

On the Robustness of Kolmogorov-Arnold Networks: An Adversarial Perspective

Anonymous authors

Paper under double-blind review

Abstract

Kolmogorov-Arnold Networks (KANs) have recently emerged as a novel approach to function approximation, demonstrating remarkable potential in various domains. Despite their theoretical promise, the robustness of KANs under adversarial conditions has yet to be thoroughly examined. In this paper we explore the adversarial robustness of KANs, with a particular focus on image classification tasks. We assess the performance of KANs against standard white-box and black-box adversarial attacks, comparing their resilience to that of established neural network architectures. Our experimental evaluation encompasses a variety of standard image classification benchmark datasets and investigates both fully connected and convolutional neural network architectures, of three sizes: small, medium, and large. We conclude that small- and medium-sized KANs (either fully connected or convolutional) are not consistently more robust than their standard counterparts, but that large-sized KANs are, by and large, more robust. This comprehensive evaluation of KANs in adversarial scenarios offers the first in-depth analysis of KAN security, laying the groundwork for future research in this emerging field.

1 Introduction

In the rapidly evolving field of deep learning, the robustness of neural networks against adversarial attacks has emerged as a cornerstone of research, driven by the ubiquity of their application, often in sensitive areas. As a result, a multitude of methods have been developed to detect and mitigate adversarial attacks, underscoring the critical need for robust defenses in real-world scenarios. These methods span several areas, including adversarial training (Madry et al., 2017; Bai et al., 2021; Shafahi et al., 2019; Andriushchenko & Flammarion, 2020), defensive distillation (Papernot et al., 2016; Papernot & McDaniel, 2016; Carlini & Wagner, 2016), feature squeezing (Xu et al., 2017), input transformations (Guo et al., 2017; Harder et al., 2021), and using auxiliary models for detection purposes (Zheng & Hong, 2018; Pinhasov et al., 2024; Lapid et al., 2024a). The proliferation of these techniques reflects the ongoing efforts within the research community to fortify neural networks against increasingly sophisticated attacks.

Fully connected neural networks (FCNNs) and convolutional neural networks (CNNs)—foundational models in deep learning—have been extensively studied for their vulnerability and defense mechanisms against adversarial attacks (Dey et al., 2017; Wang et al., 2020; Mohammed et al., 2020; Kerlirzin & Vallet, 1993; Kalina et al., 2022). Recently, the introduction of Kolmogorov-Arnold Networks (KANs) has opened new avenues in function approximation, with theoretical promises of enhanced performance and efficiency (Liu et al., 2024; Bozorgasl & Chen, 2024; Genet & Inzirillo, 2024; Vaca-Rubio et al., 2024; Kiamari et al., 2024).

Many studies have compared the performance of KANs across a range of tasks. For example, Yu et al. (2024) found that FCNNs generally perform better in areas such as machine learning and computer vision. Conversely, KANs excel in representing symbolic formulas.

Another interesting comparison was done by Zeng et al. (2024), focusing on the performance of KANs and FCNNs in handling irregular or noisy functions. They found that KANs exhibited superior performance both for regular and irregular functions. However, in cases involving functions with jump discontinuities or singularities, FCNNs demonstrated greater efficacy.

Dong et al. (2024) studied KANs for time-series classification. They concluded that KANs exhibit significant robustness advantages, attributed to their lower Lipschitz constants. Note that they focused on time-series problems, while we focus herein on image-related tasks; further, they deployed a single attack algorithm (Projected Gradient Descent, PGD; Madry et al. (2017)), whereas we will deploy seven attack algorithms.

This paper delves into the robustness of KANs when faced with adversarial attacks under white-box and black-box conditions. Using the MNIST, KMNIST, FashionMNIST, CIFAR-10, and SVHN datasets, our analysis not only sheds light on the robustness of KANs but also contributes to the broader understanding of neural network security in adversarial settings.

The next section provides a comprehensive overview of the architectures employed in this study, accompanied by a detailed discussion of their robustness. Section 3 outlines the experimental setup used to evaluate our proposed approach. Section 4 examines the robustness of the different architectures against adversarial attacks. Finally, we present concluding remarks in Section 5.

2 Background

2.1 Model Architectures

Our study evaluates two primary types of layers used in artificial neural networks: fully connected layers and convolutional layers. Fully connected layers are the foundational building blocks of feedforward neural networks, where each neuron is connected to every neuron in the next layer. FCNNs are classical feedforward networks that learn floating-point weights for these connections (Cybenko, 1989). In contrast, FCKANs, inspired by the Kolmogorov-Arnold theorem (Kolmogorov, 1961), decompose functions into simpler univariate components, learning these components instead of conventional weights, as detailed by Liu et al. (2024).

Convolutional layers extend the capabilities of fully connected layers by leveraging spatial hierarchies in data such as images. These layers are the first and main layers introduced in CNNs (LeCun et al., 1998), where they typically learn floating-point weights for their convolutional filters. CKANs (Drokin, 2024) apply their unique decomposition principle to learn univariate functions for the filters rather than conventional weights. For more detailed descriptions, please refer to the cited literature.

2.2 Robustness in Neural Networks

Robustness in neural networks refers to their ability to maintain performance when subjected to various forms of perturbations, such as noise, adversarial attacks, or distributional shifts. A robust neural network should not only perform well on the training and validation datasets but also generalize effectively to unseen data, including those with slight modifications or corruptions.

Mathematically, robustness can be defined as follows: Consider a neural network $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^m$, parameterized by θ , and let $\mathbf{x} \in \mathbb{R}^n$ be an input to the network. The network is said to be robust to perturbations if, for a small perturbation, $\delta \in \mathbb{R}^n$, the output of the network remains unchanged, that is:

$$\|f_\theta(\mathbf{x}) - f_\theta(\mathbf{x} + \delta)\| \leq \epsilon,$$

for some small $\epsilon > 0$ and for all $\|\delta\| \leq \eta$, where η represents the maximum allowable perturbation norm.

Adversarial robustness in particular has gained significant attention in recent years. Adversarial attacks involve making small—often imperceptible—perturbations to input data that can cause a neural network to produce incorrect outputs with high confidence (Goodfellow et al., 2014). These attacks exploit the vulnerability of neural networks to specific input manipulations, raising concerns about their reliability in real-world applications (Szegedy et al., 2013).

For example, the Fast Gradient Sign Method (FGSM) demonstrates how slight changes to image pixels can deceive a classifier into misclassification (Goodfellow et al., 2014). Additionally, the Carlini-Wagner (C&W) attack shows that even robust defenses can be circumvented through optimized perturbations (Carlini & Wagner, 2017).

Such examples underscore the critical need for developing networks with enhanced robustness against adversarial attacks. Numerous other attack methodologies further highlight this necessity, revealing the extensive landscape of adversarial vulnerability in neural networks (Goodfellow et al., 2014; Lapid et al., 2024b; Szegedy et al., 2013; Tamam et al., 2023; Wei et al., 2022; Lapid et al., 2022; Chen et al., 2024; Carlini & Wagner, 2017; Lapid & Sipper, 2023b; Andriushchenko et al., 2020; Lapid & Sipper, 2023a).

Several techniques have been proposed to enhance the robustness of neural networks, including adversarial training, regularization methods, and architectural modifications. Adversarial training involves training the network on adversarially perturbed examples, thereby improving its ability to withstand attacks (Madry et al., 2017). Regularization methods—such as dropout (Srivastava et al., 2014) and weight decay—help prevent overfitting and improve generalization. Architectural modifications—such as incorporating attention mechanisms (Vaswani et al., 2017) and leveraging alternative network structures—also contribute to robustness.

The upshot is that robustness is a critical aspect of neural network performance, particularly in the context of adversarial attacks and other perturbations. By understanding and improving the robustness of neural networks, we can develop more reliable and secure models for various applications.

