



A Survey on Adversarial Perturbations and Attacks on CAPTCHAs

Suliman A. Alsuhibany 匝

Department of Computer Science, College of Computer, Qassim University, Buraydah 51452, Saudi Arabia; salsuhibany@qu.edu.sa

Abstract: The Completely Automated Public Turing test to tell Computers and Humans Apart (CAPTCHA) technique has been a topic of interest for several years. The ability of computers to recognize CAPTCHA has significantly increased due to the development of deep learning techniques. To prevent this ability from being utilised, adversarial machine learning has recently been proposed by perturbing CAPTCHA images. As a result of the introduction of various removal methods, this perturbation mechanism can be removed. This paper, thus, presents the first comprehensive survey on adversarial perturbations and attacks on CAPTCHAs. In particular, the art of utilizing deep learning techniques with the aim of breaking CAPTCHAs are reviewed, and the effectiveness of adversarial CAPTCHAs is discussed. Drawing on the reviewed literature, several observations are provided as part of a broader outlook of this research direction. To emphasise adversarial CAPTCHAs as a potential solution for current attacks, a set of perturbation techniques have been suggested for application in adversarial CAPTCHAs.

Keywords: cyber security; CAPTCHAs; deep learning; adversarial examples; robustness; usability

1. Introduction

The protection of websites and electronic services against potential attacks by using numerous forms has received increasing attention due to their importance in our lives. One of these forms is a CAPTCHA (Completely Automated Public Turing test to tell Computers and Humans Apart); a CAPTCHA is one of the Human Interaction Proofs (HIP) systems [1]. CAPTCHAs may come in different forms, including video-based, text-based and image-based CAPTCHAs. For example, the text-based CAPTCHA prompts users to recognize a text, a task that state-of-the-art text recognition programs cannot complete successfully; due to its various advantages, this form of CAPTCHA is the most commonly deployed type used by websites to date [2].

Deep learning (DL) and machine learning (ML) are considered the newest artificial intelligence (AI) revolutions in how we live, work, study and discover [3]. The capabilities of ML and DL are continually opening up new opportunities for progress in different areas, such as health, education, energy, economics, and other environments. Furthermore, in the security domain, the AI has many beneficial applications, particularly in the cyber-security domain. It can be adopted in both the defensive and offensive cyber categories. On the one hand, the defensive aspect monitors the software and any potentially suspicious behaviour. On the other hand, the artificial intelligence-based offensive aspect can be seen in malicious actions adopted by attackers that change anomaly behaviours. One aspect found within the defensive category is an adversarial perturbation. The adversarial perturbation entails the noise added to an original image in order to create a modified version of the original image (i.e., an adversarial example). This modified version would fool an attack that uses a machine learning technique, such as the deep learning networks referred to in [3–6]. On the other hand, one aspect found in the offensive category is the adversarial attack. Adversarial attacks are inputs specifically crafted to fool detection systems.



Citation: Alsuhibany, S.A. A Survey on Adversarial Perturbations and Attacks on CAPTCHAs. *Appl. Sci.* **2023**, *13*, 4602. https://doi.org/ 10.3390/app13074602

Academic Editors: Nawin Raj and Jason Brown

Received: 15 October 2022 Revised: 11 February 2023 Accepted: 17 February 2023 Published: 5 April 2023



Copyright: © 2023 by the author. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). Over the past decade, several studies have been proposed to break CAPTCHAs. Early works which are summarized in [7] typically follow three main steps: (1) using heuristics methods to filter the background patterns such as the line noise (i.e., a pre-processing); (2) using character segmentation techniques (i.e., the edge corners and fuzzy logic segmentation technique); and (3) using a ML model in order to recognize the segmented character. As an enhancement to the attack approaches, recent works have combined the segmentation and recognition steps together, such as in [8–13]. Most of these works utilize a deep neural network (DNN) model to directly recognize CAPTCHA characters without the segmentation step. To defend against them, several technologies have been proposed, such as in [3–6,14–18].

Owing to the potential of adversarial attacks on CAPTCHAs, several research efforts have been accomplished in this area, as shown in Figure 1. Since there has not yet been a survey that shows these research efforts, this paper provides a review for these efforts and discusses their major contributions while identifying issues for further researches.



Figure 1. Summarized historically-distributed studies for adversarial perturbations and attacks on CAPTCHAs.

The literature search was conducted using Engineering Village (Engineering Village is an engineering literature search tool, which provides access to 14 engineering literature and patent databases), where we collected 218 related papers. The papers were filtered according to defined terms shown in Section 2, and information was extracted from these papers according to their contributions to adversarial attacks on CAPTCHAs. A total of 51 papers published between 2017 and 2022 were selected.

The main contributions of this paper are as follows. First, the paper presents a comprehensive survey of published studies in the topic area of adversarial attacks and perturbations on CAPTCHAs, highlighting their contributions and technological advances. Second, the paper critically analyses these related works from different perspectives, in which knowledge gaps can be identified as well as issues for further researches. Third, the references cited in this paper can be a useful guide into this topic area.

The structure of this paper is as follows. Section 2 describes the common terms related to adversarial CAPTCHAs. Deep learning attacks against CAPTCHAs are reviewed in Section 3. Section 4 explains the adversarial CAPTCHAs. The effectiveness of adversarial CAPTCHAs is discussed in Section 5. Several observations based on the analysed studies are presented in Section 6. Section 7 reviews other techniques for generating adversarial CAPTCHAs. The paper is concluded in Section 8.

2. Terms Definition

This section describes the main technical terms used in the literature related to adversarial CAPTCHAs.

- An adversarial CAPTCHA is the changed version of a clean CAPTCHA that is intentionally perturbed in order to confuse a machine learning method.
- An adversarial perturbation represents the noise that reshapes the CAPTCHA image to make it adversarial.
- An adversarial training utilizes the adversarial CAPTCHA aside from the CAPTCHA images to train the machine learning models.
- An adversarial example detector is a mechanism that detects whether or not an image is an adversarial example.
- A perturbation domain refers to either: (1) the frequency domain, which inserts perturbations into a single character image and combines different character images into one CAPTCHA or (2) the space domain, which directly injects perturbations into CAPTCHAs.
- DL algorithms include the convolutional neural network (CNN), the recurrent neural network (RNN), the recursive cortical network (RCN), the deep neural network (DNN), and the artificial neural network (ANN).
- ML algorithms include the support vector machine (SVM), the decision tree (DT), the random forest (RF), the logistic regression (LR) and the k-nearest neighbour (kNN).

