *29th USENIX security symposium (USENIX Security 20)*. pp. 1291–1308 (2020) [2](#)
67. Sannai, A.: Reconstruction of training samples from loss functions. *arXiv preprint arXiv:1805.07337* (2018) [8](#)
68. Sanyal, S., Addepalli, S., Babu, R.V.: Towards data-free model stealing in a hard label setting. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 15284–15293 (2022) [1](#)
69. Schneiderman, H.: Feature-centric evaluation for efficient cascaded object detection. In: *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004. vol. 2*, pp. II–II. IEEE (2004) [2](#)
70. Shokri, R., Stronati, M., Song, C., Shmatikov, V.: Membership inference attacks against machine learning models. In: *2017 IEEE symposium on security and privacy (SP)*. pp. 3–18. IEEE (2017) [23](#)
71. Song, L., Mittal, P.: Systematic evaluation of privacy risks of machine learning models. In: *30th USENIX Security Symposium (USENIX Security 21)*. pp. 2615–2632 (2021) [4](#)
72. Struppek, L., Hintersdorf, D., Correia, A.D.A., Adler, A., Kersting, K.: Plug & play attacks: Towards robust and flexible model inversion attacks. In: *International Conference on Machine Learning*. pp. 20522–20545. PMLR (2022) [1](#), [2](#), [4](#), [7](#), [9](#), [10](#), [23](#), [25](#), [26](#), [27](#), [29](#), [31](#)
73. Sun, B., Tsai, N.h., Liu, F., Yu, R., Su, H.: Adversarial defense by stratified convolutional sparse coding. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 11447–11456 (2019) [2](#)
74. Teti, M., Kenyon, G., Migliori, B., Moore, J.: Lcanets: Lateral competition improves robustness against corruption and attack. In: *International Conference on Machine Learning*. pp. 21232–21252. PMLR (2022) [2](#), [5](#), [6](#), [7](#), [21](#)
75. Titcombe, T., Hall, A.J., Papadopoulos, P., Romanini, D.: Practical defences against model inversion attacks for split neural networks. *arXiv preprint arXiv:2104.05743* (2021) [2](#), [4](#), [7](#), [8](#), [10](#), [11](#), [12](#), [13](#), [22](#), [24](#), [25](#), [26](#), [27](#), [28](#), [29](#), [30](#), [31](#), [32](#)
76. Tramèr, F., Shokri, R., San Joaquin, A., Le, H., Jagielski, M., Hong, S., Carlini, N.: Truth serum: Poisoning machine learning models to reveal their secrets. In: *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*. pp. 2779–2792 (2022) [23](#)
77. Vepakomma, P., Gupta, O., Swedish, T., Raskar, R.: Split learning for health: Distributed deep learning without sharing raw patient data. *arXiv preprint arXiv:1812.00564* (2018) [4](#)
78. Vhaduri, S., Cheung, W., Dibbo, S.V.: Bag of on-phone anns to secure iot objects using wearable and smartphone biometrics. *IEEE Transactions on Dependable and Secure Computing* (2023) [23](#)

79. Vhaduri, S., Dibbo, S.V., Chen, C.Y.: Predicting a user’s demographic identity from leaked samples of health-tracking wearables and understanding associated risks. In: 2022 IEEE 10th International Conference on Healthcare Informatics (ICHI). pp. 309–318. IEEE (2022) [23](#)
80. Vhaduri, S., Dibbo, S.V., Cheung, W.: Hiauth: A hierarchical implicit authentication system for iot wearables using multiple biometrics. *IEEE Access* **9**, 116395–116406 (2021) [23](#)
81. Wang, K.C., Fu, Y., Li, K., Khisti, A., Zemel, R., Makhzani, A.: Variational model inversion attacks. *Advances in Neural Information Processing Systems* **34**, 9706–9719 (2021) [24](#)
82. Wang, T., Zhang, Y., Jia, R.: Improving robustness to model inversion attacks via mutual information regularization. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 35, pp. 11666–11673 (2021) [2](#), [7](#), [10](#), [11](#), [12](#), [13](#), [26](#), [27](#), [28](#), [31](#), [32](#)
83. Wang, X., Wang, W.H.: Group property inference attacks against graph neural networks. In: *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*. pp. 2871–2884 (2022) [4](#)
84. Wang, Y., Qian, H., Miao, C.: Dualcf: Efficient model extraction attack from counterfactual explanations. In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. pp. 1318–1329 (2022) [1](#), [24](#)
85. Wang, Z., Song, M., Zhang, Z., Song, Y., Wang, Q., Qi, H.: Beyond inferring class representatives: User-level privacy leakage from federated learning. In: *IEEE INFOCOM 2019-IEEE conference on computer communications*. pp. 2512–2520. IEEE (2019) [2](#)
86. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004) [8](#)
87. Wei, W., Liu, L., Loper, M., Chow, K.H., Gursoy, M.E., Truex, S., Wu, Y.: A framework for evaluating gradient leakage attacks in federated learning. *arXiv preprint arXiv:2004.10397* (2020) [2](#)
88. Wu, P., Liu, J., Li, M., Sun, Y., Shen, F.: Fast sparse coding networks for anomaly detection in videos. *Pattern Recognition* **107**, 107515 (2020) [6](#)
89. Wu, Y., Yu, N., Li, Z., Backes, M., Zhang, Y.: Membership inference attacks against text-to-image generation models. *arXiv preprint arXiv:2210.00968* (2022) [8](#)
90. Xu, T., Goossen, G., Cevahir, H.K., Khodeir, S., Jin, Y., Li, F., Shan, S., Patel, S., Freeman, D., Pearce, P.: Deep entity classification: Abusive account detection for online social networks. In: *30th {USENIX} Security Symposium ({USENIX} Security 21)* (2021) [24](#)
91. Xu, Y., Liu, X., Hu, T., Xin, B., Yang, R.: Sparse black-box inversion attack with limited information. In: *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 1–5. IEEE (2023) [23](#)
92. Yang, Z., Zhang, J., Chang, E.C., Liang, Z.: Neural network inversion in adversarial setting via background knowledge alignment. In: *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. pp. 225–240 (2019) [2](#)
93. Yuan, X., Ding, L., Zhang, L., Li, X., Wu, D.O.: Es attack: Model stealing against deep neural networks without data hurdles. *IEEE Transactions on Emerging Topics in Computational Intelligence* **6**(5), 1258–1270 (2022) [1](#), [24](#)
94. Zhang, J., Peng, S., Gao, Y., Zhang, Z., Hong, Q.: Apmsa: adversarial perturbation against model stealing attacks. *IEEE Transactions on Information Forensics and Security* **18**, 1667–1679 (2023) [1](#)