2.3 Adversarial Attack Types

Adversarial attacks are techniques used to evaluate the robustness of neural networks by introducing small perturbations to input data, causing the model to misclassify the input. These attacks are broadly categorized into white-box attacks and black-box attacks, based on the adversary’s level of access to the model.

White-Box Attacks. In white-box attacks, the adversary has full knowledge of the target model, including its architecture, weights, and gradients. This access allows the adversary to compute perturbations directly by exploiting the model’s vulnerabilities. White-box attacks represent the strongest adversarial threat model, as they assume complete transparency of the system.

Black-Box Attacks. In black-box attacks, the adversary has no direct access to the model’s internal parameters or gradients. Instead, the adversary generates adversarial examples by querying the model or leveraging transferability across architectures. Black-box attacks are more practical in real-world scenarios, as they mimic situations where the adversary does not have access to proprietary models or data.

3 Experimental Setup

Do KANs offer enhanced security compared to classical neural network architectures, when confronted with an adversary?

To address this question we conduct an experimental evaluation of KANs. While it may appear unconventional to begin with experiments, direct assessment of the security of KANs is most effectively achieved through targeted attacks and subsequent analyses. We evaluate the robustness of FCKANs, CKANs, FCNNs and CNNs under white-box and black-box adversarial attacks. Both fully connected and convolutional architectures are tested across multiple datasets, and the results are analyzed to compare the robustness of the two architectures.

3.1 Classifier Models

We utilized six models for each architecture (fully connected and convolutional), categorized into small, medium, and large configurations. The configurations for fully connected models are detailed in Table 1, while those for convolutional models are provided in Table 2. All the KAN models were configured with a uniform setting of $num\ knots = 5$ and $spline\ order = 3$.

All the models were trained using the AdamW (Loshchilov & Hutter, 2017) optimizer for 20 epochs. The fully connected models were trained on the MNIST (Deng, 2012), FashionMNIST (Xiao et al., 2017), and KMNIST (Prabhu, 2019) datasets with a learning rate of 1×10^{-4} , weight decay of 5×10^{-4} , and a batch size of 64. The convolutional models were trained on the MNIST, SVHN (Netzer et al., 2011), and CIFAR-10

Table 1: Model configurations for fully connected architectures.

Model	#Params	Hidden Layers
FCKAN _{small}	508,160	[64]
FCNN _{small}	508,810	[640]
FCKAN _{medium}	4,730,880	[256, 1024]
FCNN _{medium}	5,043,210	[1024, 4096]
FCKAN _{large}	35,676,160	[128, 128, 256, 256, 256, 512, 512, 512, 1024, 1024, 1024]
FCNN _{large}	33,678,986	[128, 256, 256, 512, 512, 1024, 1024, 2048, 2048, 4096, 4096]

Table 2: Model configurations for convolutional architectures. All convolutional layers use a kernel size of 3×3 , stride of 1, and padding of 1. Max pooling is applied after the last convolutional layer, followed by two fully connected layers with hidden dimensions [1024, 256].

(a) Grayscale Datasets			(b) Multi-Channel Datasets		
Model	Output Channels	#Params	Model	Output Channels	#Params
CKAN _{small}	[32, 64, 128]	27,056,173	CKAN _{small}	[32, 64, 128]	34,925,677
CNN _{small}	[32, 64, 128]	26,316,810	CNN _{small}	[32, 64, 128]	34,181,706
CKAN _{medium}	[64, 64, 64, 128, 128, 128, 256]	58,554,961	CKAN _{medium}	[64, 64, 64, 128, 128, 128, 256]	74,293,969
CNN _{medium}	[64, 64, 128, 128, 256, 256, 256]	53,648,586	CNN _{medium}	[64, 64, 128, 128, 256, 256, 256]	69,378,378
CKAN _{large}	[64, 64, 64, 128, 128, 128, 256, 512]	120,552,018	CKAN _{large}	[64, 64, 64, 128, 128, 128, 256, 512]	152,019,666
CNN _{large}	[64, 128, 256, 256, 256, 512, 512, 512]	110,744,074	CNN _{large}	[64, 128, 256, 256, 256, 512, 512, 512]	142,202,506

(Krizhevsky, 2009) datasets using a learning rate of 1×10^{-4} , weight decay of 1×10^{-4} , and a batch size of 32.

3.2 Adversarial Attacks

We evaluated the robustness of our models against four distinct white-box attacks. The evaluation began with FGSM (Goodfellow et al., 2014), a foundational attack used as an initial test of model resilience. Subsequently, we assessed the models using more advanced iterative attacks, including PGD (Madry et al., 2017), C&W (Carlini & Wagner, 2017), and the Momentum Iterative Method (MIM) (Dong et al., 2018).

Furthermore, we evaluated the robustness of all the models against three different black-box attack methodologies: transfer-based attacks (Section B.2.1), simple black-box attack (SimBA, (Tu et al., 2019)), square attack (Andriushchenko et al., 2020), and natural evolutionary strategies attack (NES, (Ilyas et al., 2018)). White-box attacks were evaluated on all test sets, while black-box attacks were evaluated on a random sample of 1,000 images from the corresponding test set. Transferability metrics were calculated using 5,000 images from the test sets.

The hyperparameter configurations for these experiments are provided in Table 3. For the grayscale datasets (MNIST, FashionMNIST, and KMNIST), an ϵ value of 32/255 was utilized, whereas for the multi-channel; datasets (CIFAR-10 and SVHN), an ϵ value of 8/255 was employed. All experiments were conducted with a batch size of 16.

Table 3: Hyperparameters used in experiments with white-box and black-box attacks.

Attack Type	Attack	Parameters (Grayscale Datasets)	Parameters (Multi-channel Datasets)
White-box	PGD	$k = 40, \alpha = 0.01$	$k = 30, \alpha = 0.01$
	C&W	$k = 40, \alpha = 0.01, c = 1$	$k = 30, \alpha = 0.01, c = 1$
	MIM	$k = 40, \alpha = 0.01, \mu = 1$	$k = 30, \alpha = 0.01, \mu = 1$
Black-box	Square	$k = 3000$	$k = 2000$
	SimBA	$k = 5000$	$k = 3000$
	NES	$k = 300, \alpha = 0.025, \sigma = 0.0625, n = 40$	$k = 200, \alpha = 0.01, \sigma = 0.01, n = 30$

A comprehensive description of each attack can be found in Appendix B.

4 Results

4.1 Fully Connected Models

In the context of iterative white-box adversarial attack methods (PGD, C&W, MIM), small and medium FCKANs demonstrated lower robustness compared to small and medium FCNNs, as shown in [Table 4](#). Across all examined datasets, small and medium FCKANs were notably vulnerable to these attacks. In contrast, small and medium FCNNs demonstrated a measurable—albeit limited—degree of robustness. Nonetheless, even these models remain highly susceptible to adversarial exploitation, rendering them far from secure in practice.

Table 4: White-box and black-box attacks on fully connected models evaluated across MNIST, FashionMNIST, and KMNIST datasets. Robust accuracy is given for each corresponding attack. The “Acc” column refers to the clean accuracy of the model. Boldface values indicate the most robust model in each experiment.

