

Robust Adversarial Defense by Tensor Factorization

Manish Bhattarai
Theoretical Division, LANL
Los Alamos, USA
ceodsppectrum@lanl.gov

Mehmet Cagri Kaymak
Ryan Barron
Theoretical Division, LANL
Los Alamos, USA

Ben Nebgen
Kim Rasmussen
Boian S. Alexandrov
Theoretical Division, LANL
Los Alamos, USA

Abstract—As machine learning techniques become increasingly prevalent in data analysis, the threat of adversarial attacks has surged, necessitating robust defense mechanisms. Among these defenses, methods exploiting low-rank approximations for input data preprocessing and neural network (NN) parameter factorization have shown potential. Our work advances this field further by integrating the tensorization of input data with low-rank decomposition and tensorization of NN parameters to enhance adversarial defense. The proposed approach demonstrates significant defense capabilities, maintaining robust accuracy even when subjected to the strongest known auto-attacks. Evaluations against leading-edge robust performance benchmarks reveal that our results not only hold their ground against the best defensive methods available but also exceed all current defense strategies that rely on tensor factorizations. This study underscores the potential of integrating tensorization and low-rank decomposition as a robust defense against adversarial attacks in machine learning.

Index Terms—adversarial defense, tensor factorizations, tensorial denoising, tensor train, tucker decomposition

I. INTRODUCTION

The recent advances in machine learning and deep learning have empowered a multitude of applications across various domains, from image recognition to recommender systems. However, these models are susceptible to adversarial attacks, wherein small, carefully crafted perturbations to the input can cause them to output incorrect results [1]. This vulnerability represents a significant challenge to the reliable application of machine learning models, particularly in critical areas such as cybersecurity [2] and healthcare [3]. This raises two primary concerns: the impact on the credibility of current machine learning systems and the danger of malevolent adversarial attacks in real-world applications.

Motivated by the above-mentioned challenges, our study seeks to enhance our understanding of adversarial attacks and provide efficient and robust defense mechanisms against them. We are specifically interested in exploring the capabilities of tensor factorization as a means to defend against adversarial incursions in the image domain. We propose effective defenses that can be incorporated without significantly altering the original model structure or performance. We leverage an extensive parameter search for tensor factorization method to counter attacks, with a focus on the preservation of core data features in the process of eliminating adversarial perturbations.

Dataset (Metric, ϵ)	Method	Clean	AA
CIFAR-10 ($l_\infty, \epsilon = 8/255$)	Rank #1	93.25	70.69
	Ours	85.59	70.24
CIFAR-10 ($l_2, \epsilon = 128/255$)	Rank #1	95.54	84.86
	Ours	86.61	77.73
CIFAR-100 ($l_\infty, \epsilon = 8/255$)	Rank #1	75.22	42.67
	Ours	60.12	42.68

TABLE I: Comparison of test accuracy(%) from our tensorial denoiser to the state-of-the-art model, as in RobustBench [4].

II. RELATED WORK

Tensor decompositions as a defense against adversarial attacks on deep learning (DL) models were first introduced in [5] to demonstrate a simple yet effective approach to resist the attacks without significant degradation compared to the model’s original performance on clean data. Another work, “defensive tensorization”, presented in [6], proposed a novel adversarial defense technique that employed a latent high-order factorization of network layers, then applied tensor dropout in the latent subspace to yield dense reconstructed weights. This work demonstrated effective versatility across multitudinous domains, reaching low-precision architectures. Entezari and Papalexakis investigated adversarial attacks on recommendation systems to present defense methods such as low-rank reconstructions as well as a transformation of attacked data. [7]

Block-term Dropout (BT-Dropout) was proposed in [8] whereby a network was factorized into a latent high-order representation, imposing a low-rank block-term tensor structure on the weights of the fully-connected layer. They applied BT-Dropout in the latent subspace without directly pruning the weights. He et al. [9] presented an adversarial defense method based on tensor decomposition (TDNN), which decomposed then reconstructed images to maintain critical features. This process removed adversarial example perturbations, exhibiting improved defense along with lower run times relative to traditional tensor decomposition.

Tensor layers plus tensor dropout implemented in convolutional neural networks (CNNs) as a means to improve inductive bias, robustness, and efficiency by using low-rank tensor structures on the weights of tensor regression layers was introduced in Ref. [10]. This work established superior performance post-model-modification, furthering robustness against noise and adversarial attacks.

Samangouei et al. [11] presented Defense-GAN. This new framework utilized generative models' capabilities to protect deep neural networks against adversarial attacks by modeling the distribution of unperturbed images paired with the removal of adversarial changes during inference. Similarly, Ilyas et. al [12] posits that adversarial examples in machine learning are due to non-robust features expressed as predictive data distribution patterns but fragile for model performance, even unintuitive to humans. These features, identified within a theoretical framework and shown to be prevalent in standard datasets, reflect a disconnect between human-defined concepts of robustness due to inherent data geometry.

Current tensor-based denoising frameworks lack strategies for selecting optimal hyper-parameters, like tensor ranks for optimal decomposition, and their defense mechanisms are tested only against relatively simple attacks such as FGSM and PGD. Our work addresses this gap, inspired by previous research on low-rank approximation tensor decomposition defense strategies. Our approach combines the low-rank approximated tensorized image, which reshapes into a higher-order tensor by stacking image patches, and the reparametrization of the neural networks (NN) with tensors. We leverage both Tucker [13] and Tensor-Train-based [14] decompositions and highlight the efficacy of our optimal rank selection strategy. Furthermore, to facilitate direct comparisons with existing state-of-the-art defense frameworks, we assess our model's effectiveness using AutoAttack [15], a well-recognized standard in the field. This innovative approach enhances the robustness and efficiency of the denoising process, thereby improving the overall system performance.