95. Zhang, Y., Jia, R., Pei, H., Wang, W., Li, B., Song, D.: The secret revealer: Generative model-inversion attacks against deep neural networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 253–261 (2020) [2](#), [23](#), [24](#)
96. Zhao, B.Z.H., Agrawal, A., Coburn, C., Asghar, H.J., Bhaskar, R., Kaafar, M.A., Webb, D., Dickinson, P.: On the (in) feasibility of attribute inference attacks on machine learning models. In: 2021 IEEE European Symposium on Security and Privacy (EuroS&P). pp. 232–251. IEEE (2021) [24](#)
97. Zhao, X., Zhang, W., Xiao, X., Lim, B.: Exploiting explanations for model inversion attacks. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 682–692 (2021) [24](#)
98. Zhong, D., Sun, H., Xu, J., Gong, N., Wang, W.H.: Understanding disparate effects of membership inference attacks and their countermeasures. In: Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security. pp. 959–974 (2022) [1](#), [23](#)

## A Appendix

This is the supplementary document containing the additional results and details of our proposed Sparse Coding Architecture (SCA) formulations, as well as cluster details and additional preliminaries.

### A.1 Reproducibility

In order to promote further research and standardize the evaluations of new defenses, we provide full cluster-ready PyTorch [\[61\]](#) implementations of SCA and all benchmarks as well as replication codes for all experiments on our project page at: <https://sayantondibbo.github.io/SCA>.

We provide full details of the cluster hardware and all parameter choices used in our experiments in *Appendix A.3* and *A.4*, and in *Appendix Tables 4* and *5*.

### A.2 Adapting Rozell LCA to Convolutional Networks

Although the original LCA formulation [\[64\]](#) was introduced for the non-convolutional case, it is based on the general principle of feature-similarity-based competition between neurons within the same layer, which can be adapted to the convolutional setting via only two minimal changes to Equation [3](#) [\[43, 74\]](#). In Rozell’s original formulation,  $\Psi(t)$  can simply be recast from a matrix multiplication to a convolution between the input and dictionary. Second, the lateral interaction tensor,  $\mathcal{G}$  in Equation [3](#), can also be recast from a matrix multiplication to a convolution between the dictionary and its transpose. Neuron membrane potential works as follows:

$$\dot{\mathcal{P}}(t) = \frac{1}{\tau}[\Psi(t) - \mathcal{P}(t) - \mathcal{R}_x(t) * \mathcal{G}] \quad (3)$$

where  $\tau$  is a time constant,  $\Psi(t) = \mathcal{X} * \Omega$  is the neuron’s bottom-up drive from the input computed by taking the convolution,  $*$ , between the input,  $\mathcal{X}$ , and the dictionary,  $\Omega$ , and  $-\mathcal{P}(t)$  is the leak term [\[43, 74\]](#).

### A.3 Cluster Details

We run all our experiments using the slurm batch jobs on industry-standard high-performance GPU clusters with 40 cores and 4 nodes. Details of the hardware and architecture of our cluster are described in Table 4. We note that noise-based GAUSSIAN and Titcombe et al. [75] defenses are typically fastest on this architecture (though they are the least-performant). We emphasize that our sparse coding implementations are ‘research-grade’, unlike the optimized torch GAN implementations available for [25]. See also *Appendix G*. Note that for large scale applications, SCA’s sparse coding updates can be accelerated such that they can be computed extremely efficiently (see the training complexity discussion in the main paper body).

**Table 4:** Hardware Details of the Cluster in our Experiments.

<i>Parameter</i> MEASUREMENTS	
Core	40
RAM	565GB
GPU	Tesla V100
Nodes	p01-p04
Space	1.5TB

### A.4 Parameters and architecture of SCA

We implement SCA using two Sparse Coding Layers (SCL): One following the input image, and one following a downstream dense batch normalization layer. Finally, we follow these two pairs of dense-then-sparse layers with downstream fully connected (linear) layers before the classification layer. In the case of end-to-end network experiments, we use 5 downstream linear layers, which is a reasonable default. In the split network setting, we are careful to use 3 downstream fully connected layers in order to match the architectures used in the split network experimental setup of [75], and per our public codebase, we make every effort to make the benchmarks within each setting comparable in terms of architecture, aside from the obvious difference of SCA’s sparse layers. We train SCA’s sparse layers with 500 iterations of lateral competitions during reconstructions in SCL layers. We emphasize that SCA can be made significantly more complex, either via the addition of more sparse-dense pairs of layers, or by adding additional (convolutional, linear) downstream layers before classification. We avoid such complexity in the experiments in order to compare more directly to benchmarks and because our goal is to study an architecture that captures the essence of SCA. We give all parameter and training details in Table 5.