Dataset	Model	White-Box				Black-Box			Acc
		FGSM	PGD	C&W	MIM	Square	SimBA	NES	
MNIST	FCKAN _{small}	6.96	0.00	0.01	0.00	0.09	0.00	0.00	96.02
	FCNN _{small}	1.40	0.52	0.65	0.62	2.38	0.19	2.08	97.94
	FCKAN _{medium}	9.51	0.00	0.00	0.00	0.00	0.19	0.00	97.62
	FCNN _{medium}	10.80	0.77	0.87	1.05	3.27	18.35	19.94	98.31
	FCKAN _{large}	37.83	1.64	2.27	0.42	24.20	38.78	0.00	94.54
	FCNN _{large}	9.57	4.78	0.15	4.76	2.38	5.75	5.25	97.47
FashionMNIST	FCKAN _{small}	4.84	0.00	0.00	0.00	0.00	0.00	0.00	86.83
	FCNN _{small}	0.95	0.23	0.44	0.32	3.86	0.39	3.67	88.20
	FCKAN _{medium}	3.78	0.00	0.00	0.00	0.00	0.00	0.00	86.72
	FCNN _{medium}	5.23	0.12	0.11	0.18	1.58	11.30	11.80	89.14
	FCKAN _{large}	16.68	0.05	0.09	0.00	19.64	29.06	0.00	83.24
	FCNN _{large}	1.48	0.11	0.00	0.11	0.19	2.48	2.67	88.53
KMNIST	FCKAN _{small}	2.80	0.00	0.00	0.00	0.00	0.00	0.00	81.71
	FCNN _{small}	3.09	1.16	1.52	1.41	7.24	0.89	5.55	90.28
	FCKAN _{medium}	7.65	0.00	0.00	0.00	0.00	1.48	1.68	87.47
	FCNN _{medium}	11.74	1.01	1.11	1.35	6.54	8.82	17.85	91.15
	FCKAN _{large}	25.90	0.36	0.53	0.04	14.78	23.61	0.00	78.92
	FCNN _{large}	7.14	0.70	0.29	1.05	3.86	1.38	2.67	89.28

This trend, in which small and medium FCNNs generally outperform FCKANs in robustness under iterative white-box attacks, and for which FCKANs often failed to classify any adversarial images generated by these methods (except for the small FCKAN on MNIST), begins to shift when considering larger configurations. For large models, the robustness of FCKANs improves significantly, surpassing that of large FCNNs under the C&W attack across all datasets. Not only do large FCKANs outperform large FCNNs under the C&W attack, but they also show substantial improvement over their smaller and medium-sized counterparts. Across all iterative adversarial attacks and datasets, large FCKANs exhibited considerably greater resilience than small and medium FCKANs. In contrast, large FCNNs often demonstrated lower robustness than their smaller variants.

Further evidence of the robustness of large FCKANs under white-box iterative attacks is shown in [Figure 1](#). The figure illustrates that the loss function, which increases as the discrepancy between the predicted and true labels is maximized, remains consistently lower for large FCKANs compared to other models. This indicates that deceiving large FCKANs is more challenging. Additionally, while the losses of other models tend to converge during attacks, the loss of large FCKANs maintains a linear trajectory, further underscoring their resilience.

A particularly notable observation is the superior performance of large FCKANs compared to large FCNNs under the FGSM white-box adversarial attack, with a pronounced gap favoring FCKANs across all datasets. This is further supported by [Figure 2](#), which illustrates the clear advantage of large FCKANs. While the

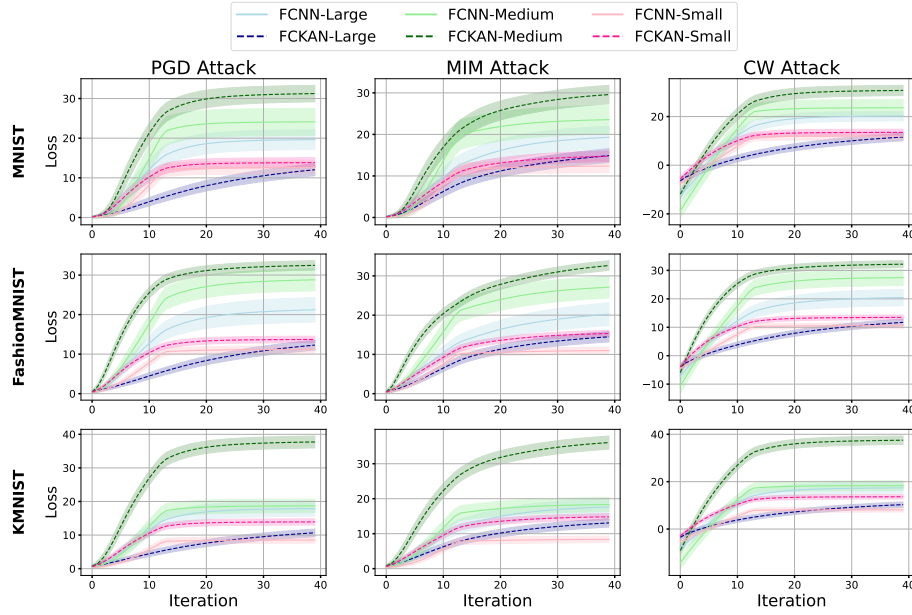


Figure 1: Comparative loss dynamics of FCNNs and FCKANs under PGD, MIM, and C&W attacks on the MNIST, FashionMNIST, and KMNIST datasets (from top to bottom, respectively). Each line represents the mean loss across batches per iteration, with shaded areas indicating the standard deviation.

differences in robustness between small and medium FCNNs and FCKANs are less definitive, the superiority of large FCKANs over large FCNNs is both clear and substantial, firmly establishing FCKANs as the more robust architecture in larger configurations.

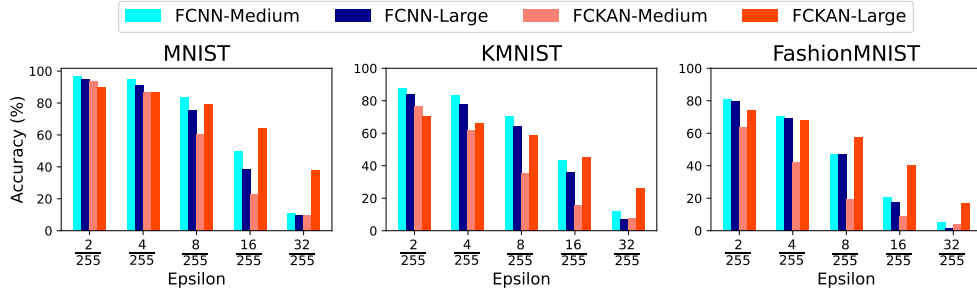


Figure 2: Robust accuracy of $\text{FCNN}_{\text{medium}}$, $\text{FCNN}_{\text{large}}$, $\text{FCKAN}_{\text{medium}}$, and $\text{FCKAN}_{\text{large}}$ as a function of FGSM attack strength, with varying ϵ values. Each bar represents a different model, and robustness is measured as the accuracy of the model against adversarial examples generated with specific ϵ values. The x-axis indicates the ϵ values used for FGSM attacks, while the y-axis shows the corresponding robust accuracy of the models.

When considering black-box attacks (Square, SimBA, and NES), the performance dynamics shift. Small FCNNs outperform FCKANs on all datasets, where small FCKANs failed to classify any image generated with those attacks. When considering the medium models, although FCKANs successfully classified a small number of images (especially on the KMNIST dataset), medium FCNNs still performed better than medium FCKANs across all datasets. However, for larger models, FCKANs demonstrate a marked improvement in robustness against SimBA and Square black-box attacks, significantly surpassing FCNNs. In contrast, similar to the small and medium sizes under NES attack, FCKANs failed to classify any adversarial image. The fact that FCKANs failed under the NES attack could result from the gradient approximation for these architectures being much more accurate than for the FCNN architecture.

In our analysis of the transferability of adversarial white-box attacks (Table 5), we observed that across all datasets, the large FCKAN consistently exhibited the lowest percentage of misclassified images among

all six models. This suggests that fooling the large FCKAN is the most challenging task. Moreover, except for the small and medium models on the KMNIST dataset, FCKANs consistently demonstrated greater resistance to adversarial attacks than their FCNN counterparts for each model size (small, medium and large), highlighting the superior robustness of FCKANs at every scale.

Table 5: Maximum transferability of our fully connected models, for FashionMNIST, KMNIST, and MNIST. Each row represents the model that generates the adversarial example and each column represents the model that evaluates those examples. The value in row i , column j represents the maximum transferability between row- i model and column- j model across four attacks—FGSM, PGD, C&W, MIM—calculated using Equation 11. Unlike robustness tables, the values here indicate the percentage of images misclassified. The ‘Average’ row shows the average transferability for the attacking model. Boldfaced values denote the lowest transferability, where a lower value indicates better robustness.