III. BACKGROUND

This section discusses our defense model's mathematical foundation. Our denoiser model, as illustrated in Figure 1, is comprised of patch extract, patch merge, and tensor factorization.

A. Tensorization of images

1) *Patch Extraction*: Given an input image as $I \in \mathbb{R}^{C \times W \times H}$, where C is the number of channels, W is the width, and H is the height. Similarly, the kernel size (patch size) is represented as K , the stride as S , the padding as P , and the dilation as D . With the above defined, the Image I is tensorized into a 4-D tensor $O \in \mathbb{R}^{\frac{W-K+2P}{S} \times \frac{H-K+2P}{S} \times C \times K \times K}$, where each $K \times K$ patch in the input tensor becomes a column in the output tensor as:

$$\begin{aligned} O_{w,h,c,k_1,k_2} &= I_{c,S \cdot w + D \cdot k_1, S \cdot h + D \cdot k_2} \\ &\forall k_1, k_2 \in [0, K] \\ &\forall w \in [0, \frac{W-K+2P}{S}] \\ &\forall h \in [0, \frac{H-K+2P}{S}] \end{aligned} \quad (1)$$

Here, O_{w,h,c,k_1,k_2} is the pixel value at the position (k_1, k_2) in the patch at position (w, h) for the channel c in the output

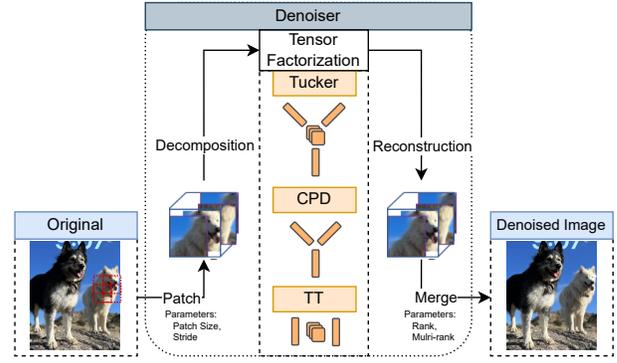


Fig. 1: Overview of the Tensorial denoiser

tensor. $I_{c,S \cdot w + D \cdot k_1, S \cdot h + D \cdot k_2}$ is the corresponding pixel value in the input tensor. s is the stride and d is the dilation.

2) *Patch Merge*: This operation does the inverse of the operation performed in patch extraction, folding the patches back to their original locations to reform the image. Mathematically, this operation can be expressed as follows:

$$\begin{aligned} I_{c,w,h} &= \frac{1}{C_{w,h}} \sum_{k_1=0}^K \sum_{k_2=0}^K O_{\frac{w-D \cdot k_1}{S}, \frac{h-D \cdot k_2}{S}, c, k_1, k_2} \\ &\cdot \mathbb{I}_{S|(w-D \cdot k_1), S|(h-D \cdot k_2)} \\ &\forall w \in [0, W], h \in [0, H], \end{aligned} \quad (2)$$

where

$$\begin{aligned} C_{w,h} &= \sum_{k_1=0}^K \sum_{k_2=0}^K \mathbb{I}_{S|(w-D \cdot k_1), S|(h-D \cdot k_2)} \\ &\forall w \in [0, W], h \in [0, H]. \end{aligned} \quad (3)$$

Here, $I_{c,w,h}$ is the resulting image, O_{w,h,c,k_1,k_2} are the patches to be folded back into the image, and \mathbb{I} is the indicator function that ensures the conditions inside the brackets hold. The function returns 1 if the condition is true and 0 otherwise.

The indicator function $\mathbb{I}_{S|(w-D \cdot k_1), S|(h-D \cdot k_2)}$ ensures the original patch extraction positions in the image are correctly restored. The positions w and h are chosen to be multiples of the stride S and are positioned correctly for the dilation D and the patch indices k_1 and k_2 . This patch and merge strategies are shown in Figure 1.

B. Tensor Decomposition

Once the image is transformed into patches, we perform a low-rank approximation of the tensors with tensor decomposition. We apply Tucker Decomposition [13] and Tensor train decomposition [14] in this work. The details about these methods are presented in the following sections.

1) *Tucker Decomposition*: Given the tensor patch $O \in \mathbb{R}^{M \times N \times C \times K \times K}$, the Tucker decomposition is:

$$O \approx \mathcal{G} \times_1 A^{(1)} \times_2 A^{(2)} \times_3 A^{(3)} \times_4 A^{(4)} \times_5 A^{(5)}, \quad (4)$$

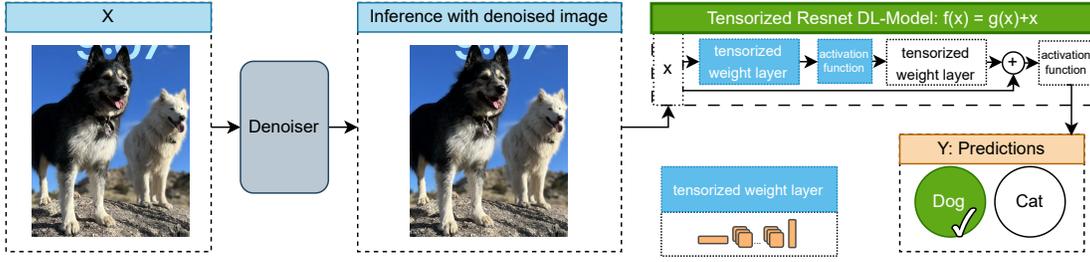


Fig. 2: Overview of the Adversarial denoising setup

where $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times R_3 \times R_4 \times R_5}$ is the core tensor, $A^{(n)} \in \mathbb{R}^{I_n \times R_n}$ are the factor matrices for each mode ($n = 1, 2, 3, 4, 5$), and \times_n denotes the n -mode product.