**Table 5:** Architecture and Parameters of SCA implementation.

<i>Parameter</i>	VALUE
Sparse Layers	2
Batch Norm Layers	2
Fully Connected Layers	5
$\lambda$	0.5
Learning rate $\eta$	0.01
Time constant $\tau$	1000
Kernel size	5
Stride	1,1
Lateral competition iterations	500

### A.5 Attack details

In the Plug-&-Play attack experiments, we follow the authors’ attack exactly [72], except we update their approach to use the latest **StyleGAN3** [41] for high-resolution image generation. For the end-to-end and split-network attacks, we consider a recent state-of-the-art surrogate model training attack optimized via SGD [4, 91]. This attack works by querying the target model with an externally obtained dataset. To capture a well-informed ‘worst-case’ attacker, we set this dataset to a holdout set from the true training dataset. The attack then uses the corresponding model high-dimensional intermediate outputs to train an inverted surrogate model that outputs actual training data.

### A.6 SCA sparsity vs. robustness

We vary the sparsity, i.e.,  $\lambda$  parameter and run the SPARSE-STANDARD, as well as our SCA. We observe that increasing  $\lambda$  helps improve the robustness, without significant accuracy drops. For example, Table 6 shows this comparison for MNIST in the end-to-end setting.

## B Model Inversion Attack Methodology: Additional discussion

Because privacy attacks are an emerging field, we feel it is relevant to include additional context and discussion here. Recent work has highlighted a variety of attack vectors targeting sensitive training data of machine learning models [12, 13, 16, 18, 19, 19, 24, 47, 49, 65, 70, 76, 78–80, 95, 98]. These attacks not only target centralized models but also can make the federated learning models

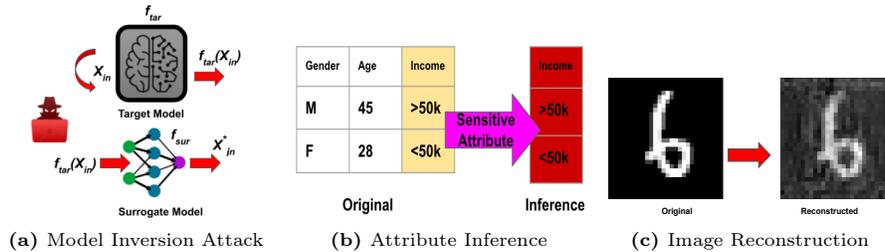
**Table 6:** SPARSE-STANDARD and SCA performance with  $\lambda \in \{0.1, 0.25, 0.5, 0.75\}$ 

$\lambda$	PSNR $\downarrow\downarrow$		SSIM $\downarrow\downarrow$		FID $\uparrow\uparrow$		Accuracy	
	SP-STD	SCA	SP-STD	SCA	SP-STD	SCA	SP-STD	SCA
0.1	23.45	19.54	0.650	0.502	111.5	178.5	0.984	0.984
0.25	21.34	18.81	0.438	0.340	142.9	174.1	0.986	0.983
0.5	22.16	17.85	0.598	0.164	136.9	<b>335.4</b>	0.985	0.977
0.75	22.39	<b>14.65</b>	0.593	<b>0.086</b>	142.0	214.1	0.981	0.971

vulnerable to attacks [21, 29]. Adversaries with different access (i.e., black-box, white-box) to these models perform different attacks leveraging a wide range of capabilities, e.g., knowledge about the target model confusion matrix and access to blurred images of that particular class [16, 23, 29, 37, 81]. Such attacks commonly fall under the umbrella of privacy attacks, which include specific attacker goals such as membership inference, model stealing, model inversion, etc. [34, 52, 84, 93]. Defending against privacy attacks is a core task of mainstream technology platforms ranging from public social networks to private medical research [9, 57, 90].

Our focus is model inversion attack, where an adversary aims to infer sensitive training data attributes  $X_s$  or reconstruct training samples  $X_{in}$ , a severe threat to the privacy of training data  $D_{Tr}$  [52, 75]. In Figure 7a, we present the pipelines of the model inversion attack. Depending on data types and purpose, model inversion attacks can be divided into two broader categories: (i) attribute inference (AttrInf) and (ii) image reconstruction (ImRec) attacks [18]. In AttrInf attacks, it is assumed the adversary can query the target model  $f_{tar}$  and design a surrogate model  $f_{sur}$  to infer some sensitive attributes  $X_s$  in training data  $D_{Tr}$ , with or without knowing all other non-sensitive attributes training data  $X_{ns}$  in the training data  $D_{Tr}$ , as presented in Figure 7b. In ImRec attacks the adversary reconstructs entire training samples  $D_{Tr}$  using the surrogate model  $f_{sur}$  with or without having access to additional information like blurred, masked, or noisy training samples  $D_s$ , as shown in Figure 7c [23, 95, 97]. To contextualize our SCA setting, recall that we suppose the attacker has only black-box access to query the model  $f_{tar}$  without knowing the details of the target model  $f_{tar}$  architecture or parameters like gradient information  $\nabla_{Tr}$ . The attacker attempts to compute training data reconstruction (i.e., ImRec) attack without having access to other additional information, e.g., blurred or masked images  $D_s$ .