3. Deep Learning Attacks against CAPTCHAs

This section reviews the literature that discusses the art of utilizing DNNs with the aim of breaking CAPTCHAs. It also reviews the weaknesses of a deep neural network in the context of the image classification. Despite this fact, it is worth noting that the DL performs various tasks for the computer vision with a high accuracy.

3.1. Deep Learning Attack

DNN models have been improved in terms of visual recognition tasks, and they have become the centre of attention since the impressive performance of CNNs. Recent studies (e.g., [8–12]) have demonstrated threats from the automated CAPTCHA attacks using different DL techniques that reveal remarkable accuracy. This leads to much difficulty in designing usable and secure CAPTCHAs. Despite these high accuracies, the DL networks are surprisingly vulnerable to adversarial perturbations, even with small perturbations to CAPTCHA images that are still readable to the human visual system.

In particular, a novel approach to solving CAPTCHAs using the ML is proposed in [19]. This approach attacks the segmentation and recognition problems simultaneously, which allows the exploitation of the information and context that are not available when they are done sequentially. Also, this approach can be generalized due to the automation of the segmentation and recognition processes. Moreover, a probabilistic generative model for the vision is introduced in [11] in which the message-passing-based inference handles the recognition, segmentation, and reasoning together. The results showed excellent generalization and occlusion-reasoning capabilities. The performance of this model is outstanding using DNNs. Furthermore, a framework based on the Generative Adversarial Network (GAN) is introduced in [20,21] in which the character segmentation algorithm is improved.

In [10], a comprehensive study of reCaptcha is conducted in order to explore how the risk analysis process is influenced by each aspect of the request. Based on this study, a novel low-cost attack that leverages DL technologies for the semantic annotation of images is designed. The results showed an interesting accuracy for solving reCaptcha challenges.

In [22], novel segmentation and recognition methods are proposed which apply simple image processing techniques such as the pixel count methods along with an ANN for textbased CAPTCHAs. Popular CCT (Crowded Characters Together) based CAPTCHAs are targeted for evaluating the proposed method. The overall accuracy was 53.2%. This study explores not only the flaws in the text-based CAPTCHA design, but also finds an approach to segment and recognize the connected characters from images.

A synthetic training data approach is used in [23] in order to train a neural network for breaking the text-based CAPTCHAs. The results showed a remarkable recognition performance on the real-world CAPTCHAs currently used such as on Facebook. Likewise, a generic yet effective text-based CAPTCHA solver based on the GANs is proposed in [8]. This solver requires significantly fewer real CAPTCHAs as it is the first learning for a CAPTCHA synthesizer to automatically generate synthetic CAPTCHAs. The results demonstrated a significantly higher accuracy on all selected schemes. Moreover, a simple, generic, and fast attack on text-based CAPTCHAs is proposed in [24] using DL techniques. This attack demonstrates a high success rate in breaking not only English CAPTCHAs, but also some Chinese CAPTCHAs that use a larger character set. Moreover, this attack is enhanced in [25]. Similarly, an automatic attacking method is proposed in [26] to deal with the variable-length Chinese character CAPTCHAs with noises. The results of evaluating the proposed method showed the ability of breaking the mixed character CAPTCHAs.

A study in [27] achieved a tremendous progress in breaking accuracy cracking CAPTCHAs using the conditional deep convolutional generative adversarial networks (cDCGAN) and CNN. Similarly, a generic solver combining unsupervised learning and representation learning to automatically remove the noisy background and solve text-based CAPTCHAs is introduced in [28]. Further, a customized DNN model is developed in [29] that results a high cracking accuracy rate of 98.94% and 98.31% for the numerical and the alpha-numerical test datasets, respectively. Based on these results, this study identified some efficient techniques to enhance the robustness of the text-based CAPTCHAs.

A dynamic approach is proposed in [30] that predicts the text-based CAPTCHAs. In particular, this approach firstly uses a pre-processing step through several techniques like Erosion, Dilation, and Binarization in order to remove the noise from the CAPTCHA. This CAPTCHA is then fed to the CNN that generates a feature vector. This feature vector is then passed to the long short-term memory (LSTM) which generates a sequence of characters that reflect the outcome to the users. Additionally, in [31], an efficient CNN model is introduced that uses attached binary images to recognize CAPTCHAs without the segmentation of CAPTCHAs into individual characters. The results revealed the strength of the introduced model in recognizing CAPTCHAs' characters. Additionally, a study in [12] proposed an efficient CAPTCHA solver that periodically retrains the solver model when its accuracy drops using an incremental learning. This proposed solver requires a small amount of data while achieving a high accuracy. The results of evaluating the proposed solver demonstrated that the existing defense methods based on a text-based CAPTCHA scheme and an image-based CAPTCHA scheme can be bypassed.

A transfer learning-based approach is proposed in [32] in which the attack complexity and the cost of labeling samples are reduced by pre-training the model. This model randomly produced samples and fine-tuning the pre-trained model with a small number of real-world samples. Furthermore, a GAN is applied to refine the samples sequentially. The results of evaluating this approach showed that the cost of data preparation is reduced while preserving the model's attack accuracy.

For enhancing the accuracy of CAPTCHA attacks, a simple preprocessing approach is introduced in [33]. This approach includes a data selector, which automatically filters out a training data set with training significance, and a data augmenter, which applies four different image noises to generate different CAPTCHA images. The results showed that the attack accuracy rate is improved after applying this approach. In addition, the brute-force attack with transfer learning is combined in [34] for breaking the text-based CAPTCHAs. This achieves 80% classification accuracy for a five-digit text-based CAPTCHA scheme.

An efficient end-to-end attack method based on cycle-consistent generative adversarial networks (Cycle-GANs) is proposed in [9]. This method focuses on reducing the cost of data labeling. The results demonstrated efficiently the performance of this method in terms of breaking CAPTCHAs.

A recent study in [35] proposed a fast CAPTCHA solver based on GANs to simplify the CAPTCHA images before segmenting and recognizing characters. This effectively breaks the text-based CAPTCHAs with complex security features by a small amount of labeled data. Results showed that the proposed solver achieved a high success rate of over 96% character accuracy.