FashionMNIST						
	FCKAN _{small}	FCKAN _{medium}	FCKAN _{large}	FCNN _{small}	FCNN _{medium}	FCNN _{large}
FCKAN _{small}	100.00	57.32	38.97	59.98	56.23	48.76
FCKAN _{medium}	40.55	100.00	28.11	39.89	39.99	35.81
FCKAN _{large}	32.34	29.91	100.00	31.90	26.03	27.18
FCNN _{small}	81.34	81.57	59.43	99.76	98.46	90.64
FCNN _{medium}	64.82	68.91	45.16	90.19	99.86	85.33
FCNN _{large}	57.13	59.31	39.45	75.18	82.21	100.00
Average (\downarrow)	62.69	66.16	51.85	66.15	67.13	64.61

KMNIST						
	FCKAN _{small}	FCKAN _{medium}	FCKAN _{large}	FCNN _{small}	FCNN _{medium}	FCNN _{large}
FCKAN _{small}	100.00	29.58	21.11	21.60	13.77	19.71
FCKAN _{medium}	32.87	100.00	17.47	16.19	10.94	17.43
FCKAN _{large}	19.59	11.82	99.94	7.99	4.69	11.30
FCNN _{small}	70.68	62.29	45.95	98.72	89.77	71.58
FCNN _{medium}	51.04	42.42	30.76	83.49	98.88	59.40
FCNN _{large}	45.75	35.72	31.89	50.22	43.74	99.66
Average (\downarrow)	53.32	46.97	41.18	46.36	43.46	46.51

MNIST						
	FCKAN _{small}	FCKAN _{medium}	FCKAN _{large}	FCNN _{small}	FCNN _{medium}	FCNN _{large}
FCKAN _{small}	100.00	28.10	22.94	55.10	30.55	32.32
FCKAN _{medium}	35.33	100.00	17.71	38.83	22.67	26.74
FCKAN _{large}	22.78	08.94	99.76	18.50	08.36	14.56
FCNN _{small}	81.35	68.48	48.68	99.78	95.84	88.31
FCNN _{medium}	64.17	54.40	33.83	95.48	99.62	84.39
FCNN _{large}	58.89	44.78	31.77	78.42	72.29	99.92
Average (\downarrow)	60.41	50.78	42.44	64.35	54.88	57.70

Qualitative results corresponding to these attacks are presented in Appendix A.

4.2 Convolutional Models

In the context of iterative white-box adversarial attacks, CKANs and CNNs exhibit contrasting behaviors across datasets. As observed in Table 6, both CKANs and CNNs exhibit poor performance on the complex datasets (CIFAR-10 and SVHN), demonstrating significant vulnerability to adversarial perturbations. On CIFAR-10, none of the models successfully classify any images generated by those attacks. On the SVHN dataset, the medium and large models manage to classify a small number of adversarially generated images. This suggests that as the dataset complexity increases, the models become more susceptible to adversarial attacks, likely due to the intricate decision boundaries required for accurate classification.

Table 6: White-box and black-box attacks on convolutional models, evaluated across MNIST, CIFAR-10, and SVHN datasets. The table reports robust accuracy for each attack, while the ‘Acc’ column indicates the model’s accuracy on clean, unperturbed inputs. Boldface values indicate the most robust model in each experiment.

Dataset	Model	White-Box				Black-Box			Acc
		FGSM	PGD	C&W	MIM	Square	SimBA	NES	
MNIST	CKAN _{small}	72.36	4.52	5.70	8.98	38.39	34.82	49.90	99.47
	CNN _{small}	68.20	0.74	0.98	2.26	27.38	9.82	33.73	99.23
	CKAN _{medium}	86.33	29.88	33.12	42.72	57.93	41.76	66.07	99.55
	CNN _{medium}	28.36	0.00	0.00	0.00	0.39	7.24	8.92	99.36
	CKAN _{large}	84.72	31.28	34.32	39.64	56.84	36.21	58.92	99.52
	CNN _{large}	55.87	0.01	0.13	0.19	5.15	2.38	5.45	99.49
CIFAR-10	CKAN _{small}	0.01	0.00	0.00	0.00	0.00	2.48	0.99	74.70
	CNN _{small}	0.01	0.00	0.00	0.00	0.00	2.28	0.69	73.03
	CKAN _{medium}	0.57	0.00	0.00	0.00	0.39	4.26	1.88	80.72
	CNN _{medium}	0.27	0.00	0.00	0.00	0.09	2.28	0.79	76.43
	CKAN _{large}	0.85	0.00	0.00	0.00	0.19	4.06	1.98	82.82
	CNN _{large}	1.94	0.00	0.00	0.00	1.48	4.46	2.38	79.26
SVHN	CKAN _{small}	1.04	0.00	0.00	0.00	7.44	10.41	7.44	92.34
	CNN _{small}	0.43	0.00	0.00	0.00	4.76	5.15	2.87	90.69
	CKAN _{medium}	6.74	0.02	0.02	0.02	16.96	16.36	13.59	94.37
	CNN _{medium}	5.00	0.05	0.05	0.04	14.78	9.52	8.03	93.41
	CKAN _{large}	8.33	0.09	0.10	0.09	19.84	16.26	13.09	94.41
	CNN _{large}	6.67	0.04	0.05	0.04	14.08	8.23	7.83	92.72

Conversely, for the simple dataset (MNIST), CKANs exhibit significantly better robustness than CNNs across all model sizes. Notably, medium and large CKAN models demonstrate remarkable resilience, successfully classifying a significant proportion of adversarial images. In contrast, CNNs perform poorly under these conditions, managing to classify only a negligible number of adversarial images.

This disparity is further corroborated by [Figure 3](#), which reveals a faster loss convergence on complex datasets (CIFAR-10 and SVHN) compared to MNIST during the attack process. The rapid convergence of the loss implies that adversarial attacks “fool” models more easily on these complex datasets.

For the non-iterative white-box attack (FGSM), the results are more definitive. On CIFAR-10, CKANs exhibit greater robustness than CNNs for medium-sized models, whereas CNNs perform better for large models. On the SVHN dataset, CKANs demonstrate slightly higher robustness across all model sizes, although the gap between CKANs and CNNs is relatively small. On the MNIST dataset, which is less complex, CKANs show significantly higher robustness than CNNs, with a noticeable gap favoring CKANs. As shown in [Figure 4](#), this gap is not consistently large across different ϵ values.

When analyzing black-box attacks, CKANs consistently demonstrated superior robustness compared to CNNs across most datasets and model sizes. CKANs outperformed CNNs in nearly all scenarios, except for the black-box attacks on large models for the CIFAR-10 dataset, where CNNs exhibited slightly better robustness, and on the small models for the CIFAR-10 dataset, where both of the configurations failed to classify any adversarial image generated with the Square attack.

In our analysis of the transferability of adversarial white-box attacks ([Table 7](#)), we observed that the results vary across datasets and model configurations, with no clear overall trend. On the MNIST dataset, fooling the medium CKAN proves to be the most challenging task. On the SVHN dataset, the large CKAN is the hardest model to fool. Conversely, on the CIFAR-10 dataset, the large CKAN is the easiest model to fool, right after the medium CNN model. Similar to the transferability results for the fully connected models, CKANs consistently demonstrated greater resistance to adversarial attacks than their CNN counterparts for each model size, with the exception of the large size on the CIFAR-10 dataset.