Two major components of the model inversion attack workflow are the target model  $f_{tar}$  and the surrogate attack model  $f_{sar}$  [18, 35, 96]. Training data reconstruction (i.e., ImRec) attack in the literature considers the target model  $f_{tar}$  to be either the split network [75] or the end-to-end network [25, 95]. In the split network  $f_{tar}$  model, the output of a particular layer  $l$  in the network, i.e.,  $a^{[l]}$ , where  $1 \leq l < L$  is accessible to the adversary, whereas, for the end-to-end network, the adversary does not have access to intermediate layer outputs; rather,



**Fig. 7:** Illustration of Model Inversion attack along with (a.) pipelines—an adversary queries target model  $f_{tar}$  with inputs  $\mathcal{X}_{in}$  to obtain output  $f_{tar}(\mathcal{X}_{in})$ . Then adversary trains a surrogate attack model  $f_{sur}$ , where the  $f_{tar}(\mathcal{X}_{in})$  is the input and  $\mathcal{X}^*$  is the output; and (b.) categories, i.e., attribute inference (AttrInf) attack, where the adversary infers sensitive attribute  $\mathcal{X}_s$  with or without knowing non-sensitive attribute values, i.e.,  $\mathcal{X}_{n_s} \rightarrow \mathcal{X}_s$  and (c.) image reconstruction (ImRec) attack, where adversary reconstructs similar to original images, i.e.,  $\mathcal{X}_{in} \approx \mathcal{X}_{in}^*$ .

the adversary only has access to the output from the last hidden layer before the classification layer  $a^{[L]}$ .

## C Results of extra {threat model, dataset} experiments

We experiment all 3 attack setups: *Plug- $\mathcal{E}$ -Play* model inversion attack [72], *end-to-end*, and *split* on three additional datasets: MNIST, Fashion MNIST, and CIFAR10. We experiment with all benchmarks and present the results on Tables 7, 8, and 9. In all of these additional datasets, SCA consistently outperforms all benchmarks.

## D Additional baseline tuning

We also attempt to improve the Laplace noise-based defense of Titcombe et al. [75] by increasing the noise scale parameter  $b$  from  $\mathcal{L}(\mu=0, b=0.5)$  to  $\mathcal{L}(\mu=0, b=1.0)$ . Tables 10, 11, and 12 compare these results to SCA for in all 3 attack settings. Observe that the additional noise significantly degrades classification accuracy in all but one case, yet it does not result in reconstruction metrics that rival those of SCA’s. In Figure 8, we present the reconstructed images in the Split network attack setting on MNIST data. We also include the Laplace noise-based defense with higher noise parameter  $\mathcal{L}(\mu=0, b=1.0)$ .

## E Stability analysis of SCA

Tables 13 and 14 show mean metrics and std. deviation error bars taken over *multiple runs* of each defense. Observe that SCA is at least as stable (and in some cases significantly more stable) than alternatives.

**Table 7:** Experiments set 1 **Additional Datasets:** Performance in Plug-&-Play Model Inversion Attack [72] setting (*lower rows=better defense*).

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
MNIST	NO-DEFENSE	7.24	0.783	23.6	0.971
	GAUSSIAN-NOISE	6.94	0.686	31.22	0.958
	GAN	6.83	0.734	89.38	0.968
	Gong et al. [25]++	6.69	0.716	92.21	0.987
	Titcombe et al. [75]	6.34	0.744	131.8	0.980
	Gong et al. [25]	6.76	0.681	99.53	0.985
	Peng et al. [62]	6.89	0.704	283.8	0.941
	Hayes et al. [28]	7.03	0.672	396.1	0.871
	Wang et al. [82]	7.14	0.752	261.2	0.937
	SPARSE-STANDARD	6.24	0.631	158.6	0.986
	<b>SCA0.1</b>	<b>6.19</b>	<b>0.633</b>	<b>287.9</b>	<b>0.984</b>
	<b>SCA0.25</b>	<b>5.83</b>	<b>0.607</b>	<b>289.3</b>	<b>0.983</b>
<b>SCA0.5</b>	<b>5.74</b>	<b>0.604</b>	<b>299.6</b>	<b>0.977</b>	
Fashion MNIST	NO-DEFENSE	8.91	0.147	235.5	0.886
	GAUSSIAN-NOISE	8.67	0.132	239.8	0.815
	GAN	8.66	0.147	243.3	0.883
	Gong et al. [25]++	8.73	0.130	220.2	0.906
	Titcombe et al. [75]	8.56	0.134	229.8	0.905
	Gong et al. [25]	8.57	0.143	244.3	0.888
	Peng et al. [62]	8.85	0.147	227.5	0.845
	Hayes et al. [28]	8.63	0.139	218.4	0.752
	Wang et al. [82]	8.90	0.119	210.3	0.880
	SPARSE-STANDARD	8.71	0.135	223.3	0.879
	<b>SCA0.1</b>	<b>8.49</b>	<b>0.039</b>	<b>222.8</b>	<b>0.897</b>
	<b>SCA0.25</b>	<b>8.49</b>	<b>0.032</b>	<b>229.9</b>	<b>0.887</b>
<b>SCA0.5</b>	<b>8.45</b>	<b>0.047</b>	<b>233.5</b>	<b>0.876</b>	
CIFAR10	NO-DEFENSE	11.94	0.381	39.38	0.821
	GAUSSIAN-NOISE	11.88	0.365	77.92	0.626
	GAN	11.86	0.369	88.39	0.596
	Titcombe et al. [75]	10.89	0.346	79.19	0.792
	Gong et al. [25]++	11.06	0.339	78.48	0.773
	Gong et al. [25]	11.21	0.334	92.33	0.682
	Peng et al. [62]	11.96	0.354	120.5	0.752
	Hayes et al. [28]	11.12	0.342	142.1	0.626
	Wang et al. [82]	11.02	0.346	142.6	0.756
	SPARSE-STANDARD	10.74	0.303	137.4	0.790
	<b>SCA0.1</b>	<b>10.59</b>	<b>0.305</b>	<b>144.1</b>	<b>0.787</b>
	<b>SCA0.25</b>	<b>10.27</b>	<b>0.279</b>	<b>189.9</b>	<b>0.772</b>
<b>SCA0.5</b>	<b>10.23</b>	<b>0.276</b>	<b>189.7</b>	<b>0.744</b>	

## F Compute time

Our basic SCA research implementation completes in comparable or less compute time than highly optimized implementations of benchmarks. In the ‘worst-

**Table 8:** Experiments set 2 **Additional Datasets:** Performance in *end-to-end* network setting (*lower rows=better defense*).