A study in [36] analyzed the security level of exist audio-based CAPTCHA schemes against such attacks using ML and DL models. The experimental results reveal that audio-based CAPTCHAs that had no or medium background noise could be broken with nearly 99% to 100% accuracy, whereas the attack accuracy is decreased to 85% with high background noise.

The aforementioned studies are summarized in Table 1.

Table 1. Summarizing recent automatic CAPTCHA attack studies armed with different DL techniques.

Study	Year	Approach	Accuracy *
[19]	2014	ML	51.09%
[10]	2016	Leveraging DL technologies for the semantic annotation	83.5%
[11]	2017	Recursive cortical network (RCN)	94.3%
[22]	2017	Simple image processing techniques for the segmentation and recognation	53%
[23]	2017	A neural network trained using synthetic data	90%
[8]	2018	A GAN	96%
[24]	2018	A generic attack based on DL techniques	90%
[27]	2018	Conditional deep convolutional generative adversarial networks and CNN	98.4%
[26]	2019	Pre-processing, segmentation and recognition	96.8%
[25]	2019	Utilising a CNN and an attention-based RNN	97.3%
[28]	2020	Unsupervised learning and representation learning	94.5%
[21]	2020	Generative adversarial network (GAN)	92.08%
[29]	2020	Using an automated DL based solution	98.94%
[30]	2020	CNN and LSTM	85.97%
[31]	2020	An efficient CNN model that uses attached binary images to recognize CAPTCHAs	92.68%
[12]	2020	Incremental learning	87.37%
[32]	2020	A transfer learning-based approach	96.9%
[34]	2021	Combining a brute-force attack with transfer learning	80%
[33]	2021	Filtering and enhancing the accuracy attack rate	8.31%
[9]	2021	Cycle-consistent generative adversarial networks	97%
[35]	2021	A CAPTCHA transformation model based on GAN	96%
[37]	2022	The adversarial training strategy	87%
[38]	2022	CNN and RNN-based automatic speech recognition systems	49.76%

* If the study contains multiple experiments, the highest accuracy was used.

In the extant literature, the minimal security level standard for designing a CAPTCHA is varied. That is, the requirements for the possibility of an automatic bypassing of the CAPTCHA system [i.e., the false positive rate (FPR)] range from 0.6% to approximately 5% [3].

This section obviously demonstrates the impact of DL for breaking CAPTCHAs with high accuracy. This is accomplished by applying different DL techniques from 2014 to the present date. Although most of the works are from 2020, they may increase in number in the coming years, resulting in probably more robust CAPTCHA schemes.

3.2. Utilised Deep Learning Algorithms

There are different DL algorithms used in the literature to break CAPTCHAs. In particular, the DL algorithms that are most commonly used to break CAPTCHAs are shown in Table 2. This table also includes ML algorithms that are stated in the literature for the same purpose.

Table 2. Most commonly DL and ML algorithms used to break CAPTCHAs.

Study	DL Algorithms					ML Algorithms				
Study	CNN	RNN	RCN *	DNN	ANN	SVM	DT	RF	LR	KNN
[8]	\checkmark									
[9]	\checkmark	\checkmark								
[10]		\checkmark								
[11]			\checkmark							
[12]				\checkmark						
[19]						\checkmark				\checkmark
[28]	\checkmark									
[33]	\checkmark					\checkmark	\checkmark	\checkmark	\checkmark	
[35]	\checkmark									
[21]	\checkmark									
[34]	\checkmark									
[29]	\checkmark									
[30]	\checkmark									
[26]	\checkmark									
[24]	\checkmark									
[25]	\checkmark	\checkmark								
[31]	\checkmark									
[27]	\checkmark									
[23]	\checkmark	\checkmark								
[22]					\checkmark					
[32]	\checkmark	\checkmark								
[37]				\checkmark						
[38]	\checkmark	\checkmark								

* Recursive cortical network (RCN) inspired from system neuroscience.

It seems that the DL algorithms have been utilized more than the ML algorithms, especially in recent years. The main reason might be the high accuracy level of breaking CAPTCHAs using the DL algorithms. However, the DL algorithms are highly vulnerable to perturbations as will be discussed in the following section.

3.3. Deep Learning Attack and Adversarial Perturbation

The gap between the human and machine in terms of solving problems that have been typically used in CAPTCHAs has reduced by DL technique. A study in [19] insinuated that this marks the end of CAPTCHAs. However, most, if not all, of the existing DL attacks have various disadvantages; it is particularly time-consuming and expensive to form an attack process with high complexity and to manually collect and label a huge number of samples to train a DL recognition model [24]. More importantly, the DL technique is vulnerable to small perturbations of input that are still readable by humans; this can cause

misclassification [3]. Such adversarial perturbations can cause a neural network classifier to completely change its prediction regarding a given image. This has been demonstrated in [39], which is a ground-breaking study demonstrating an intriguing weakness of deep neural learning networks in the context of the image classification.

4. Adversarial CAPTCHAs

This section explains in details the foundation of adversarial CAPTCHAs, techniques of generating adversarial CAPTCHAs, the domain of perturbations, the security of adversarial CAPTCHAs, and the usability of adversarial CAPTCHAs.

4.1. Foundation of Adversarial CAPTCHAs

Generally speaking, ML models have a limitation in regard to adversarial manipulations, particularly in terms of distinguishability measures among classes [40]. Aside from ML, deep neural networks possess the same limitation, as discussed previously in [39]. That is, adding a trivial tailored noise piece causes misclassification with a high confidence. Figure 2 shows a sample of a normal CAPTCHA and its adversarial version.



Figure 2. A sample of a normal CAPTCHA and its adversarial version. (**a**) A normal CAPTCHA. (**b**) The adversarial version of (**a**).

This limitation renders the idea of applying adversarial examples as the source for adversarial CAPTCHAs. Although this idea can be applied in CAPTCHAs as well as other security applications, the resistance to removal attacks—which can remove the added noise as emphasized in [3]—should be taken into consideration. As such, the techniques utilised to generate an adversarial CAPTCHA should be sufficiently robust to these removal attacks. In the following section, this is discussed in further details.