2) *Tensor Train Decomposition*: Tensor Train decomposition, given the tensor patch $O \in \mathbb{R}^{M \times N \times C \times K \times K}$, is as follows:

$$O(i_1, i_2, i_3, i_4, i_5) \approx \sum_{r_1, r_2, r_3, r_4} G_{i_1, r_1}^{(1)} G_{r_1, i_2, r_2}^{(2)} G_{r_2, i_3, r_3}^{(3)} G_{r_3, i_4, r_4}^{(4)} G_{r_4, i_5}^{(5)} \quad (5)$$

where $G^{(n)}$ are the TT cores for each mode, and the indices r_n (called ranks) represent the connections between the cores [14].

The low-rank, compressed representation of the patched image tensor is then reconstructed using the Patch Merge algorithm. The decompressed image, where high-frequency noise has been removed through tensor factorization, is then classified with the DL model. An overview of utilizing this denoiser model for adversarial denoising is shown in Figure 2.

C. Tensorizing Neural Network

In addition to applying tensor factorizations as pre-processing steps for input images, we also execute a low-rank re-parameterization on the NN layers. While low-rank approximation of NNs was initially designed for significantly reducing the parameters yielding great accelerations [16]–[20], they have also demonstrated promising results for adversarial defense [6], [9]. In an NN with a convolutional kernel parameterized as $\mathbf{S} \in \mathbb{R}^{d \times d \times P \times Q}$, where d is the filter size, and P and Q are the number of input and output channels respectively, the weight tensor can be factorized using the Tucker Decomposition. The factorization is expressed as:

$$\tilde{\mathbf{S}}_{i,j,p,q} = \sum_{r_p=1}^{R_p} \sum_{r_q=1}^{R_q} \mathbb{G}_{i,j,r_p,r_q} \mathbb{A}_{p,r_p}^P \mathbb{A}_{q,r_q}^Q \quad (6)$$

Here, \mathbb{G} is the reduced kernel tensor, R_p and R_q are the ranks of input and output feature map dimensions, respectively, and \mathbb{A}^P and \mathbb{A}^Q are the factor matrices corresponding to the input and output feature maps. This factorization transforms a single convolutional layer into three distinct convolutional layers: two (1×1) layers for \mathbb{A}^P and \mathbb{A}^Q , and a $d \times d$ convolutional layer for \mathbb{G} . Consequently, the complexity of

the original layer is significantly reduced, leading to potential computational and storage efficiency.

Alongside Tucker, we employ tensor-train based reparameterization of the NN. Here, the TT factorization is given as:

$$\tilde{\mathbf{S}}_{i,j,p,q} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} \mathbb{G}_{i,r_1}^1 \mathbb{G}_{r_1,j,r_2}^2 \mathbb{G}_{r_2,p,r_3}^3 \mathbb{G}_{r_3,q}^4 \quad (7)$$

Each \mathbb{G}^k is a TT-core, and R_k are the TT-ranks. This factorization transforms the 4D tensor \mathbf{S} into a sequence of matrices and vectors, which can be stored and manipulated more efficiently. The TT-decomposition is mainly used on high-order (> 3) tensors, providing efficient tensor operations yet more compression than the Tucker decomposition.

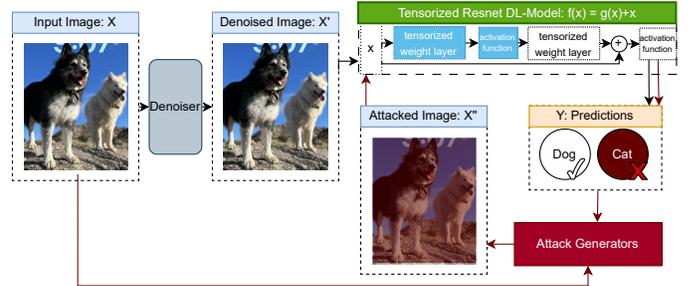


Fig. 3: Overview of the attack generation

D. Adversarial Attacks

The recent advancements in adversarial defense frameworks are commonly evaluated for their robustness using an ensemble of adversarial attacks bundled in a tool called AutoAttack [15]. We, too, utilize this tool to assess the effectiveness of our proposed defense mechanism. AutoAttack integrates white-box and black-box attack methodologies tailored to challenge a model's performance under adversarial conditions. The white-box attack methods include Auto-PGD (Projected Gradient Descent) and APGD-DLR (Auto PGD with Difference of Logits Ratio). In contrast, the black-box attack techniques encompass FAB (Fast Adaptive Boundary attack) and Square Attack.

Auto-PGD and APGD-DLR constitute gradient-based approaches. These methods strategically leverage gradient data to iteratively adjust inputs until a misclassification results.

They utilize constraints by limiting the alterations within a predefined threshold, denoted as ϵ , governed by the attack norms (l_1 or l_∞). On the other hand, FAB and Square Attack are decision-based attacks that generate adversarial examples by making small, ϵ -bounded adjustments to the input. These modifications are fine-tuned to observe whether they lead to misclassification by the model. AutoAttack enables the generation of adversarial samples that probe the performance of our denoising model under severe black-box and white-box adversarial attacks. This helps us determine how our framework fares when subjected to these rigorous tests. The overview of the attack generation process is shown in Figure 3. Initially, the adversarial attack is crafted for one DL model, such as ResNet or WideResNet. Once generated, these adversarial samples are processed through the denoiser before passage to the DL model.