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
MNIST	No-DEFENSE	40.87	0.982	16.31	0.971
	GAUSSIAN-NOISE	40.88	0.983	15.88	0.958
	GAN	40.69	0.981	16.59	0.968
	Titcombe et al. [75]	31.18	0.863	47.32	0.980
	Gong et al. [25] <sup>++</sup>	30.37	0.838	72.99	0.987
	Gong et al. [25]	29.05	0.817	75.39	0.985
	Peng et al. [62]	18.44	0.354	111.6	0.968
	Hayes et al. [28]	19.75	0.488	298.8	0.871
	Wang et al. [82]	27.26	0.862	72.66	0.962
	SPARSE-STANDARD	21.34	0.439	142.9	0.986
	<b>SCA0.1</b>	<b>19.54</b>	<b>0.502</b>	<b>178.5</b>	<b>0.984</b>
<b>SCA0.25</b>	<b>18.81</b>	<b>0.340</b>	<b>174.1</b>	<b>0.983</b>	
<b>SCA0.5</b>	<b>17.85</b>	<b>0.164</b>	<b>335.5</b>	<b>0.977</b>	
Fashion MNIST	No-DEFENSE	37.86	0.975	13.91	0.886
	GAUSSIAN-NOISE	36.54	0.969	16.49	0.815
	GAN	37.68	0.974	19.26	0.883
	Gong et al. [25] <sup>++</sup>	27.71	0.794	41.35	0.906
	Titcombe et al. [75]	26.66	0.759	53.76	0.905
	Gong et al. [25]	21.24	0.523	93.08	0.888
	Peng et al. [62]	17.98	0.368	70.53	0.880
	Hayes et al. [28]	21.13	0.297	223.3	0.752
	Wang et al. [82]	25.98	0.806	41.87	0.838
	SPARSE-STANDARD	19.35	0.446	128.4	0.879
	<b>SCA0.1</b>	<b>17.92</b>	<b>0.209</b>	<b>196.1</b>	<b>0.897</b>
<b>SCA0.25</b>	<b>17.03</b>	<b>0.186</b>	<b>195.2</b>	<b>0.887</b>	
<b>SCA0.5</b>	<b>14.51</b>	<b>0.069</b>	<b>423.2</b>	<b>0.876</b>	
CIFAR10	No-DEFENSE	21.17	0.477	70.96	0.821
	GAUSSIAN-NOISE	20.26	0.220	77.42	0.626
	GAN	19.71	0.259	132.0	0.596
	Titcombe et al. [75]	18.62	0.174	171.9	0.792
	Gong et al. [25] <sup>++</sup>	18.27	0.209	149.1	0.773
	Gong et al. [25]	19.10	0.150	133.8	0.682
	Peng et al. [62]	17.20	0.002	130.3	0.717
	Hayes et al. [28]	17.95	0.002	142.4	0.626
	Wang et al. [82]	17.08	0.002	136.1	0.793
	SPARSE-STANDARD	18.01	0.003	168.6	0.790
	<b>SCA0.1</b>	<b>17.09</b>	<b>0.001</b>	<b>172.0</b>	<b>0.787</b>
<b>SCA0.25</b>	<b>16.78</b>	<b>0.001</b>	<b>189.3</b>	<b>0.772</b>	
<b>SCA0.5</b>	<b>16.24</b>	<b>0.001</b>	<b>197.0</b>	<b>0.744</b>	

case’ across all of our experiments, SCA is faster than the best performing baseline (Peng et al. [62]) but slower than other baselines. Table 16 shows the compute times (in seconds) for this ‘worst-case’ experiment below (The MNIST dataset under the Plug-&-Play attack [72]).

**Table 9:** Experiments set 3 **Additional Datasets:** Performance in *split* network setting (*lower rows=better defense*).

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
MNIST	No-DEFENSE	31.21	0.923	19.64	0.963
	GAUSSIAN-NOISE	31.07	0.922	23.27	0.972
	GAN	28.39	0.894	27.26	0.969
	Gong et al. [25]	28.30	0.806	69.38	0.986
	Titcombe et al. [75]	25.40	0.713	76.88	0.952
	Gong et al. [25]++	21.94	0.591	97.33	0.991
	Peng et al. [62]	16.90	0.475	103.2	0.960
	Hayes et al. [28]	17.23	0.030	288.1	0.856
	Wang et al. [82]	21.87	0.696	53.09	0.903
	SPARSE-STANDARD	18.71	0.288	188.4	0.981
	<b>SCA0.1</b>	<b>16.17</b>	<b>0.109</b>	<b>227.4</b>	<b>0.988</b>
<b>SCA0.25</b>	<b>17.40</b>	<b>0.058</b>	<b>301.6</b>	<b>0.980</b>	
<b>SCA0.5</b>	<b>14.98</b>	<b>0.044</b>	<b>307.7</b>	<b>0.975</b>	
Fashion MNIST	No-DEFENSE	29.66	0.911	14.33	0.868
	GAUSSIAN-NOISE	29.49	0.909	14.81	0.871
	GAN	26.03	0.849	19.33	0.885
	Gong et al. [25]	23.70	0.631	97.52	0.884
	Titcombe et al. [75]	20.48	0.565	81.01	0.872
	Gong et al. [25]++	25.77	0.726	57.72	0.908
	Peng et al. [62]	20.67	0.583	46.48	0.865
	Hayes et al. [28]	20.10	0.256	200.6	0.748
	Wang et al. [82]	24.53	0.588	81.79	0.881
	SPARSE-STANDARD	19.54	0.405	200.5	0.882
	<b>SCA0.1</b>	<b>18.11</b>	<b>0.154</b>	<b>171.1</b>	<b>0.904</b>
<b>SCA0.25</b>	<b>17.74</b>	<b>0.188</b>	<b>203.8</b>	<b>0.896</b>	
<b>SCA0.5</b>	<b>17.15</b>	<b>0.134</b>	<b>270.4</b>	<b>0.879</b>	
CIFAR10	No-DEFENSE	16.48	0.709	47.77	0.823
	GAUSSIAN-NOISE	14.79	0.311	149.5	0.598
	GAN	14.87	0.296	13.01	0.675
	Titcombe et al. [75]	14.68	0.244	157.3	0.779
	Gong et al. [25]++	13.32	0.003	162.4	0.691
	Gong et al. [25]	14.55	0.291	152.1	0.644
	Peng et al. [62]	17.18	0.002	169.1	0.707
	Hayes et al. [28]	15.44	0.005	204.5	0.596
	Wang et al. [82]	14.73	0.001	176.3	0.820
	SPARSE-STANDARD	13.22	0.003	167.9	0.769
	<b>SCA0.1</b>	<b>13.18</b>	<b>0.002</b>	<b>174.2</b>	<b>0.758</b>
<b>SCA0.25</b>	<b>13.07</b>	<b>0.002</b>	<b>181.2</b>	<b>0.742</b>	
<b>SCA0.5</b>	<b>12.88</b>	<b>0.002</b>	<b>375.3</b>	<b>0.739</b>	