4.2. Used Datasets

In the literature, various datasets have been used to generate adversarial CAPTCHAs. In particular, there are three main types of datasets: real CAPTCHA, image-based and textbased datasets, as shown in Table 3. The most used datasets for training the image-based CAPTCHAs were ImageNet datasets; while MNIST datasets were most commonly used for the text-based CAPTCHAs. Although an appeal was presented in using MNIST datasets for the text-based CAPTCHAs, empirical results from a study in [3] showed that images consisting of two colours are poor sources of adversarial examples, such as those in MNIST.

Table 3. Used datasets in the literature to generate adversarial CAPTCHAs.

Study -	Datasets						
	Real CAPTCHA	Image-Based	Text-Based	Audio-Based			
[3]	-	ILSVRC-2012		-			
[5]	-	ImageNet	MNIST	-			
[14]	-	-		-			
[15]		-	-	-			
[4]	Real CAPTCHA	-	-	-			
[6]	-	-	-	-			

Study -	Datasets						
	Real CAPTCHA	Image-Based	Text-Based	Audio-Based			
[18]	-	ImageNat	-	-			
[13]	-	Imageiver	-	-			
[16]	-	-	MNIST & EMNIST	-			
[17]	-	Caltech-200 bird	-	-			
[37]	-	-	-	LibriSpeech			
[38]	-	-	-	LibriSpeech, Google Speech Commands and LDC: ISOLET Spoken Letter			

Table 3. Cont.

Furthermore, the features of used datasets to generate the adversarial CAPTCHAs are detailed in Table 4.

Table 4. Characteristics of used Datasets.

Dataset	Characteristics
Real CAPTCHA	A collection of real CAPTCHA set
ILSVRC-2012	The ImageNet Large Scale Visual Recognition Challenge (ILSVRC) is used to evaluate algorithms for image classification at large scale. It contains more than 1.25 million training images and 50 thousand validation images.
ImageNet	It is used to evaluate algorithms for image classification at large scale. It contains more than 14 million images and 20,000 categories.
Caltech-200 bird	It includes more than 11,000 images for bird species, and 200 bird categories.
MNIST	The Modified National Institute of Standards and Technology (MNIST) is a large database of handwritten digits that can be used for different image processing system.
EMNIST	It extends the MNIST dataset by including much data and has a set of handwritten digits with a 28×28 format.
LibriSpeech	It is a large-scale database of roughly 1000 h of 16 kHz read English speech
Google Speech Commands	It has 65,000 one-second-long utterances of 30 short words by various people.
ISOLET Spoken Letter	It contains two samples of each English letter being spoken by 75 males and 75 females of varying ages.

4.3. Techniques for Generating Adversarial CAPTCHA

There are several techniques used for generating the adversarial CAPTCHAs. Each technique carries several different aspects. In the data analysis from [15], it can be seen that humans and algorithms demonstrate different vulnerabilities to visual distortions. That is, adversarial perturbation is considerably bothersome to an algorithm yet friendly to a human.

Since a normal adversarial noise is not enough to achieve a secure CAPTCHA scheme, a study in [3] introduced immutable adversarial noise (IAN) that is resistant to the removal attempts. The results showed that this IAN offers a good security and usability levels.

In 2018, the influence of adversarial examples on CAPTCHA robustness is analyzed in [13] using two generation algorithms: FGSM and UAPM. The former is fast, convenient and widely used, whereas the latter attempts to discover a universal perturbation vector by aggregating atomic perturbation vectors that guide the successive data points to the decision boundary of the classifier.

In 2020, an attempt to design secure CAPTCHA questions that are smoothly solvable to humans is introduced in [15]. Specifically, while analyzing the data, it observes that adversarial perturbation is meaningfully annoying to the algorithm yet friendly to humans.

Four modules of multi-target attacks are proposed to address the characteristics of the character-based CAPTCHA cracking. The results demonstrated the effectiveness of the proposed solution. An approach to synthesize a robust CAPTCHA generator that is resistant and cannot be easily attacked by such recognition algorithms is proposed in [14]. Additionally, a practical adversarial CAPTCHA generation system is introduced in [6] that can defend against DL-based CAPTCHA solvers, and apply it on a comprehensive online platform with close to a billion users. By applying adversarial learning techniques in a novel manner, the proposed generation system can make an effective adversarial CAPTCHA to significantly reduce the success rate of attackers. The results showed that the proposed approach can serve as a key enabler for generating robust CAPTCHAs in practice. Also, a novel CAPTCHA scheme based on adversarial examples is proposed in [18]. Typically, adversarial examples are used to lead an ML model astray. The basic idea behind this novel scheme is to make a "good use" of such mechanisms in order to increase the robustness and security of existing CAPTCHA schemes. The results showed that the proposed scheme generates CAPTCHA samples that are usable, whilst being efficiently resistant against such sophisticated ML-based bot solvers.

In 2021, a text-based CAPTCHA generation technique named Robust Text CAPTCHA (RTC) is proposed in [16]. The evaluation results showed that the proposed method has a high usability level and robust level against different defensive techniques. Furthermore, a CAPTCHA approach relies on the cognition process and semantic reasoning and a novel model to generate the CAPTCHA are introduced in [17]. Three features are synthesized by this approach: sentence, object, and location towards a multi-conditional CAPTCHA that resists the attack of the CNN classification. The results of evaluating this approach revealed that the classification of ResNet-50 only achieves 3.38% accuracy. Besides, a structure for text-based and image-based adversarial CAPTCHA generation on top of state-of-the-art adversarial image generation techniques is proposed in [5]. Based on this framework, an adversarial CAPTCHA generation technique named Jacobian-based Saliency Map Attack (JSMA) is designed and implemented. The results demonstrate that the security of normal CAPTCHAs is significantly improved while maintaining similar usability.

In 2022, an audio adversarial CAPTCHA scheme is designed and implemented in [37] to improve the security level of audio CAPTCHAs, since they have seen highly vulnerable to automatic speech recognition systems. This scheme exploits the audio adversarial examples as a security feature against automatic speech recognition systems. The usability and security aspects of this scheme are evaluated. The results highlighted that the proposed scheme enhances the security level of traditional audio CAPTCHA schemes while preserving the usability level. It is interesting to note that this scheme has a high security level even when the attackers have a complete knowledge about existing attacks. Moreover, a study in [38] proposed a secure audio CAPTCHA approach by modifying current audio samples in order to mitigate automatic speech recognition systems. By applying a new algorithm for high transferability called Yeehaw Junction, the evaluation of this approach showed a good robustness level against automatic speech recognition attacks as well as being highly usable.