IV. METHOD

Figure 4 offers a comprehensive overview of our training and inference pipeline. We use an eight fold cross-validation approach to optimally select the parameters for the denoising block, such as *patch size*, *stride*, and *tensor decomposition ranks*. The training data is divided into eight partitions. Each training fold consists of seven partitions used for training, leaving one partition for validation. This process results in eight distinct training partitions, thus yielding eight trained DL models. While we utilized CIFAR10 and CIFAR100 datasets for the above experiments, we employed Resnet18 and WideResnet28-10 DL models for adversarial robustness evaluation. We utilized the PyTorch Lightning framework with appropriate seeding to ensure reproducibility across experiments. We set the learning rate of $1e-2$, batch size of 256, and maximum training epoch of 200. The experiments were run on NVIDIA A100 GPUs of 80GB memory. For hyperparameter tuning, we exploited Ray Tune [21] framework for handling parallel workloads.

Autoattack is subsequently employed to create attacks corresponding to the validation sets for each model. For a one-to-one comparison with state-of-the-art adversarial robustness framework, we exploited L_1 and L_∞ attacks with perturbation ϵ equals to $8/255$ and $128/255$ respectively. Almost all the datasets subjected to Autoattack resulted in 0 evaluation accuracy for the DL model.

The generated attacks are then fed into the denoiser block. Here, the block is evaluated using different combinations of hyperparameters such as *patch size*, *stride*, and *tensor decomposition ranks*. These hyperparameters were selected using an optimization technique known as "Tree-Structured Parzen Estimator" [22] specifically implemented through Optuna [23], a Python library for hyperparameter optimization. Specifically, we have four hyperparameters to tune: *patch_size*, *stride*, *rank_p*, and *rank_k*. *patch_size* is the size of the patches into which the image is divided, and *stride* determines the overlap between these patches. These parameters are selected from categorical lists of potential values, specifically [4, 8, 16, 24] for patch size and [1, 2, 4] for stride. Although a

patch is 2D, the same rank (*rank_p*) is chosen to control the tensor decomposition's rank along both patch axes to simplify the search space for tucker decomposition. Lastly, parameter *rank_k* for tucker decomposition controls the decomposition's rank along the number of patches. The rank corresponding to the channel dimension is kept as it is. Our approach is defined-by-run, meaning the hyperparameter search space is dynamically constructed during the optimization process since the feasible ranges for *rank_k* and *rank_p* are determined by the patch size and stride parameters (i.e., the decomposition rank cannot be larger than the size of the corresponding tensor dimension). A step size of 4 is used for the rank hyperparameters to prune the search space further. The ranks' maximum and minimum allowed values are carefully calculated based on the patch size and stride, ensuring an optimal balance between computational efficiency and performance. For the TT decomposition, the hyperparameters *rank_k* and *rank_p* have a different interpretation. They determine the sizes of the "connecting" dimensions in the train of tensors (TT-cores), controlling the complexity of the multilinear relationships between the tensor dimensions. Building upon Equation 1, let's consider a tensorized image, O , in its 5D form. The first two dimensions of this tensor are contracted to form a simpler 4D tensor, denoted as $O \in \mathbb{R}^{\frac{(W-K+2P)(H-K+2P)}{S^2} \times C \times K \times K}$. For the Tucker decomposition process, the multi-ranks are arranged in the configuration [*rank_k*, 3, *rank_p*, *rank_p*]. In contrast, the TT ranks were configured as [1, *rank_k*, *rank_p*, 3, 1] for the reshaped tensor $O \in \mathbb{R}^{\frac{(W-K+2P)(H-K+2P)}{S^2} \times K \times K \times C}$. This reshaping of the tensor is required for the appropriate selection of TT ranks.

The specific configuration of the DL model and the corresponding attacks on the validation dataset determine the optimal parameters for the denoiser. The denoiser is evaluated based on a fitness score metric, the average clean and adversarial accuracy. Then, the optimal denoiser is selected based on the parameter configuration that yields the highest fitness score across the independent folds. For the DL model, we adopt the model that is trained on the entire training set. Early stopping based on validation loss is employed to find the optimal model. Once we determine the optimal denoiser and DL model, an attack is launched on the test set and evaluated using the denoiser.

Moreover, we also incorporate low-rank approximation of the neural network weights to accelerate the DL model and provide additional robustness across adversarial attacks. We exploited Tucker/TT based tensor decomposition to compress the weight matrices of NN without significant loss of useful information. The reduced complexity of the model aided in quicker processing and efficient memory usage, accelerating the model's performance. Simultaneously, low-rank approximation also helped in minimizing the impact of adversarial perturbations as it aids in capturing the dominant, most important features of the data while potentially discarding minute, less meaningful perturbations that are commonly used in adversarial attacks. The tensor decomposition was performed on the final DL model, where the ranks for decomposition