## G Ablations: Tuning SCA

Observe that our SCA outperforms SOTA defense baselines in robustness even without any tuning of parameters. However, tuning the hyper-parameters can

**Table 10:** Experiments set 1: additional Laplace noise benchmark with larger 1.0 noise parameter: Performance in Plug-&-Play Model Inversion Attack [72] setting (*lower rows=better defense*).

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
MNIST	Titcombe et al. [75]-1.0	6.60	0.685	280.1	0.938
	<b>SCA0.1</b>	<b>6.19</b>	<b>0.633</b>	<b>287.9</b>	<b>0.984</b>
	<b>SCA0.25</b>	<b>5.83</b>	<b>0.607</b>	<b>289.3</b>	<b>0.983</b>
	<b>SCA0.5</b>	<b>5.74</b>	<b>0.604</b>	<b>299.6</b>	<b>0.977</b>
Fashion MNIST	Titcombe et al. [75]-1.0	8.72	0.1412	232.1	0.823
	<b>SCA0.1</b>	<b>8.49</b>	<b>0.039</b>	<b>222.8</b>	<b>0.897</b>
	<b>SCA0.25</b>	<b>8.49</b>	<b>0.032</b>	<b>229.9</b>	<b>0.887</b>
	<b>SCA0.5</b>	<b>8.45</b>	<b>0.047</b>	<b>233.5</b>	<b>0.876</b>
CIAFR10	Titcombe et al. [75]-1.0	10.75	0.335	112.7	0.779
	<b>SCA0.1</b>	<b>10.59</b>	<b>0.305</b>	<b>144.1</b>	<b>0.787</b>
	<b>SCA0.25</b>	<b>10.27</b>	<b>0.279</b>	<b>189.9</b>	<b>0.772</b>
	<b>SCA0.5</b>	<b>10.23</b>	<b>0.276</b>	<b>189.7</b>	<b>0.744</b>

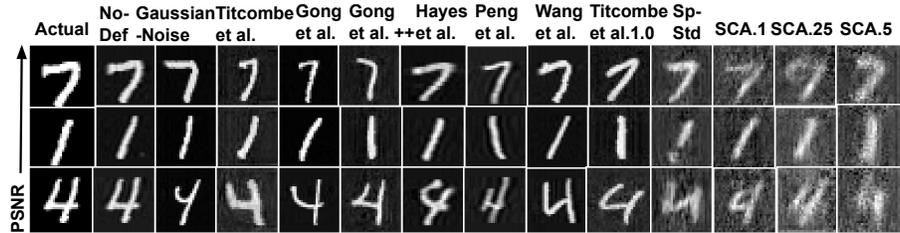
**Table 11:** Experiments set 2 additional Laplace noise benchmark with larger 1.0 noise parameter: Performance in *end-to-end* network setting (*lower rows=better defense*).

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
MNIST	Titcombe et al. [75]-1.0	24.89	0.664	50.64	0.938
	<b>SCA0.1</b>	<b>19.54</b>	<b>0.502</b>	<b>178.5</b>	<b>0.984</b>
	<b>SCA0.25</b>	<b>18.81</b>	<b>0.340</b>	<b>174.1</b>	<b>0.983</b>
	<b>SCA0.5</b>	<b>17.85</b>	<b>0.164</b>	<b>335.5</b>	<b>0.977</b>
Fashion MNIST	Titcombe et al. [75]-1.0	20.21	0.567	80.55	0.823
	<b>SCA0.1</b>	<b>17.92</b>	<b>0.209</b>	<b>196.1</b>	<b>0.897</b>
	<b>SCA0.25</b>	<b>17.03</b>	<b>0.186</b>	<b>195.2</b>	<b>0.887</b>
	<b>SCA0.5</b>	<b>14.51</b>	<b>0.069</b>	<b>423.2</b>	<b>0.876</b>
CIFAR10	Titcombe et al. [75]-1.0	18.71	0.672	170.8	0.779
	<b>SCA0.1</b>	<b>17.09</b>	<b>0.001</b>	<b>172.0</b>	<b>0.787</b>
	<b>SCA0.25</b>	<b>16.78</b>	<b>0.001</b>	<b>189.3</b>	<b>0.772</b>
	<b>SCA0.5</b>	<b>16.24</b>	<b>0.001</b>	<b>197.0</b>	<b>0.744</b>

boost the accuracy further, e.g., we use kernel size as default 5 for all experiments. Increasing the kernel from 5 to 7 can improve SCA accuracies beyond. While our goal is to capture the essence of the SCA itself in terms of robustness, we explore a little bit on further possible improvements on accuracy scores. We consider the lowest robust SCA, i.e., SCA0.1 for the tuning of kernel size, and we present the comparisons of accuracies between SCA0.1 and TUNED SCA0.1 in Table 15.