The aforementioned techniques are summarized in Table 5.

Table 5. Summarizing the used techniques for generating adversarial CAPTCHAs.

Study	Year	Technique
[3]	2017	Immutable adversarial noise (IAN)
[13]	2018	FGSM and the universal adversarial perturbation method (UAPM)
[15]	2020	Two types of visual distortions: (1) Gaussian white noise and (2) FGSM
[14]	2020	The EOT (Expectation Over Transformation) algorithm and a generative adversarial network (GAN) as a generative model
[6]	2020	Connectionist Temporal Classification (CTC)

Study	Year	Technique
[18]	2020	Semantic image generator
[4]	2021	Multi-label classification training text CAPTCHAs as a pre-training model
[16]	2021	Scaled Gaussian translation with channel shifts attack (SGTCS)
[17]	2021	Cognition to tackle the emerging challenge generative adversarial network (GAN) as a generative model
[5]	2021	Jacobian-based Saliency Map Attack (JSMA)
[37]	2022	FGSM and Projection gradient descent (PGD)
[38]	2022	Yeehaw Junction approach

Table 5. Cont.

4.4. Domain of Perturbations

The perturbation is an essential part in generating adversarial CAPTCHAs, specifically in terms of adding noise. Thus, there are two main domains that can be applied in the target CAPTCHA image: the frequency and space domains. The frequency domain inserts perturbations into a single character image, and subsequently combines different character images into one CAPTCHA image, while the space domain directly injects perturbations into CAPTCHAs. The used domains in each study are shown in Table 6; most of these studies utilise the space perturbation domain rather than the frequency domain. According to [5], the first study to add perturbations in the frequency domain, perturbations added in the space domain seem easier to remove when using a perturbation removal method, as this is considered as a local change. In contrast, perturbations added in the frequency domain are a global change; hence, they are hard to remove. Despite the advantage of the frequency domain, improvements are needed regarding the inconsistency of the transferability property among various ML models.

Study	CAPTCHA Scheme	Perturbations Domain	
[3]	Image-based	Space	
[1.4]	Image-based	Crocco.	
[14]	Text-based	- Space	
[15]	Text-based	Space	
	Image-based		
[13]	Text-based	Space	
	Click-based	_	
[4]	Text-based	Space	
[5]	Image-based	Space	
[18]	Text-based	Frequency	
[6]	Text-based	Space	
[16]	Text-based	Frequency	
[17]	Text-image-based	Space	
[18]	Image-based	Space	
[37]	Audio-based	Space	
[41]	Image-based	Frequency	
[38]	Audio-based	Space	

Table 6. The used perturbations domain in the state-of-the-art of CAPTCHAs.

4.5. Robustness of Adversarial CAPTCHA

The robustness of adversarial CAPTCHAs reflects their resistance to perturbation removal methods. That is, it should be difficult for any computationally efficient methods—such as filtering, ML or DL approaches as shown in Section 3.1—to remove the added noise. Moreover, a study in [3] stated that some of the current techniques for adversarial example constructions are not satisfactorily robust to such attacks. On the other hand, most recently used methods demonstrate a robustness level against such perturbation removal methods. Table 7 summarized both recommended and not recommended perturbation techniques based on state-of-the-art adversarial machine learning.

Perturbation Techniques	Recommendation	Remarks
Optimization methods [39]	Natao and a d	Too slow and easy to remove as demonstrated in [39]
FGSM [42]	Not recommended	Easy to remove, as demonstrated in [42]
JSMA [5]		Its robustness level is empirically evaluated in [5] and shows a sufficient resistance level. Also, no attack has been reported yet.
IAN [3]		Its robustness level is empirically evaluated in [3] and shows a sufficient resistance level. Also, no attack has been reported yet.
EOT [14]		Its robustness level is empirically evaluated in [14] and shows a sufficient resistance level. Also, no attack has been reported yet.
CTC [6]	Recommended	Its robustness level is empirically evaluated in [6] and shows a sufficient resistance level. Also, no attack has been reported yet.
SGTCS [16]		Its robustness level is empirically evaluated in [16] and shows a sufficient resistance level. Also, no attack has been reported yet.
GAN [14,17]		Its robustness level is empirically evaluated in [14,17] and shows a sufficient resistance level. Also, no attack has been reported yet.
Semantic Image Generator [18]		Its robustness level is empirically evaluated in [18] and shows a sufficient resistance level. Also, no attack has been reported yet.

Table 7. Recommended and not recommended perturbation techniques.

4.6. Usability of Adversarial CAPTCHA

The usability of adversarial CAPTCHAs reflects user satisfaction in terms of avoiding an effect on the human perception of the image content. To date, most of the conducted studies have experimentally evaluated the usability aspect of the proposed adversarial CAPTCHA scheme. The results of our study have shown a generally high satisfaction level with the developed schemes. However, the sample size of some experiments, such as in [4,17], was too small; this decreases the power of the study and rises the margin of error, which can render the study meaningless. Figure 3 shows the sample size in each study to date. In particular, the samples sizes in [6,14] are not reported, while [16] used Amazon Mechanical Turk (MTurk), https://mturk.com (accessed on 22 December 2022). In addition, there was no usability study in [13].



Figure 3. The sample size used in each study to date: [3–5,15,17,18,37,38].

5. Effectiveness of Adversarial CAPTCHA

The increased improvements in attack methods—particularly those using deep neural network to break CAPTCHAs—has led to the development of adversarial CAPTCHAs using various generation algorithms. This is due to the fact that adversarial examples pose a real threat to DL in practice. Moreover, the adversarial examples are often readable by humans. This improvement has significantly reduced the success rate of attackers, as demonstrated in [6].

However, filtering attacks may completely remove the adversarial noise in specific domains. For example, a study in [3] empirically demonstrated the possibility of removing adversarial noise that was generated using some adversarial CAPTCHA generation techniques. Moreover, a comprehensive study in [43] conducted a survey of adversarial attacks on DL in computer vision, and it listed a set of defence approaches that can detect adversarial perturbations.

Therefore, a secure CAPTCHA cannot be achieved using a plain adversarial noise. Based on this, developers must focus on the used perturbation algorithm in order to generate a high-quality adversarial CAPTCHA that is resistant to the removal of adversarial noise attacks.