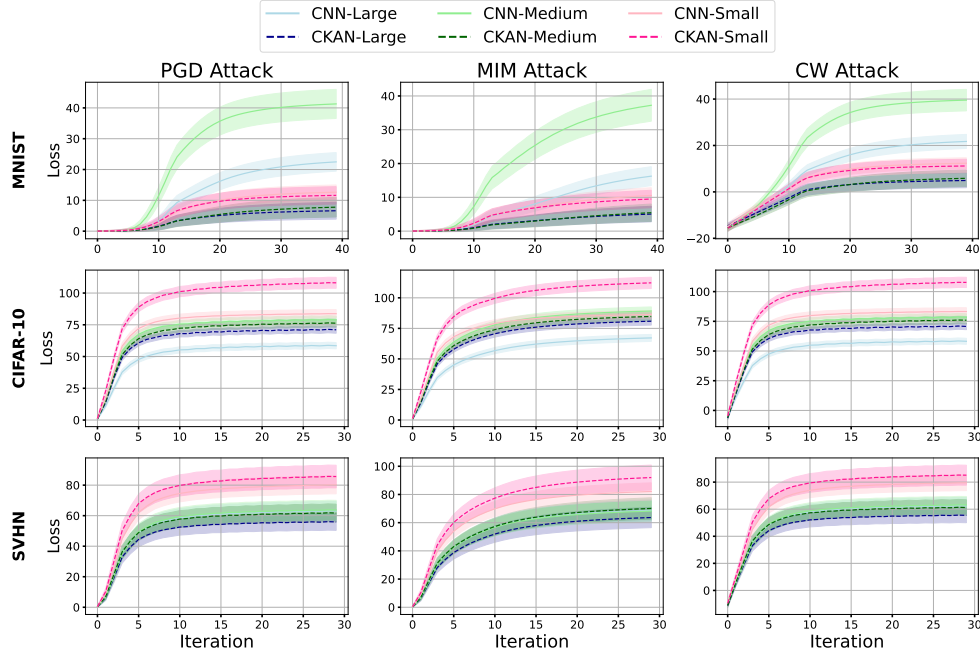


Figure 3: Comparative loss dynamics of CNNs and CKANs under PGD, MIM, and C&W attacks on the MNIST, CIFAR-10, and SVHM datasets. Each line represents the mean loss across batches per iteration, with shaded areas indicating the standard deviation.

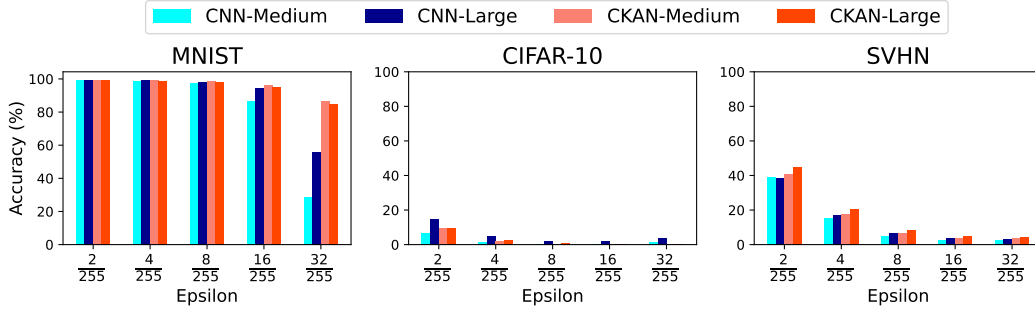


Figure 4: Robust accuracy of $\text{CNN}_{\text{medium}}$, $\text{CNN}_{\text{large}}$, $\text{CKAN}_{\text{medium}}$, and $\text{CKAN}_{\text{large}}$ as a function of FGSM attack strength, with varying ϵ values. Each bar represents a different model, and robustness is measured as the accuracy of the model against adversarial examples generated with specific ϵ values. The x-axis indicates the ϵ values used for FGSM attacks, while the y-axis shows the corresponding robust accuracy of the models.

5 Concluding Remarks

The introduction of KANs opens new avenues in adversarial machine learning, presenting both challenges and opportunities for improving model security. In this study, we explored several key aspects of KANs, focusing on their behavior in convolutional and fully connected layers, which have not been extensively investigated before.

First, we assessed the inherent robustness of FCKANs and CKANs against both white-box and black-box adversarial attacks. Our findings indicate that while FCKANs exhibit vulnerability levels similar to traditional FCNNs in smaller configurations, they demonstrate significant improvements in robustness as network size increases. This suggests that FCKANs could be particularly effective in scenarios where larger models are feasible and robustness is a priority. In contrast to the trends observed for FCKANs and FCNNs,

Table 7: Maximum transferability of our convolutional models for MNIST, CIFAR-10, and SVHN. Each row represents the model that generates the adversarial example, while each column represents the model that evaluates those examples. The value in row i , column j represents the maximum transferability between row- i model and column- j model across four attacks—FGSM, PGD, C&W, MIM—calculated using Equation 11. Unlike robustness tables, the values here indicate the percentage of images misclassified. The ‘Average’ row shows the average transferability for the attacking model. Boldface values denoting the lowest transferability, where a lower value indicates better robustness.

MNIST						
	CKAN _{small}	CKAN _{medium}	CKAN _{large}	CNN _{small}	CNN _{medium}	CNN _{large}
CKAN _{small}	95.21	2.94	3.26	4.70	20.6	4.80
CKAN _{medium}	5.30	74.81	17.78	3.32	18.95	12.01
CKAN _{large}	3.75	14.95	72.80	2.51	15.16	9.93
CNN _{small}	3.43	1.71	1.55	99.22	33.47	6.72
CNN _{medium}	3.82	2.78	3.48	7.80	100.00	46.28
CNN _{large}	3.35	4.49	5.28	5.64	6.98	100.00
Average (\downarrow)	19.14	16.94	17.35	20.53	32.52	29.95

CIFAR-10						
	CKAN _{small}	CKAN _{medium}	CKAN _{large}	CNN _{small}	CNN _{medium}	CNN _{large}
CKAN _{small}	100.00	53.39	45.62	62.95	60.19	35.45
CKAN _{medium}	70.25	100.00	93.48	79.12	90.40	68.52
CKAN _{large}	62.80	92.13	100.00	75.85	89.65	73.39
CNN _{small}	70.47	71.27	69.74	100.00	85.70	59.78
CNN _{medium}	63.23	82.31	82.52	82.50	100.00	78.68
CNN _{large}	47.53	68.93	76.86	65.07	84.74	100.00
Average (\downarrow)	69.04	78.00	78.03	77.58	85.15	69.30

SVHN						
	CKAN _{small}	CKAN _{medium}	CKAN _{large}	CNN _{small}	CNN _{medium}	CNN _{large}
CKAN _{small}	100.00	63.08	58.62	72.23	58.73	57.70
CKAN _{medium}	77.44	100.00	83.76	69.97	75.40	75.40
CKAN _{large}	76.54	86.43	99.90	69.42	76.23	75.67
CNN _{small}	68.52	55.79	52.85	100.00	63.33	62.15
CNN _{medium}	72.11	74.21	71.98	78.26	99.96	85.87
CNN _{large}	70.93	74.52	72.70	78.62	85.89	99.96
Average (\downarrow)	77.59	75.67	73.30	78.08	76.58	76.12

CKANs consistently exhibit greater robustness than CNNs of comparable sizes across most attacks. This highlights a distinct advantage of CKANs in adversarial settings, even when network sizes remain constant.

Despite these promising findings, it is crucial to recognize that KANs remain vulnerable to adversarial perturbations. Our analysis reveals that while larger configurations enhance robustness, KANs are not impervious to adversarial attacks. This ongoing fragility underscores the need for continued research to develop more effective strategies for improving the security of KANs. Addressing these vulnerabilities is essential for the reliable deployment of KANs in security-critical applications.

As the field continues to advance, further investigation is required to better understand the strengths and limitations of FCKANs and CKANs. This research will be instrumental in paving the way for the practical application of KANs in domains where security and reliability are paramount.