**Table 12:** Experiments set 3: additional Laplace noise benchmark with larger 1.0 noise parameter: Performance in *split* network setting (*lower rows=better defense*).

Dataset	Defense	PSNR ↓↓	SSIM ↓↓	FID ↑↑	Accuracy
MNIST	Titcombe et al. [75]-1.0	22.63	0.503	66.40	0.980
	<b>SCA0.1</b>	<b>16.17</b>	<b>0.109</b>	<b>227.4</b>	<b>0.988</b>
	<b>SCA0.25</b>	<b>17.40</b>	<b>0.058</b>	<b>301.6</b>	<b>0.980</b>
	<b>SCA0.5</b>	<b>14.98</b>	<b>0.044</b>	<b>307.7</b>	<b>0.975</b>
Fashion MNIST	Titcombe et al. [75]-1.0	18.36	0.408	80.80	0.878
	<b>SCA0.1</b>	<b>18.11</b>	<b>0.154</b>	<b>171.1</b>	<b>0.904</b>
	<b>SCA0.25</b>	<b>17.74</b>	<b>0.188</b>	<b>203.8</b>	<b>0.896</b>
	<b>SCA0.5</b>	<b>17.15</b>	<b>0.134</b>	<b>270.4</b>	<b>0.879</b>
CIAFR10	Titcombe et al. [75]-1.0	14.27	0.259	171.6	0.786
	<b>SCA0.1</b>	<b>13.18</b>	<b>0.002</b>	<b>174.2</b>	<b>0.758</b>
	<b>SCA0.25</b>	<b>13.07</b>	<b>0.002</b>	<b>181.2</b>	<b>0.742</b>
	<b>SCA0.5</b>	<b>12.88</b>	<b>0.002</b>	<b>375.3</b>	<b>0.739</b>

**Fig. 8:** Qualitative comparisons among actual and reconstructed images under SCA and additional Laplace noise defense benchmark with larger 1.0 noise parameter.

## H Robustness of sparse coding layers: UMap

In Figure 9, we present the UMap representation of linear, convolutional, and sparse coding layers on the other datasets, i.e., CelebA and Medical MNIST datasets. Observe that, the data points are more scattered in the sparse coding layer UMap (Figure 9c and Figure 9f) representations compared to the linear (Figure 9a and Figure 9d) and convolutional layers (Figure 9b and Figure 9e), which provide more robustness to models with sparse coding layers, i.e., our proposed SCA, against the privacy attacks.

**Table 13: Stability analysis 1:** Performance comparison (mean $\pm$  standard deviations) across multiple runs in Plug-&Play Model Inversion Attack [72] setting (*lower rows=better defense*) on high-res CelebA dataset.

Dataset	Defense	PSNR $\downarrow\downarrow$	SSIM $\downarrow\downarrow$	FID $\uparrow\uparrow$	Accuracy
CelebA	NO-DEFENSE	11.42 $\pm$ 2.44	0.613 $\pm$ 0.29	292.9 $\pm$ 81.5	0.721 $\pm$ 0.04
	GAUSSIAN-NOISE	10.87 $\pm$ 2.25	0.614 $\pm$ 0.30	296.5 $\pm$ 73.3	0.624 $\pm$ 0.03
	GAN	11.02 $\pm$ 1.82	0.600 $\pm$ 0.29	301.4 $\pm$ 92.4	0.613 $\pm$ 0.02
	Gong et al. [25]++	10.84 $\pm$ 1.94	0.556 $\pm$ 0.28	301.0 $\pm$ 81.4	0.658 $\pm$ 0.02
	Titcombe et al. [75]	10.76 $\pm$ 2.37	0.557 $\pm$ 0.24	345.5 $\pm$ 86.1	0.643 $\pm$ 0.01
	Gong et al. [25]	10.91 $\pm$ 1.88	0.560 $\pm$ 0.29	304.5 $\pm$ 82.5	0.616 $\pm$ 0.01
	Peng et al. [62]	10.17 $\pm$ 2.32	0.491 $\pm$ 0.24	399.1 $\pm$ 55.3	0.667 $\pm$ 0.04
	Hayes et al. [28]	10.16 $\pm$ 1.95	0.535 $\pm$ 0.25	320.8 $\pm$ 79.0	0.601 $\pm$ 0.02
	Wang et al. [82]	10.39 $\pm$ 2.55	0.505 $\pm$ 0.24	341.7 $\pm$ 74.2	0.669 $\pm$ 0.05
	SPARSE-STD	9.78 $\pm$ 2.13	0.485 $\pm$ 0.24	367.3 $\pm$ 44.7	0.663 $\pm$ 0.03
	<b>SCA0.1</b>	<b>9.56 <math>\pm</math> 2.30</b>	<b>0.454 <math>\pm</math> 0.25</b>	<b>396.6 <math>\pm</math> 45.0</b>	<b>0.659 <math>\pm</math> 0.04</b>
	<b>SCA0.25</b>	<b>9.27 <math>\pm</math> 2.06</b>	<b>0.452 <math>\pm</math> 0.25</b>	<b>412.8 <math>\pm</math> 53.7</b>	<b>0.661 <math>\pm</math> 0.05</b>
	<b>SCA0.5</b>	<b>9.12 <math>\pm</math> 2.68</b>	<b>0.368 <math>\pm</math> 0.24</b>	<b>421.7 <math>\pm</math> 49.9</b>	<b>0.653 <math>\pm</math> 0.04</b>

**Table 14: Stability analysis 2:** Performance comparison (mean $\pm$  standard deviations) across multiple runs in *end-to-end* network setting (*lower rows=better defense*) on Medical MNIST dataset.