6. Observation Based on the Analysed Studies

Based on the reviewed studies, some observations have been noted and summarized as follows:

- Although new adversarial CAPTCHA generation algorithms have recently been introduced, researchers have not evaluated the robustness of the generated samples against the new sophisticated Google OCR [44]. Thus, it might be interesting to evaluate the robustness of these samples. However, this is out out of the scope of our study and can be a future work.
- The sample size used for evaluating the usability aspects was generally too small. Therefore, it is recommended that the sample size should be taken into consideration when conducting a usability study for developing adversarial CAPTCHAs. Furthermore, since there have been few scientific studies supporting a systematic design or tuning for users, a parametric study may be an interesting focus of future research in analysing adversarial CAPTCHAs at the parameter level by accompanying an experimental study.
- CAPTCHA-solving services that employ human users (e.g., [45]) are still an attack to adversarial CAPTCHAs, since CAPTCHAs are designed to be recognized by humans. Despite the fact that a recent study in [2] proposed a CAPTCHA system against human-based attacks by exploiting the keystroke dynamics authentication system, this should be investigated in the context of adversarial CAPTCHA schemes.

- There have been few methodologies for generating adversarial CAPTCHAs. Hence, there are still various methods for generating adversarial examples that have not been used yet. These methods are highlighted in the following section.
- The DL algorithms that are most commonly used for breaking CAPTCHAs are the RNN and CNN. Thus, it may be interesting to evaluate other DL algorithms' performance in order to identify the best algorithm.
- The most used domain for adding perturbation is the space domain. However, it
 appears that this domain is more susceptible to perturbation removal methods. On
 the other hand, perturbations added in the frequency domain are difficult to remove.
 This may encourage more investigations to be conducted in terms of inconsistency of
 the transferability property among various ML models.
- We have observed the potential significance of introducing requirements for generating an adversarial CAPTCHA. Accordingly, the recommended requirements for generating adversarial CAPTCHAs are as follows:
 - Perturbation: The applied noise should be effective in misleading the targeted system at least 98.5% [3].
 - Security: The perturbation should be considered as a global change in order to be resistant to any removable algorithms such as filtering or ML algorithms.
 - Usability: The perturbation should not affect the readability of the generated CAPTCHA.
 - Scalability: The CAPTCHA generator should be computationally efficient in terms of generating thousands of CAPTCHAs per second.
 - Repeatability: The generated adversarial CAPTCHAs should not be repeated.
 - O Predictability: The generated adversarial CAPTCHAs should not be predicted.
- We have observed that most, if not all, attacks against adversarial CAPTCHAs are accomplished based on real CAPTCHA samples downloaded directly from the system. As such, it is highly recommended to investigate a new protection approach, not only for the CAPTCHA itself but also against downloading the generated sample.

7. Other Techniques for Generating Adversarial CAPTCHA

In this paper, we have discussed some techniques for generating adversarial CAPTCHAs in Section 4.3. However, there are still various techniques that have not been used yet. Thus, this section presents a set of techniques for generating adversarial examples that might be suitable for generating adversarial CAPTCHA.

7.1. Carlini and Wagner (CW)

This approach is a new technique based on the L-BFGS technique [46] to define the issue of finding adversarial samples. The main goal is to identify the smallest changes on the original data for the purpose of changing the classification.

7.2. Deep Fool

This technique proposed by [47] generates untargeted adversarial examples. This technique attempts to minimize the distance measure between changed samples and original samples.

7.3. Zeroth Order Optimization (ZOO)

This approach is proposed by [48] to estimate the gradient of the classifiers, which is an appropriate option for a black-box attack, without accessing the classifier for generating adversarial examples. The approach consists of applying noise data added to each feature of the original example in an iterative process and asking the classifier to measure the gradient of these different features.

7.4. Basic Iteration Method (BIM)

This approach is a type of FGSM method introduced by [49] to add a simple way to expand the "FGSM" strategy. This is implemented many times with a small step size.

7.5. Projected Gradient Descent (PGD)

This strategy is also a kind of FGSM technique proposed by [50], and it is used to find the perturbation that optimizes the loss function by keeping the perturbation small enough to be in the permitted range.

7.6. Particle Swarm Optimization (PSO)

This approach is a method of computing that optimizes a problem by iteratively attempting to improve solutions with respect to a given measure of quality [51].

7.7. Genetic Algorithm (GA)

This approach demonstrated by [52] relies on biologically motivated operators to produce excellent solutions to the optimization and search problems. Genetic algorithms usually operate on data structures defined as chromosomes, where a chromosome is a reflection of the problem data. A chromosome in our field is composed of a data feature from the original data sets.

8. Conclusions

This paper provides a comprehensive survey on adversarial perturbations and attacks in CAPTCHAs. This survey can help new researchers with up-to-date knowledge, current trends and field progress. Based on our analysis of current studies, the development of adversarial CAPTCHAs is a promising research line in creating more robust CAPTCHAs, especially for the text-based scheme and the impact of existing developed attacks. Moreover, we have indicated several observations for a broader outlook in this research direction. Finally, we suggest the investigation of the appropriateness of applying various perturbation techniques to adversarial CAPTCHAs.

Funding: This research received no external funding.

Data Availability Statement: The data presented in this study are available in the cited studies.

Acknowledgments: The researchers would like to thank the Deanship of Scientific Research, Qassim University for funding the publication of this project.

Conflicts of Interest: The author declare no conflict of interest.

References

- 1. Von Ahn, L.; Blum, M.; Langford, J. Telling humans and computers apart automatically. Commun. ACM 2004, 47, 56–60. [CrossRef]
- Alsuhibany, S.A.; Alreshoodi, L.A. Detecting human attacks on text-based CAPTCHAs using the keystroke dynamic approach. IET Inf. Secur. 2021, 15, 191–204. [CrossRef]
- Osadchy, M.; Hernandez-Castro, J.; Gibson, S.; Dunkelman, O.; Pérez-Cabo, D. No bot expects the DeepCAPTCHA! Introducing immutable adversarial examples with applications to CAPTCHA generation. *IEEE Trans. Inf. Forensics Secur.* 2017, 12, 2640–2653. [CrossRef]
- 4. Wang, S.; Zhao, G.; Liu, J. Text Captcha Defense Algorithm Based on Overall Adversarial Perturbations. In *Journal of Physics: Conference Series*; IOP Publishing: Bristol, UK, 2021; Volume 1744, p. 042243.
- Shi, C.; Xu, X.; Ji, S.; Bu, K.; Chen, J.; Beyah, R.; Wang, T. Adversarial captchas. *IEEE Trans. Cybern.* 2021, 52, 6095–6108. [CrossRef] [PubMed]
- Shi, C.; Ji, S.; Liu, Q.; Liu, C.; Chen, Y.; He, Y.; Liu, Z.; Beyah, R.; Wang, T. Text captcha is dead? A large scale deployment and empirical study. In Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, New York, NY, USA, 9–13 November 2020; pp. 1391–1406.
- Chen, J.; Luo, X.; Guo, Y.; Zhang, Y.; Gong, D. A Survey on Breaking Technique of Text-Based CAPTCHA. Secur. Commun. Netw. 2017, 2017, 6898617. [CrossRef]