References

- Maksym Andriushchenko and Nicolas Flammarion. Understanding and improving fast adversarial training. *Advances in Neural Information Processing Systems*, 33:16048–16059, 2020.
- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *European conference on computer vision*, pp. 484–501. Springer, 2020.
- Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang. Recent advances in adversarial training for adversarial robustness. *arXiv preprint arXiv:2102.01356*, 2021.
- Zavareh Bozorgasl and Hao Chen. Wav-kan: Wavelet kolmogorov-arnold networks. *arXiv preprint arXiv:2405.12832*, 2024.
- Nicholas Carlini and David Wagner. Defensive distillation is not robust to adversarial examples. *arXiv preprint arXiv:1607.04311*, 2016.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57. Ieee, 2017.
- Zhaoyu Chen, Bo Li, Shuang Wu, Kaixun Jiang, Shouhong Ding, and Wenqiang Zhang. Content-based unrestricted adversarial attack. *Advances in Neural Information Processing Systems*, 36, 2024.
- George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
- Prasenjit Dey, Kaustuv Nag, Tandra Pal, and Nikhil R Pal. Regularizing multilayer perceptron for robustness. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 48(8):1255–1266, 2017.
- Chang Dong, Liangwei Zheng, and Weitong Chen. Kolmogorov-arnold networks (kan) for time series classification and robust analysis, 2024. URL <https://arxiv.org/abs/2408.07314>.
- Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9185–9193, 2018.
- Ivan Drokin. Kolmogorov-arnold convolutions: Design principles and empirical studies. *arXiv preprint arXiv:2407.01092*, 2024.
- Remi Genet and Hugo Inzirillo. Tkan: Temporal kolmogorov-arnold networks. *arXiv preprint arXiv:2405.07344*, 2024.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens Van Der Maaten. Countering adversarial images using input transformations. *arXiv preprint arXiv:1711.00117*, 2017.
- Paula Harder, Franz-Josef Pfrendt, Margret Keuper, and Janis Keuper. Spectraldefense: Detecting adversarial attacks on cnns in the fourier domain. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2021.
- Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. *arXiv preprint arXiv:1804.08598*, 2018.
- Jan Kalina, Jiří Tumpach, and Martin Holeňa. On combining robustness and regularization in training multilayer perceptrons over small data. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2022.

- P Kerlirzin and F Vallet. Robustness in multilayer perceptrons. *Neural computation*, 5(3):473–482, 1993.
- Mehrdad Kiamari, Mohammad Kiamari, and Bhaskar Krishnamachari. Gkan: Graph kolmogorov-arnold networks. *arXiv preprint arXiv:2406.06470*, 2024.
- Andrei Nikolaevich Kolmogorov. *On the representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables*. American Mathematical Society, 1961.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical Report TR-2009, University of Toronto, Toronto, Canada, 2009.
- Raz Lapid and Moshe Sipper. Patch of invisibility: Naturalistic black-box adversarial attacks on object detectors. *arXiv preprint arXiv:2303.04238*, 2023a.
- Raz Lapid and Moshe Sipper. I see dead people: Gray-box adversarial attack on image-to-text models. In *5th Workshop on Machine Learning for Cybersecurity, part of ECMLPKDD 2023*, 2023b.
- Raz Lapid, Zvika Haramaty, and Moshe Sipper. An evolutionary, gradient-free, query-efficient, black-box algorithm for generating adversarial instances in deep convolutional neural networks. *Algorithms*, 15(11):407, 2022.
- Raz Lapid, Almog Dubin, and Moshe Sipper. Fortify the guardian, not the treasure: Resilient adversarial detectors. *arXiv preprint arXiv:2404.12120*, 2024a.
- Raz Lapid, Ron Langberg, and Moshe Sipper. Open sesame! universal black-box jailbreaking of large language models. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*, 2024b. URL <https://openreview.net/forum?id=OSuyN0ncxX>.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle, James Halverson, Marin Soljačić, Thomas Y Hou, and Max Tegmark. Kan: Kolmogorov-arnold networks. *arXiv preprint arXiv:2404.19756*, 2024.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Abdulrahman Jassam Mohammed, Muhanad Hameed Arif, and Ali Adil Ali. A multilayer perceptron artificial neural network approach for improving the accuracy of intrusion detection systems. *IAES International Journal of Artificial Intelligence*, 9(4):609, 2020.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- Nicolas Papernot and Patrick McDaniel. On the effectiveness of defensive distillation. *arXiv preprint arXiv:1607.05113*, 2016.
- Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE symposium on security and privacy (SP)*, pp. 582–597. IEEE, 2016.
- Ben Pinhasov, Raz Lapid, Rony Ohayon, Moshe Sipper, and Yehudit Aperstein. XAI-based detection of adversarial attacks on deepfake detectors. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=7pBKrcn199>.

- Vinay Uday Prabhu. Kannada-mnist: A new handwritten digits dataset for the kannada language. *arXiv preprint arXiv:1908.01242*, 2019.
- Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *Advances in neural information processing systems*, 32, 2019.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1): 1929–1958, 2014.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Snir Vitrack Tamam, Raz Lapid, and Moshe Sipper. Foiling explanations in deep neural networks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=wwLQMHyLk>.
- Jiawei Tu, Bo Li, Daniel Kushner, and Dawn Song. Simple black-box adversarial attacks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 2019.
- Cristian J Vaca-Rubio, Luis Blanco, Roberto Pereira, and Màrius Caus. Kolmogorov-arnold networks (kans) for time series analysis. *arXiv preprint arXiv:2405.08790*, 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Meng Wang, Yiqin Lu, and Jiancheng Qin. A dynamic mlp-based ddos attack detection method using feature selection and feedback. *Computers & Security*, 88:101645, 2020.
- Zhipeng Wei, Jingjing Chen, Micah Goldblum, Zuxuan Wu, Tom Goldstein, and Yu-Gang Jiang. Towards transferable adversarial attacks on vision transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 2668–2676, 2022.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. *arXiv preprint arXiv:1704.01155*, 2017.
- Runpeng Yu, Weihao Yu, and Xinchao Wang. Kan or mlp: A fairer comparison. *arXiv:2407.16674*, 2024.
- Chen Zeng, Jiahui Wang, Haoran Shen, and Qiao Wang. Kan versus mlp on irregular or noisy functions. *arXiv:2408.07906*, 2024.
- Zhihao Zheng and Pengyu Hong. Robust detection of adversarial attacks by modeling the intrinsic properties of deep neural networks. *Advances in neural information processing systems*, 31, 2018.

A Image Samples of Attacks

This section provides visual examples of adversarial attacks applied to a single MNIST image, demonstrating the unique perturbation patterns across different model architectures and sizes. These visualizations compare FCNNs with FCKANs and CNNs with CKANs under white-box and black-box attack scenarios.

Fully Connected Models. Figure 5 and Figure 6 depict adversarial images and their corresponding perturbations for FCNNs and FCKANs under white-box and black-box attacks, respectively. The results highlight distinct perturbation behaviors, underscoring the unique adversarial behaviors of KANs.

Convolutional Models. Figure 7 and Figure 8 illustrate adversarial images and perturbations for CNNs and CKANs under white-box and black-box attacks, showcasing architectural influences on adversarial patterns.



Figure 5: Visualization of the adversarial images on a single MNIST image using fully connected models. Rows represent model types and sizes: $\text{FCNN}_{\text{large}}$, $\text{FCKAN}_{\text{large}}$, $\text{FCNN}_{\text{medium}}$, $\text{FCKAN}_{\text{medium}}$, $\text{FCNN}_{\text{small}}$, and $\text{FCKAN}_{\text{small}}$.