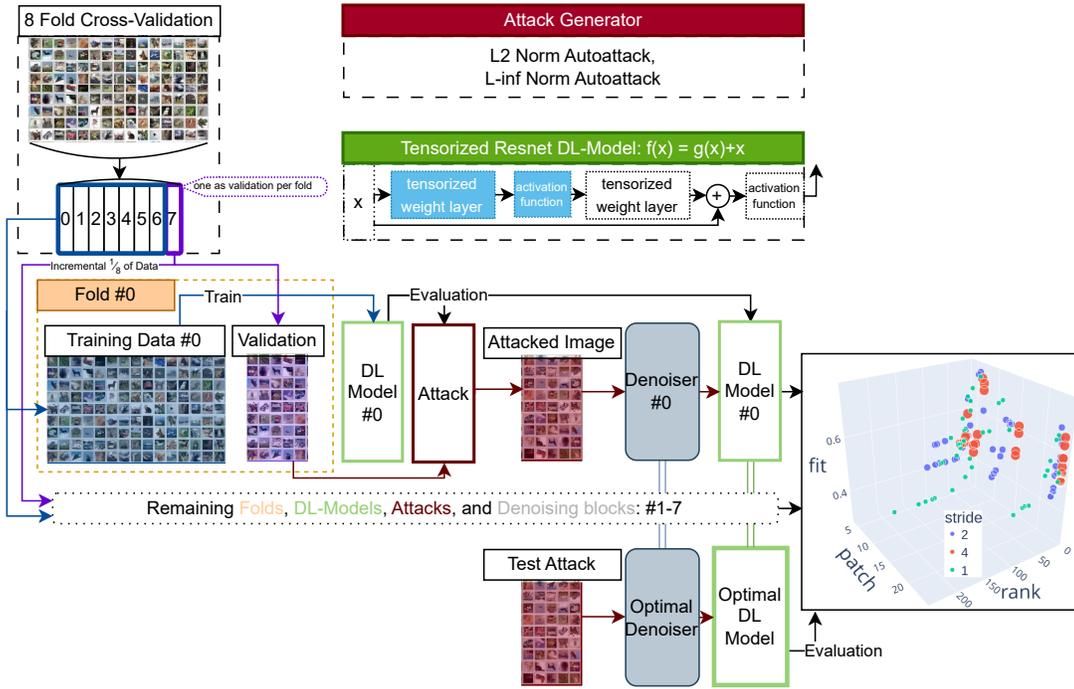


Fig. 4: Overview of comprehensive training and inference pipeline encapsulating the process of optimal parameter selection through eight-fold cross-validation, adversarial attack generation using Autoattack, the employment of the denoiser block equipped with optimized hyperparameters, and the final tensor decomposition for computational efficiency and robustness.

were achieved using the Bayesian approach presented in [24].

In summary, our comprehensive framework offers a robust and efficient solution for DL models facing adversarial attacks. The solution is achieved through cross-validation to select optimal parameters, a denoiser block for robustness against adversarial attacks, and low-rank approximations for computational efficiency. This approach sets a solid foundation for future studies to enhance the resilience and performance of DL models in adversarial scenarios.

V. RESULTS

Figure 5 presents the comprehensive performance of our cross-validated model when evaluated against the test dataset. The first 10 denoiser hyperparameter setups, which deliver the maximum test performance, are specifically reported. Both clean and adversarial accuracy scores are provided, following the procedure of passing the images through the denoiser and feeding them into the Deep Learning (DL) models.

The clean accuracy results demonstrate remarkable consistency, as indicated by the minimal variance in the accuracy plot. In contrast, the adversarial accuracy shows considerable variability. However, relative high-performance levels are maintained, indicating our denoising model’s effective and consistent ability to eliminate adversarial perturbations. Detailed performance scores can be found in tables II, III, IV, and V, where we provide a summary of the hyperparameters leading to the optimal fitness score. Tables II and III provide the results of our denoising model when applied to

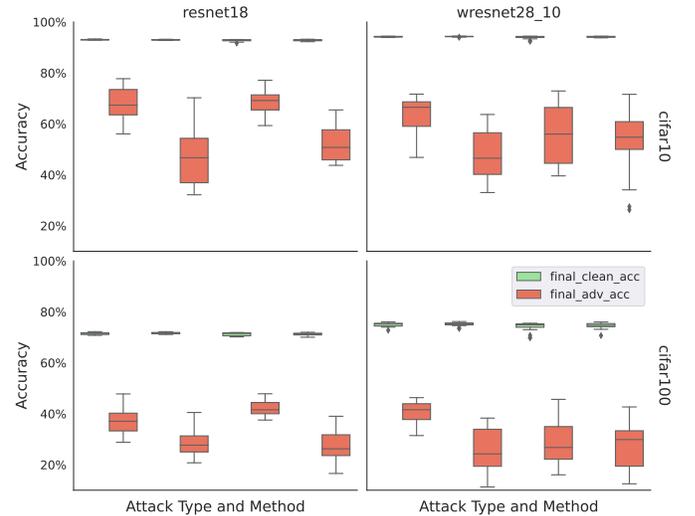


Fig. 5: Distribution of clean and adversarial accuracy scores achieved for top 10 denoiser hyperparameter configurations for test dataset.

the CIFAR10 and CIFAR100 datasets, respectively, using the WideResnet18-10 model. Meanwhile, tables IV and V present the denoising model’s performance on the same datasets but using the Resnet18 model instead. In all scenarios, the models and datasets were evaluated using the AutoAttack method with an L_∞ norm of 8/255 and an L_2 norm of 128/255. We present