- Ye, G.; Tang, Z.; Fang, D.; Zhu, Z.; Feng, Y.; Xu, P.; Chen, X.; Wang, Z. Yet another text captcha solver: A generative adversarial network based approach. In Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, New York, NY, USA, 15–19 October 2018; pp. 332–348.
- Li, C.; Chen, X.; Wang, H.; Wang, P.; Zhang, Y.; Wang, W. End-to-end attack on text-based CAPTCHAs based on cycle-consistent generative adversarial network. *Neurocomputing* 2021, 433, 223–236. [CrossRef]
- 10. Sivakorn, S.; Polakis, I.; Keromytis, A.D. I am robot: (Deep) learning to break semantic image CAPTCHAs. In Proceedings of the 2016 IEEE European Symposium on Security and Privacy (EuroS&P), Saarbrücken, Germany, 21–24 March 2016; pp. 388–403.
- George, D.; Lehrach, W.; Kansky, K.; Lázaro-Gredilla, M.; Laan, C.; Marthi, B.; Lou, X.; Meng, Z.; Liu, Y.; Wang, H.; et al. A generative vision model that trains with high data efficiency and breaks text-based CAPTCHAs. *Science* 2017, 358, eaag2612. [CrossRef]
- Na, D.; Park, N.; Ji, S.; Kim, J. CAPTCHAs Are Still in Danger: An Efficient Scheme to Bypass Adversarial CAPTCHAs. In Proceedings of the International Conference on Information Security Applications, Jeju Island, Republic of Korea, 26–28 August 2020; Springer: Cham, Switzerland, 2020; pp. 31–44.
- Zhang, Y.; Gao, H.; Pei, G.; Kang, S.; Zhou, X. Effect of adversarial examples on the robustness of CAPTCHA. In Proceedings of the International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC), Zhengzhou, China, 18–20 October 2018; IEEE: Piscataway, NJ, USA, 2018; pp. 1–109.
- 14. Ardhita, N.B.; Maulidevi, N.U. Robust Adversarial Example as Captcha Generator. In Proceedings of the 7th International Conference on Advance Informatics: Concepts, Theory and Applications (ICAICTA), Tokoname, Japan, 8–9 September 2020; IEEE: Piscataway, NJ, USA, 2020; pp. 1–4.
- Zhang, J.; Sang, J.; Xu, K.; Wu, S.; Zhao, X.; Sun, Y.; Hu, Y.; Yu, J. Robust CAPTCHAs Towards Malicious OCR. *IEEE Trans. Multimed.* 2020, 23, 2575–2587. [CrossRef]
- 16. Shao, R.; Shi, Z.; Yi, J.; Chen, P.-Y.; Hsieh, C.-J. Robust Text CAPTCHAs Using Adversarial Examples. arXiv 2021, arXiv:2101.02483.
- 17. Jia, X.; Xiao, J.; Wu, C. TICS: Text-image-based semantic CAPTCHA synthesis via multi-condition adversarial learning. *Vis. Comput.* 2021, 2021, 963–975. [CrossRef]
- Hitaj, D.; Hitaj, B.; Jajodia, S.; Mancini, L.V. Capture the Bot: Using Adversarial Examples to Improve CAPTCHA Robustness to Bot Attacks. *IEEE Intell. Syst.* 2021, 36, 104–112. [CrossRef]
- Bursztein, E.; Aigrain, J.; Moscicki, A.; Mitchell, J.C. The end is nigh: Generic solving of text-based CAPTCHAs. In Proceedings of the 8th USENIX Workshop on Offensive Technologies (WOOT), San Diego, CA, USA, 19 August 2014; pp. 1–15.
- 20. Ye, G.; Tang, Z.; Fang, D.; Zhu, Z.; Feng, Y.; Xu, P.; Chen, X.; Han, J.; Wang, Z. Using Generative Adversarial Networks to Break and Protect Text Captchas. *ACM Trans. Priv. Secur.* **2020**, *23*, 1–29. [CrossRef]
- Zhang, N.; Ebrahimi, M.; Li, W.; Chen, H. A generative adversarial learning framework for breaking text-based CAPTCHA in the dark web. In Proceedings of the International Conference on Intelligence and Security Informatics (ISI), Arlington, VA, USA, 9–10 November 2020; IEEE: Piscataway, NJ, USA, 2020; pp. 1–6.
- 22. Hussain, R.; Gao, H.; Shaikh, R.A. Segmentation of connected characters in text-based CAPTCHAs for intelligent character recognition. *Multimed. Tools Appl.* 2017, 76, 25547–25561. [CrossRef]
- TLe, A.; Baydin, A.G.; Zinkov, R.; Wood, F. Using synthetic data to train neural networks is model-based reasoning. In Proceedings of the International Joint Conference on Neural Networks (IJCNN), Anchorage, AK, USA, 14–19 May 2017; IEEE: Piscataway, NJ, USA, 2017; pp. 3514–3521.
- 24. Tang, M.; Gao, H.; Zhang, Y.; Liu, Y.; Zhang, P.; Wang, P. Research on Deep Learning Techniques in Breaking Text-Based Captchas and Designing Image-Based Captcha. *IEEE Trans. Inf. Forensics Secur.* **2018**, *13*, 2522–2537. [CrossRef]
- Zi, Y.; Gao, H.; Cheng, Z.; Liu, Y. An end-to-end attack on text captchas. *IEEE Trans. Inf. Forensics Secur.* 2019, 15, 753–766. [CrossRef]
- Wu, X.; Dai, S.; Guo, Y.; Fujita, H. A machine learning attack against variable-length Chinese character CAPTCHAs. *Appl. Intell.* 2019, 49, 1548–1565. [CrossRef]
- Liu, F.; Li, Z.; Li, X.; Lv, T. A text-based captcha cracking system with generative adversarial networks. In Proceedings of the 2018 IEEE International Symposium on Multimedia (ISM), Taichung, Taiwan, 10–12 December 2018; IEEE: Piscataway, NJ, USA, 2018; pp. 192–193.
- Tian, S.; Xiong, T. A generic solver combining unsupervised learning and representation learning for breaking text-based captchas. In Proceedings of the Web Conference 2020, New York, NY, USA, 20–24 April 2020; pp. 860–871.
- 29. Nouri, Z.; Rezaei, M. Deep-CAPTCHA: A deep learning based CAPTCHA solver for vulnerability assessment. *arXiv* 2020, arXiv:2006.08296. [CrossRef]
- UmaMaheswari, P.; Ezhilarasi, S.; Harish, P.; Gowrishankar, B.; Sanjiv, S. Designing a Text-based CAPTCHA Breaker and Solver by using Deep Learning Techniques. In Proceedings of the International Conference on Advances and Developments in Electrical and Electronics Engineering (ICADEE), Coimbatore, India, 10–11 December 2020; IEEE: Piscataway, NJ, USA, 2020; pp. 1–6.
- 31. Thobhani, A.; Gao, M.; Hawbani, A.; Ali, S.T.M.; Abdussalam, A. CAPTCHA Recognition Using Deep Learning with Attached Binary Images. *Electronics* **2020**, *9*, 1522. [CrossRef]
- 32. Wang, P.; Gao, H.; Shi, Z.; Yuan, Z.; Hu, J. Simple and Easy: Transfer Learning-Based Attacks to Text CAPTCHA. *IEEE Access* 2020, *8*, 59044–59058. [CrossRef]