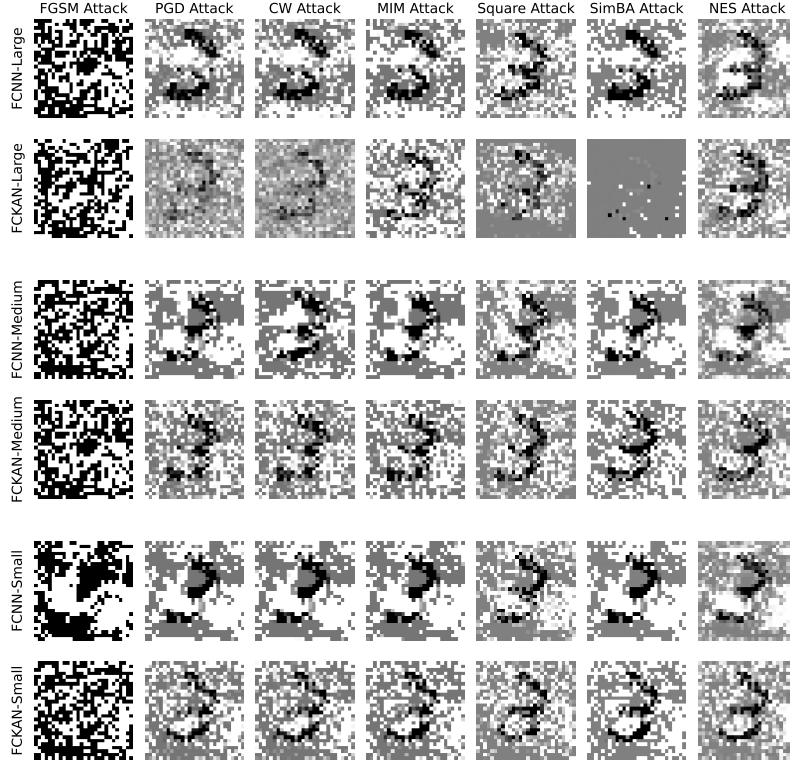


Figure 6: Visualization of the perturbations on a single MNIST image using fully connected models. Rows represent model types and sizes: $\text{FCNN}_{\text{large}}$, $\text{FCKAN}_{\text{large}}$, $\text{FCNN}_{\text{medium}}$, $\text{FCKAN}_{\text{medium}}$, $\text{FCNN}_{\text{small}}$, and $\text{FCKAN}_{\text{small}}$.



Figure 7: Visualization of the adversarial images on a single MNIST image using convolutional models. Rows represent model types and sizes: $\text{CNN}_{\text{large}}$, $\text{CKAN}_{\text{large}}$, $\text{CNN}_{\text{medium}}$, $\text{CKAN}_{\text{medium}}$, $\text{CNN}_{\text{small}}$, and $\text{CKAN}_{\text{small}}$.

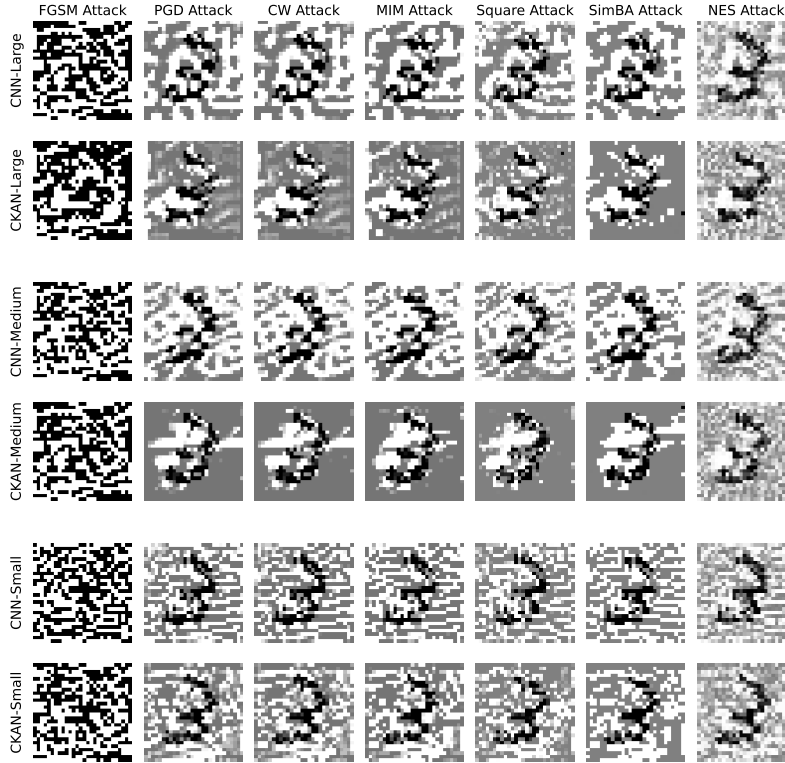


Figure 8: Visualization of the perturbations on a single MNIST image using MNIST image using convolutional models. Rows represent model types and sizes: $\text{CNN}_{\text{large}}$, $\text{CKAN}_{\text{large}}$, $\text{CNN}_{\text{medium}}$, $\text{CKAN}_{\text{medium}}$, $\text{CNN}_{\text{small}}$, and $\text{CKAN}_{\text{small}}$.

B Adversarial Attacks

In this appendix we provide a formal mathematical definition of the white-box and black box adversarial attacks.

B.1 White-Box Attacks

In the context of white-box attacks, the adversary generates an adversarial example x_{adv} from the original input x , using a specific technique, denoted as \mathcal{A} . The selection of \mathcal{A} plays a critical role in determining the attack’s success rate, which is defined as the percentage of samples that the classifier C misclassifies. The literature presents a wide array of techniques for crafting white-box attacks. In the following sections, we provide a detailed description of each attack method evaluated in this study.

B.1.1 Fast Gradient Sign Method (FGSM)

The Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014) generates adversarial examples by adding perturbations in the direction of the gradient of the loss function. Specifically, an adversarial example x_{adv} is computed as follows:

$$x_{\text{adv}} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, y; \theta)), \quad (1)$$

where \mathcal{L} is the loss function associated with a classifier parameterized by θ , $\nabla_x \mathcal{L}(x, y; \theta)$ represents the gradient of the loss with respect to the input x , and ϵ is a small perturbation. Notably, in Equation 1, only the sign of the gradient is used, with the magnitude of the perturbation controlled by ϵ . It is crucial to highlight that FGSM is a single-step attack; the adversary computes the gradient by backpropagating through the model once—and directly applies this gradient to perturb the input x .

B.1.2 Projected Gradient Descent (PGD)

The Projected Gradient Descent (PGD) method (Madry et al., 2017) is an iterative extension of the FGSM, designed to generate more robust adversarial examples. The PGD attack iteratively perturbs the input x by repeatedly applying gradient updates and projecting the perturbed image back onto an ϵ -ball around the original input. The adversarial example x_{adv} after k iterations is computed as:

$$x_{\text{adv}}^{k+1} = \Pi_{\mathcal{B}_\epsilon(x)}(x_{\text{adv}}^k + \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(x_{\text{adv}}^k, y; \theta))), \quad (2)$$

where $\Pi_{\mathcal{B}_\epsilon(x)}(\cdot)$ denotes the projection onto the ϵ -ball centered at x , α is the step size, and k indicates the iteration number.

Unlike FGSM, which performs a single gradient step, PGD iteratively refines the adversarial perturbation over multiple steps, thereby producing stronger adversarial examples. The projection step ensures that the perturbation remains within the allowable ϵ -ball, maintaining the proximity of the adversarial example to the original input.

B.1.3 Carlini & Wagner (C&W)

The Carlini & Wagner (C&W) is a powerful optimization-based adversarial attack designed to generate adversarial examples by minimizing a specially crafted loss function while ensuring that the perturbation remains imperceptible. When employing the ℓ_∞ norm, the attack seeks to find a minimal perturbation δ such that the resulting adversarial example $x_{\text{adv}} = x + \delta$ misleads the classifier. The optimization problem is formulated as follows:

$$\mathcal{L}(\delta, x) = \|\delta\|_\infty + c \cdot f(x + \delta), \quad (3)$$

where $\mathcal{L}(\delta, x)$ is the objective function to minimize, subject to:

$$x_{\text{adv}} = \text{clip}(x + \delta, 0, 1). \quad (4)$$

where $f(x_{\text{adv}})$ is a loss function that is specifically designed to ensure the misclassification of the adversarial example x_{adv} . The term $\|\delta\|_\infty$ represents the maximum perturbation applied to any individual pixel, and c is a constant that balances the trade-off between minimizing the perturbation and maximizing the loss. The clipping function ensures that the perturbed image remains within the valid input space (i.e., pixel values between 0 and 1).