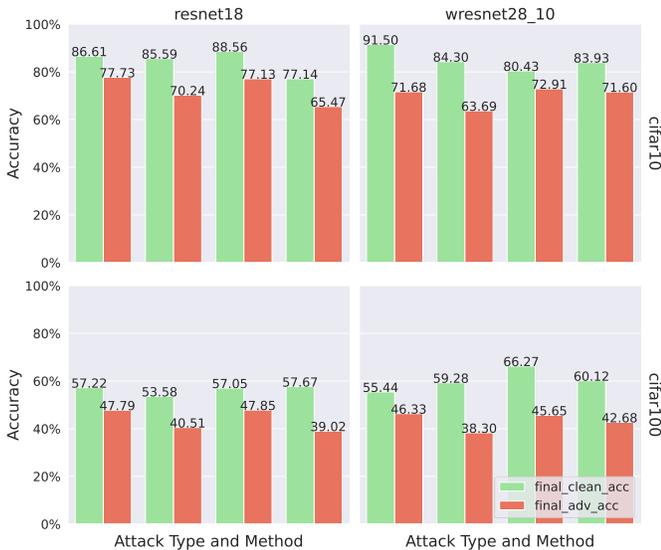


Fig. 6: Statistical representation of adversarial metrics corresponding to the optimal hyperparameter configuration that maximizes the average of clean and adversarial accuracy.

the denoising model’s performance under these configurations (DL model, datasets, attack norm), employing Tensor Train and Tucker tensor decomposition methods.

The results obtained from various model configurations highlight the influence of different adversarial attack scenarios. For instance, when the WideResnet architecture is employed on the CIFAR10 dataset, with an L_∞ attack and Tucker decomposition (as detailed in table IVa), the configuration yielding the highest adversarial accuracy (0.7160) involved a patch size of 24, a stride of 1, a rank_k of 60, and a rank_p of 12. Conversely, when the Resnet18 architecture was used with the CIFAR100 dataset, subjected to an L_2 attack and TT decomposition (refer table Vd), the top-performing model configuration (achieving an adversarial accuracy of 0.4633) had a patch size of 8, a stride of 2, a rank_k of 32, and a rank_p of 12. These variances in performance across different model configurations underscore the diverse impacts of adversarial attacks. Furthermore, in Figure 6, we present the robustness metric, highlighting the hyperparameters corresponding to the highest clean and adversarial accuracy average.

We have also compared our findings with the current best-in-class adversarial robust models, using the Robustbench benchmark [4]. This comparison is depicted in Figure 7, and Table I presents a detailed contrast of our top-performing denoising model against the state-of-the-art adversarial robustness model that is based on denoising diffusion [25].

For the CIFAR-10 dataset, the leading adversarial robust framework achieves a top accuracy of 93.25% on clean data and 70.69% on an AutoAttack dataset with L_∞ norm. In contrast, our top model performs comparably with 85.59% and 70.24% accuracy, respectively. Unlike most existing denoising models presented, our approach does not require additional datasets for training or any form of adversarial training. This

TABLE II: Top 5 Final Adversarial Accuracies for cifar10 and Wide-Resnet

(a) Linf Attack, Tucker Decomposition

patch	stride	rank_k	rank_p	clean acc	adv acc
24	1	60	12	0.8393	0.7160
8	2	24	20	0.8751	0.7081
8	4	20	8	0.8029	0.6941
8	1	36	8	0.9075	0.6164
24	1	36	20	0.7159	0.6148

(b) Linf Attack, TT Decomposition

patch	stride	rank_k	rank_p	clean acc	adv acc
4	1	12	8	0.8430	0.6369
8	1	44	8	0.9123	0.6092
8	1	40	16	0.9125	0.5815
24	1	44	16	0.8533	0.5722
16	1	64	20	0.7609	0.5681

(c) L2 Attack, Tucker Decomposition

patch	stride	rank_k	rank_p	clean acc	adv acc
8	1	52	8	0.8043	0.7291
8	2	28	8	0.9150	0.7134
4	2	12	12	0.7921	0.7118
8	1	32	12	0.8640	0.7078
8	1	40	20	0.8951	0.6791

(d) L2 Attack, TT Decomposition

patch	stride	rank_k	rank_p	clean acc	adv acc
8	1	48	8	0.8215	0.7182
8	2	32	16	0.8893	0.7026
4	1	12	12	0.7820	0.6923
24	1	40	12	0.8367	0.6801
8	1	36	16	0.9137	0.6674

may be one of the reasons why the clean accuracy of our approach trails behind the existing defense frameworks. Notably, our model surpasses the top-ranked model on the CIFAR-100 dataset, outperforming it by a margin of 0.01% on the attacked dataset for the L_∞ norm. However, our model does not perform as well under an L_2 AutoAttack, trailing behind the best model by a margin of approximately 7%. This shortfall could be attributed to utilizing the Tensor Train (TT) and Tucker tensor decomposition tools with an L_2 minimization objective, which may inadvertently allow for adversarial noise during reconstruction. To address this issue, leveraging the Tensor decomposition tool based on L_1 minimization [26], [27] may improve resistance against adversarial noise.

Despite its simplicity, our denoising model demonstrates competitive performance compared to the best existing models. The most significant advantage of our model over others lies in its real-time denoising capabilities.

VI. CONCLUSION

This study has presented a comprehensive evaluation of a Tensor factorization based denoising model and its impact on the performance of deep learning models under various adversarial attack scenarios. Tensor Train and Tucker tensor

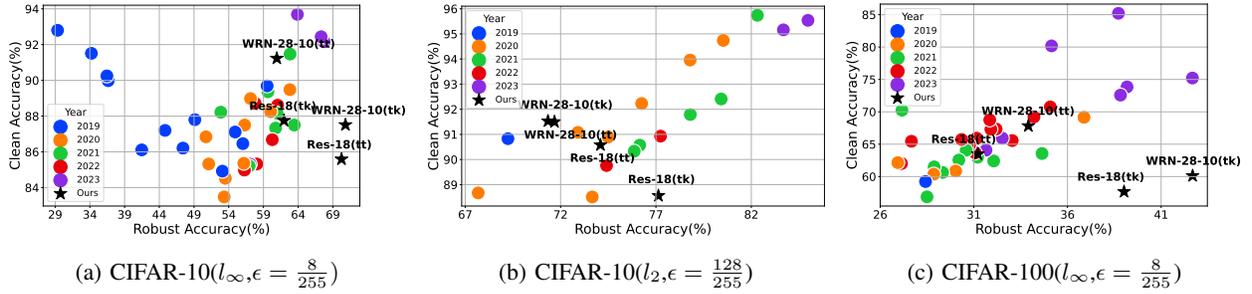


Fig. 7: Robust accuracy against AutoAttack and clean accuracy of top-rank models from RobustBench [4]. Different colors are used to indicate top-rank models for different publication years. Our denoising model was based on Tucker(tk) and TensorTrain(tt) for Resnet-18(Res-18) and WideResnet-70-16 (WRN-70-16) architectures.