- 33. Che, A.; Liu, Y.; Xiao, H.; Wang, H.; Zhang, K.; Dai, H.-N. Augmented Data Selector to Initiate Text-Based CAPTCHA Attack. *Secur. Commun. Networks* 2021, 2021, 9930608. [CrossRef]
- 34. Bostik, O.; Horak, K.; Kratochvila, L.; Zemcik, T.; Bilik, S. Semi-supervised deep learning approach to break common CAPTCHAs. *Neural Comput. Appl.* **2021**, *33*, 13333–13343. [CrossRef]
- 35. Wang, Y.; Wei, Y.; Zhang, M.; Liu, Y.; Wang, B. Make complex CAPTCHAs simple: A fast text captcha solver based on a small number of samples. *Inf. Sci.* **2021**, *578*, 181–194. [CrossRef]
- Shekhar, H.; Moh, M.; Moh, T.S. Exploring adversaries to defend audio captcha. In Proceedings of the 2019 18th IEEE International Conference on Machine Learning and Applications (ICMLA), Boca Raton, FL, USA, 16–19 December 2019; IEEE: Piscataway, NJ, USA, 2019; pp. 1155–1161.
- Hossen, I.; Hei, X. AAECAPTCHA: The design and implementation of audio adversarial captcha. In Proceedings of the 2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P), Genoa, Italy, 6–10 June 2022; IEEE: Piscataway, NJ, USA, 2022; pp. 430–447.
- Abdullah, H.; Karlekar, A.; Prasad, S.; Rahman, M.S.; Blue, L.; Bauer, L.A.; Traynor, P. Attacks as Defenses: Designing Robust Audio CAPTCHAs Using Attacks on Automatic Speech Recognition Systems. *arXiv* 2022, arXiv:2203.05408.
- 39. Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; Fergus, R. Intriguing properties of neural networks. *arXiv* **2013**, arXiv:1312.6199.
- Fawzi, A.; Fawzi, O.; Frossard, P. Analysis of classifiers' robustness to adversarial perturbations. *Mach. Learn.* 2018, 107, 481–508. [CrossRef]
- Terada, T.; Nguyen VN, K.; Nishigaki, M.; Ohki, T. Improving Robustness and Visibility of Adversarial CAPTCHA Using Low-Frequency Perturbation. In Proceedings of the International Conference on Advanced Information Networking and Applications, Sydney, Australia, 13–15 April 2022; Springer: Cham, Switzerland, 2022; pp. 586–597.
- 42. Goodfellow, I.J.; Shlens, J.; Szegedy, C. Explaining and harnessing adversarial examples. arXiv 2014, arXiv:1412.6572.
- Akhtar, N.; Mian, A. Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access* 2018, 6, 14410–14430. [CrossRef]
- 44. Google Cloud APIs. Available online: https://cloud.google.com/apis/docs/overview (accessed on 12 September 2022).
- 45. CAPTCHA Solving Service. Available online: https://anti-captcha.com/ (accessed on 15 September 2022).
- 46. Carlini, N.; Wagner, D. Towards evaluating the robustness of neural networks. In Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP), San Jose, CA, USA, 22–26 May 2017; IEEE: Piscataway, NJ, USA, 2017; pp. 39–57.
- Moosavi-Dezfooli, S.M.; Fawzi, A.; Frossard, P. Deepfool: A simple and accurate method to fool deep neural networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 2574–2582.
- Chen, P.Y.; Zhang, H.; Sharma, Y.; Yi, J.; Hsieh, C.J. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, Dallas, TX, USA, 3 November 2017; pp. 15–26.
- Kurakin, A.; Goodfellow, I.J.; Bengio, S. Adversarial examples in the physical world. In Proceedings of the 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, 24–26 April 2017; Workshop Track Proceedings. pp. 1–14. [CrossRef]
- 50. Makelov, A.; Schmidt, L.; Tsipras, D.; Vladu, A. Towards Deep Learning Models Resistant to Adversarial Attacks. *arXiv* 2017, arXiv:1706.06083.
- 51. Kennedy, J.; Eberhart, R. Particle swarm optimization PAPER—IGNORE FROM REFS. In Proceedings of the IEEE International Conference on Neural Networks, Perth, Australia, 27 November–1 December 1995; pp. 1942–1948.
- 52. Whitley, D. A Genetic Algorithm Tutorial. Stat. Comput. 1994, 4, 65–85. [CrossRef]

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.