The attack is typically solved through an iterative optimization procedure, where the perturbation δ and the adversarial image x_{adv} are updated at each iteration, as follows:

$$\begin{aligned} \delta^k &= \text{clip}(x_{\text{adv}}^k - x, -\epsilon, \epsilon), \\ x_{\text{adv}}^{k+1} &= x_{\text{adv}}^k + \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(x_{\text{adv}}^k, \delta^k)). \end{aligned} \quad (5)$$

where α is the step size, and k indicates the iteration number.

B.1.4 Momentum Iterative Method (MIM)

The Momentum Iterative Method (MIM) (Dong et al., 2018) is an enhancement of the basic iterative methods for generating adversarial examples, incorporating momentum to stabilize the gradient updates and avoid poor local maxima. The MIM attack iteratively updates the adversarial example x_{adv} by accumulating a momentum term, which helps to amplify gradients that consistently point in the same direction across iterations. The adversarial example x_{adv}^{k+1} after k iterations is computed as:

$$g^{k+1} = \mu \cdot g^k + \frac{\nabla_x \mathcal{L}(x_{\text{adv}}^k, y; \theta)}{\|\nabla_x \mathcal{L}(x_{\text{adv}}^k, y; \theta)\|_1}, \quad (6)$$

$$x_{\text{adv}}^{k+1} = \Pi_{\mathcal{B}_\epsilon(x)}(x_{\text{adv}}^k + \alpha \cdot \text{sign}(g^{k+1})), \quad (7)$$

where g^k denotes the accumulated gradient at iteration k , μ is the momentum factor, α is the step size, and $\Pi_{\mathcal{B}_\epsilon(x)}(\cdot)$ represents the projection onto the ϵ -ball centered at x .

By integrating the momentum term, MIM not only improves the convergence of the attack, but also enhances its effectiveness by consistently updating the perturbation in directions that contribute most to the increase in the loss function. This makes MIM particularly potent in generating adversarial examples that can transfer across different models.

B.2 Black-Box Attacks

In a black-box attack the adversary does not have direct access to the internal parameters or architecture of the target model. Instead, the attack relies on querying the model to gather information about its predictions or confidence scores. Adversarial examples x_{adv} are crafted based on these outputs, without explicit knowledge of the model’s gradients or weights. Below, we describe the specific black-box attack methodologies used in this study.

B.2.1 Transfer-based Attacks

Transfer-based attacks leverage the phenomenon of transferability, where adversarial examples generated for one model are also misclassified by other models. This concept of adversarial example transferability was first introduced by [Szegedy et al. \(2013\)](#).

Let \mathcal{A} denote the white-box adversarial attack method, and let f_θ represent the model under attack. The attack can be formalized as:

$$\mathcal{A}_{f_\theta}(\mathbf{x}, y) = \mathbf{x}_{\text{adv}}, \quad (8)$$

where (\mathbf{x}, y) corresponds to the input-label pairs.

Let g_ϕ denote the target model, which is employed to evaluate the effectiveness of the generated adversarial examples. An adversarial example \mathbf{x}_{adv} is considered to have successfully transferred via a given attack if and only if the following conditions hold:

$$g_\phi(\mathcal{A}_{f_\theta}(\mathbf{x}, y)) \neq y \quad \text{and} \quad f_\theta(\mathbf{x}) = g_\phi(\mathbf{x}) = y. \quad (9)$$

To quantify the transferability of adversarial examples between f_θ and g_ϕ for a specific attack, we define the transferability metric as:

$$t_{f_\theta, g_\phi}^{\mathcal{A}} = \frac{1}{m} \sum_{i=1}^m \begin{cases} 1 & \text{if } g_\phi(\mathcal{A}_{f_\theta}(\mathbf{x}_i, y_i)) \neq y_i \\ 0 & \text{otherwise} \end{cases}, \quad (10)$$

where m is the number of samples.

To identify the most effective attack in terms of transferability between the generator model and the evaluator model, we calculated the overall transferability metric, $t_{f_\theta, g_\phi}^{\text{total}}$, as follows:

$$t_{f_\theta, g_\phi}^{\text{total}} = \max_{\mathcal{A} \in \{\text{MIM, FGSM, C\&W, PGD}\}} \left(t_{f_\theta, g_\phi}^{\mathcal{A}} \right). \quad (11)$$

B.2.2 Square Attack

The Square Attack ([Andriushchenko et al., 2020](#)) is a query-efficient black-box adversarial attack via random search that operates by iteratively applying random perturbations to the input image. The target is to minimize the following loss:

$$L(f(x_{\text{adv}}), y) = f_y(x_{\text{adv}}) - \max_{k \neq y} f_k(x_{\text{adv}}), \quad (12)$$

At each iteration, the attack selects a random patch of the image and applies perturbations to it in the form of a square grid. The perturbations are scaled by a predefined ϵ value, ensuring that they stay within the specified ℓ_∞ norm constraint. The adversarial example is updated iteratively as follows:

$$x_{\text{adv}}^{k+1} = x_{\text{adv}}^k + \Delta_{\text{square}}, \quad (13)$$

where Δ_{square} represents the perturbation applied to the selected patch. The key advantages of the Square Attack are its simplicity and effectiveness, making it a strong baseline for black-box adversarial evaluation. Additionally, its random search mechanism ensures good exploration of the input space, increasing the chances of finding adversarial examples.

B.2.3 Simple Black-Box Attack (SimBA) in DCT Space

The Simple Black-Box Attack (SimBA) (Tu et al., 2019) is a query-efficient method for crafting adversarial examples. Unlike gradient-based attacks, SimBA perturbs the input image in a random manner (e.g., DCT or Fourier space) and evaluates the effect of the perturbation on the model’s confidence score. If the perturbation reduces the confidence of the correct label, it is retained; otherwise, it is discarded.

In the DCT (Discrete Cosine Transform) space, SimBA exploits the frequency-domain representation of images to prioritize perturbations on low-frequency components. This ensures that the adversarial examples remain perceptually similar to the original images while efficiently reducing the model’s confidence in its predictions.

The SimBA attack updates the adversarial example iteratively as follows:

$$x_{\text{adv}}^{k+1} = x_{\text{adv}}^k + \Delta_{\text{DCT}}, \quad (14)$$

where Δ_{DCT} represents the perturbation in the DCT basis. The step size of the perturbation is controlled by the ϵ parameter.

SimBA’s simplicity and efficiency make it a widely-used black-box attack, particularly in scenarios where query limits are enforced.

B.2.4 Natural Evolution Strategy (NES)

The Natural Evolution Strategy (NES) (Ilyas et al., 2018) is a gradient-free optimization method that estimates the gradient of a loss function by sampling noise vectors and using finite differences to approximate the effect of perturbations. The method is particularly suited for black-box attacks, where access to model gradients is unavailable. It iteratively updates the adversarial example by approximating the gradient of the loss with respect to the input.

For a given input x , the adversarial perturbation is computed by estimating the gradient using a set of n random noise vectors $u_i \sim \mathcal{N}(0, I)$, sampled from a standard normal distribution. The update rule is as follows:

$$\hat{\nabla}^k = \frac{1}{2\sigma n} \sum_{i=1}^n [\mathcal{L}(x_{\text{adv}}^k + \sigma u_i, y) - \mathcal{L}(x_{\text{adv}}^k - \sigma u_i, y)] u_i. \quad (15)$$

Using this gradient estimate $\hat{\nabla}^k$, the adversarial example x_{adv} is updated iteratively as follows:

$$x_{\text{adv}}^{k+1} = \Pi_{\mathcal{B}_\epsilon(x)} \left(x_{\text{adv}}^k + \alpha \cdot \text{sign}(\hat{\nabla}^k) \right). \quad (16)$$