TABLE III: Top 5 Final Adversarial Accuracies for cifar100 and Wide-Resnet

(a) Linf Attack, Tucker Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	1	40	8	0.6012	0.4268
8	1	32	12	0.4787	0.3444
8	2	36	8	0.4111	0.3430
24	1	44	20	0.6322	0.3429
4	2	12	8	0.6650	0.3422
(b) Linf Attack, TT Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	4	28	8	0.5928	0.3830
8	1	40	12	0.6131	0.3510
8	1	36	12	0.6661	0.3445
8	2	48	8	0.6650	0.3432
4	1	12	12	0.6690	0.3396
(c) L2 Attack, Tucker Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	2	32	12	0.6627	0.4565
8	1	40	20	0.5770	0.4217
8	1	52	8	0.6761	0.4176
8	1	32	16	0.4953	0.3644
8	1	40	12	0.7079	0.3551
(d) L2 Attack, TT Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	2	32	12	0.5544	0.4633
8	1	40	12	0.5521	0.4600
8	4	32	8	0.6630	0.4563
4	1	12	8	0.5594	0.4478
8	1	44	20	0.6690	0.4470

TABLE IV: Top 5 Final Adversarial Accuracies for cifar10 and Resnet18

(a) Linf Attack, Tucker Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	2	16	8	0.7714	0.6547
8	1	36	8	0.8775	0.6195
8	1	20	8	0.7724	0.6060
8	2	28	20	0.6583	0.5935
8	1	28	8	0.8122	0.5782
(b) Linf Attack, TT Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	4	16	8	0.8559	0.7024
8	1	24	8	0.7960	0.6788
24	2	24	12	0.8572	0.6147
8	2	20	24	0.7779	0.5818
8	4	12	16	0.8100	0.5739
(c) L2 Attack, Tucker Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	1	28	16	0.8856	0.7713
16	1	72	20	0.8899	0.7469
16	1	68	20	0.8756	0.7212
16	1	80	16	0.8479	0.7184
8	2	32	24	0.8798	0.7177
(d) L2 Attack, TT Decomposition					
patch	stride	rank _k	rank _p	clean acc	adv acc
8	1	36	20	0.8661	0.7773
8	2	32	24	0.8408	0.7633
8	1	44	8	0.9058	0.7409
8	2	32	12	0.8552	0.7398
8	4	16	12	0.8087	0.7362

decompositions have demonstrated noteworthy results with different configurations of deep learning models, datasets, and adversarial attack norms. Our analysis, encompassing several robustness metrics, has revealed the efficacy of our denoising model under both clean and adversarial conditions. Notably, we have achieved consistent results under clean conditions while effectively mitigating adversarial perturbations, as evidenced by the considerable variability in adversarial accuracy.

The benchmarking of our denoising model against state-of-the-art adversarial robust models revealed that our model’s performance is competitive, even exceeding the performance of top-ranked models in some scenarios. In particular, our model outperformed the leading model on the CIFAR-100 dataset under L_∞ norm AutoAttack. Our findings indicate the potential of the proposed denoising model to significantly enhance the robustness of deep learning models against adver-

TABLE V: Top 5 Final Adversarial Accuracies for cifar100 and Resnet18

(a) Linf Attack, Tucker Decomposition

patch	stride	rank _k	rank _p	clean acc	adv acc
24	1	32	16	0.5767	0.3902
4	1	12	8	0.4781	0.3498
24	1	56	16	0.4526	0.3399
8	4	16	8	0.6160	0.3250
8	4	16	12	0.5181	0.3232

(b) Linf Attack, TT Decomposition

patch	stride	rank _k	rank _p	clean acc	adv acc
8	2	32	12	0.5358	0.4051
116	1	68	24	0.5785	0.3617
8	1	40	24	0.5106	0.3554
24	2	20	12	0.5636	0.3169
24	2	24	12	0.4393	0.3158

(c) L2 Attack, Tucker Decomposition

patch	stride	rank _k	rank _p	clean acc	adv acc
8	1	32	8	0.5705	0.4785
24	1	48	16	0.6255	0.4769
8	1	36	16	0.6132	0.4629
8	2	24	24	0.6028	0.4502
4	1	12	8	0.6015	0.4444

(d) L2 Attack, TT Decomposition

patch	stride	rank _k	rank _p	clean acc	adv acc
8	1	40	24	0.5722	0.4779
4	2	12	12	0.4992	0.4373
16	1	64	20	0.6088	0.4310
8	1	32	24	0.5572	0.4274
8	4	24	8	0.6600	0.4089

sarial attacks.

VII. ACKNOWLEDGMENT

This manuscript has been assigned LA-UR-23-27984. This research was funded by the Los Alamos National Laboratory (LANL) Laboratory Directed Research and Development (LDRD) program under grant 20230287ER and supported by LANL’s Institutional Computing Program, and by the U.S. Department of Energy National Nuclear Security Administration under Contract No. 89233218CNA000001.