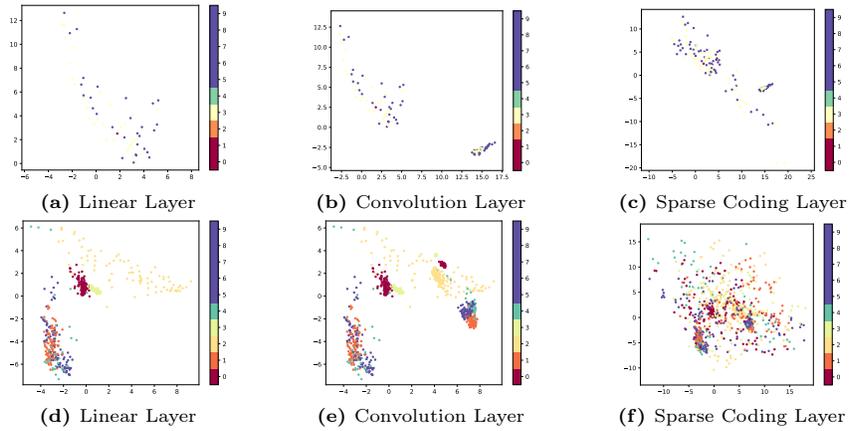
Dataset	Defense	PSNR $\downarrow\downarrow$	SSIM $\downarrow\downarrow$	FID $\uparrow\uparrow$	Accuracy
Medical	NO-DEFENSE	30.17 $\pm$ 0.90	0.912 $\pm$ 0.01	12.40 $\pm$ 8.69	0.998 $\pm$ 0.01
MNIST	GAUSSIAN-NOISE	27.00 $\pm$ 1.30	0.828 $\pm$ 0.05	17.29 $\pm$ 11.9	0.886 $\pm$ 0.06
	GAN	25.05 $\pm$ 2.78	0.699 $\pm$ 0.03	29.08 $\pm$ 20.5	0.995 $\pm$ 0.01
	Gong et al. [25]++	20.37 $\pm$ 1.65	0.451 $\pm$ 0.03	44.68 $\pm$ 30.9	0.871 $\pm$ 0.01
	Titcombe et al. [75]	20.51 $\pm$ 0.28	0.574 $\pm$ 0.01	28.23 $\pm$ 1.65	0.805 $\pm$ 0.06
	Gong et al. [25]	23.65 $\pm$ 1.07	0.605 $\pm$ 0.09	37.16 $\pm$ 26.2	0.757 $\pm$ 0.03
	Peng et al. [62]	17.42 $\pm$ 2.87	0.519 $\pm$ 0.22	65.39 $\pm$ 32.8	0.866 $\pm$ 0.08
	Hayes et al. [28]	19.57 $\pm$ 1.08	0.003 $\pm$ 0.01	155.0 $\pm$ 92.5	0.847 $\pm$ 0.08
	Wang et al. [82]	17.89 $\pm$ 2.09	0.463 $\pm$ 0.08	101.8 $\pm$ 66.5	0.829 $\pm$ 0.08
	SPARSE-STD	13.49 $\pm$ 0.29	0.158 $\pm$ 0.09	203.4 $\pm$ 92.2	0.865 $\pm$ 0.05
	<b>SCA0.1</b>	<b>12.46 <math>\pm</math> 0.30</b>	<b>0.006 <math>\pm</math> 0.01</b>	<b>231.8 <math>\pm</math> 124</b>	<b>0.858 <math>\pm</math> 0.08</b>
	<b>SCA0.25</b>	<b>11.89 <math>\pm</math> 0.35</b>	<b>0.008 <math>\pm</math> 0.01</b>	<b>254.1 <math>\pm</math> 153</b>	<b>0.850 <math>\pm</math> 0.08</b>
	<b>SCA0.5</b>	<b>11.19 <math>\pm</math> 0.11</b>	<b>0.001 <math>\pm</math> 0.01</b>	<b>276.9 <math>\pm</math> 97.0</b>	<b>0.841 <math>\pm</math> 0.08</b>

**Table 15:** Comparison of Accuracy Scores among our unoptimized SCA0.1 and TUNED SCA0.1 (kernel: 5 $\rightarrow$  7) in all 3 setups on CelebA and Medical MNIST datasets.

Dataset	Setup	SCA0.1 $\uparrow\uparrow$	TUNED SCA0.1 $\uparrow\uparrow$
CelebA	PLUG AND PLAY	0.726	<b>0.730</b>
	END TO END	0.748	<b>0.751</b>
	SPLIT	0.745	<b>0.759</b>
Medical	PLUG AND PLAY	0.888	<b>0.899</b>
MNIST	END TO END	0.888	<b>0.996</b>
	SPLIT	0.946	<b>0.967</b>

Table 16

<i>Model</i>	TIME (SEC)
NO-DEFENSE	10555.3
GAUSSIAN-NOISE	12555.3
GAN	15762.4
Titcombe et al. [75]	14390.2
Gong et al. [25]	16061.8
Gong et al. [25]++	17521.8
Peng et al. [62]	18921.2
Hayes et al. [28]	16923.9
Wang et al. [82]	15229.9
SPARSE-STANDARD	12327.5
SCA0.1	17009.8
SCA0.25	17181.2
SCA0.5	17912.9



**Fig. 9:** UMap 2D projections of input images' features by class after 2 linear layers, 2 conv. layers, or 2 sparse-coded layers on CelebA (top) & Medical MNIST (bottom).