REFERENCES

- [1] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014. **1**
- [2] I. Rosenberg, A. Shabtai, Y. Elovici, and L. Rokach, “Adversarial machine learning attacks and defense methods in the cyber security domain,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 5, pp. 1–36, 2021. **1**
- [3] S. G. Finlayson, J. D. Bowers, J. Ito, J. L. Zittrain, A. L. Beam, and I. S. Kohane, “Adversarial attacks on medical machine learning,” *Science*, vol. 363, no. 6433, pp. 1287–1289, 2019. **1**
- [4] F. Croce, M. Andriushchenko, V. Schwag, E. Debenedetti, N. Flammarion, M. Chiang, P. Mittal, and M. Hein, “Robustbench: a standardized adversarial robustness benchmark,” in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. [Online]. Available: <https://openreview.net/forum?id=SSKZPJC7B> **1, 6, 7**
- [5] S. Cho, T. J. Jun, M. Kang, and D. Kim, “Applying tensor decomposition to image for robustness against adversarial attack,” *arXiv preprint arXiv:2002.12913*, 2020. **1**
- [6] A. Bulat, J. Kossaifi, S. Bhattacharya, Y. Panagakis, T. Hospedales, G. Tzimiropoulos, N. D. Lane, and M. Pantic, “Defensive tensorization,” *arXiv preprint arXiv:2110.13859*, 2021. **1, 3**
- [7] N. Entezari and E. E. Papalexakis, “Low-rank defenses against adversarial attacks in recommender systems,” in *2022 IEEE International Conference on Big Data (Big Data)*. IEEE, 2022, pp. 5708–5714. **1**
- [8] C. Ouyang, W. Yang, and R. Hu, “Block-term dropout for robust adversarial defense,” in *2022 IEEE 34th International Conference on Tools with Artificial Intelligence (ICTAI)*. IEEE, 2022, pp. 622–629. **1**
- [9] W. He, B. Song, R. Wang, W. Peng, S. He, and W. Zhou, “Tdn: A tensor decomposition adversarial defense method based on neural network,” in *2021 5th Asian Conference on Artificial Intelligence Technology (ACAIT)*. IEEE, 2021, pp. 157–166. **1, 3**
- [10] A. Kolbeinsson, J. Kossaifi, Y. Panagakis, A. Bulat, A. Anandkumar, I. Tzoulaki, and P. M. Matthews, “Tensor dropout for robust learning,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 15, no. 3, pp. 630–640, 2021. **1**
- [11] P. Samangouei, M. Kabkab, and R. Chellappa, “Defense-gan: Protecting classifiers against adversarial attacks using generative models,” *arXiv preprint arXiv:1805.06605*, 2018. **2**
- [12] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, “Adversarial examples are not bugs, they are features,” *Advances in neural information processing systems*, vol. 32, 2019. **2**
- [13] T. G. Kolda and B. W. Bader, “Tensor decompositions and applications,” *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009. **2**
- [14] I. V. Oseledets, “Tensor-train decomposition,” *SIAM Journal on Scientific Computing*, vol. 33, no. 5, pp. 2295–2317, 2011. **2, 3**
- [15] F. Croce and M. Hein, “Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks,” in *ICML*, 2020. **2, 3**
- [16] D. Lee, S. J. Kwon, B. Kim, and G.-Y. Wei, “Learning low-rank approximation for cnns,” *arXiv preprint arXiv:1905.10145*, 2019. **3**
- [17] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin, “Compression of deep convolutional neural networks for fast and low power mobile applications,” *arXiv preprint arXiv:1511.06530*, 2015. **3**
- [18] M. Astrid and S.-I. Lee, “Cp-decomposition with tensor power method for convolutional neural networks compression,” in *2017 IEEE International Conference on Big Data and Smart Computing (BigComp)*. IEEE, 2017, pp. 115–118. **3**
- [19] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitky, “Speeding-up convolutional neural networks using fine-tuned cp-decomposition,” *arXiv preprint arXiv:1412.6553*, 2014. **3**
- [20] A. Novikov, D. Podoprikin, A. Osokin, and D. P. Petrov, “Tensorizing neural networks,” *Advances in neural information processing systems*, vol. 28, 2015. **3**
- [21] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, and I. Stoica, “Tune: A research platform for distributed model selection and training,” *arXiv preprint arXiv:1807.05118*, 2018. **4**
- [22] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, “Algorithms for hyperparameter optimization,” *Advances in neural information processing systems*, vol. 24, 2011. **4**
- [23] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2623–2631. **4**
- [24] S. Nakajima, M. Sugiyama, S. D. Babacan, and R. Tomioka, “Global analytic solution of fully-observed variational bayesian matrix factorization,” *The Journal of Machine Learning Research*, vol. 14, no. 1, pp. 1–37, 2013. **5**
- [25] Z. Wang, T. Pang, C. Du, M. Lin, W. Liu, and S. Yan, “Better diffusion models further improve adversarial training,” *arXiv preprint arXiv:2302.04638*, 2023. **6**
- [26] D. G. Chachlakis, A. Prater-Bennette, and P. P. Markopoulos, “L1-norm tucker tensor decomposition,” *IEEE Access*, vol. 7, pp. 178 454–178 465, 2019. **6**
- [27] L. Liu, D. M. Hawkins, S. Ghosh, and S. S. Young, “Robust singular value decomposition analysis of microarray data,” *Proceedings of the National Academy of Sciences*, vol. 100, no. 23, pp. 13 167–13 172, 2003. **6**