

# FAWA: Fast Adversarial Watermark Attack on Optical Character Recognition (OCR) Systems

Lu Chen<sup>1</sup>, Jiao Sun<sup>2</sup> and Wei Xu<sup>1</sup>(✉)

<sup>1</sup> Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing 100084, China

lchen17@mails.tsinghua.edu.cn, weixu@mail.tsinghua.edu.cn

<sup>2</sup> CS, University of Southern California, Los Angeles, CA 90007, USA  
jiaosun@usc.edu

**Abstract.** Deep neural networks (DNNs) significantly improved the accuracy of optical character recognition (OCR) and inspired many important applications. Unfortunately, OCRs also inherit the vulnerabilities of DNNs under adversarial examples. Different from colorful vanilla images, text images usually have clear backgrounds. Adversarial examples generated by most existing adversarial attacks are unnatural and pollute the background severely. To address this issue, we propose the *Fast Adversarial Watermark Attack (FAWA)* against sequence-based OCR models in the white-box manner. By disguising the perturbations as watermarks, we can make the resulting adversarial images appear natural to human eyes and achieve a perfect attack success rate. FAWA works with either gradient-based or optimization-based perturbation generation. In both letter-level and word-level attacks, our experiments show that in addition to natural appearance, FAWA achieves a 100% attack success rate with 60% less perturbations and 78% fewer iterations on average. In addition, we further extend FAWA to support full-color watermarks, other languages, and even the OCR accuracy-enhancing mechanism.

**Keywords:** Watermark; OCR Model; Targeted White-Box Attack

## 1 Introduction

Optical Character Recognition (OCR) has been an important component in text processing applications, such as license-plate recognition, street sign recognition and financial data analysis. Deep neural networks (DNNs) significantly improve OCR’s performance. Unfortunately, OCR inherits all counter-intuitive security problems of DNNs. OCR models are also vulnerable to *adversarial examples*, which are crafted by making human-imperceptible perturbations on original images to mislead the models. The wide application of OCR provides incentives for adversaries to attack OCR models, thus damaging downstream applications, resulting in fake ID information, incorrect metrics readings and financial data.

Prior works have shown that we can change the prediction of DNNs in image classification tasks only by applying carefully-designed perturbations [8, 19, 20, 22],

or adding a small patch [3,13] to original images. But these methods are not directly applicable to OCR attacks. 1) Most text images are on white backgrounds. Perturbations added by existing methods appear too evident for human eyes to hide, and pollute clean backgrounds. 2) Instead of classifying characters individually, modern OCR models are segmentation-free, inputting entire variable-sized image and outputting a sequence of labels. It is called the *sequential labeling* task [9]. Since modern OCR models use the CNN [15] and the LSTM [12] as feature extractors, the internal feature representation also relies on contexts (i.e., nearby characters). Therefore, besides perturbing the target character, we also need to perturb its context, resulting in more perturbations. 3) Since text images are usually dense, there is no open space to add a patch.

In this paper, we propose a new attack method, *Fast Adversarial Watermark Attack* (FAWA), against modern OCR models. Watermarks are images or patterns commonly used in documents as backgrounds to identify things, such as marking a document proprietary, urgent, or merely for decoration. Human readers are so used to the watermark that they naturally ignore it. We generate natural watermark-style adversarial perturbations in text images. Such images appear natural to human eyes but misguide OCR models. Watermark perturbations are similar to patch-based attacks [3,13] where perturbations are confined to a region. Different from patches occupying a separate region, watermarks overlay on texts but not hinder the text readability. Laplace attack [10] and HAAM [11] try to generate seemingly smooth perturbations, which work well on colored photos. However, they do not solve the problem of background pollution.

FAWA is a *white-box* and *targeted* attack. We assume adversaries have perfect knowledge of the DNN architecture and parameters (white-box model) and generate specific recognition results (targeted). There are three steps in the perturbation generation. 1) We find a good position over the target character to add an initial watermark. 2) We generate perturbations inside the watermark. 3) Optionally, we convert the gray watermark into a colored one. To generate perturbations, we leverage either gradient-based or optimization-based methods. We evaluate FAWA on a state-of-the-art open-source OCR model, Calamari-OCR [26] for English texts with five fonts. FAWA generates adversarial images with 60% less perturbations and 78% fewer iterations than existing methods, while maintaining a 100% attack success rate. Our adversarial images also look quite similar to natural watermarked images. We evaluate the effects of colored watermarks and other languages under real-world application settings. Last, we propose a positive application of the FAWA, i.e., using the perturbations to enhance the accuracy of OCR models. The contributions of this paper include:

1. We propose an attack method, FAWA, of hiding adversarial perturbations in watermarks from human eyes. We implemented FAWA as the efficient adversarial attacks against the DNN-based OCR in sequential labeling tasks;
2. Extensive experiments show that FAWA performs targeted attacks perfectly, and generate natural watermarked images with imperceptible perturbations;
3. We demonstrate several applications of FAWA, such as colored watermarks as an attack mechanism, using FAWA as an accuracy-enhancing mechanism.

## 2 Background and Related Work

**Optical Character Recognition (OCR).** We can roughly categorize existing OCR models into character-based models and end-to-end models. *Character-based recognition models* segment the image into character images first, before passing them into the recognition engine. Obviously, the OCR performance heavily relies on the character segmentation. *End-to-end recognition models* apply an unsegmented recognition engine, which recognizes the entire sequence of characters in a variable-sized image. [1,7] adopt sequential models such as Markov models and [2,24] utilize DNNs as feature extractors for sequential recognition. [9] introduces a segmentation-free approach, connectionist temporal classification (CTC), which provides a sort of loss function of enabling us to recognize variable-length sequences without explicit segmentation while training DNN models. Thus, many state-of-the-art OCR models use CTC as the loss function.

**Attacking DNN-based Computer Vision Models.** Attacking DNN-based models is a popular topic in both computer vision and security fields. Existing attack methods use the following two ways to generate perturbations. 1) Making perturbations small enough for evading human perception. For example, many projects aim at generating minor  $L_p$ -norm perturbations for the purpose of making them visually imperceptible. FGSM [8], L-BFGS [22], DeepFool [18], C&W  $L_2/L_\infty$  [4], PGD [17] and EAD [5], all perform modifications at the pixel level by a small amount bounded by  $\epsilon$ . 2) Making perturbations in a small region of an image. For example, JSMA [20], C&W  $L_0$  [4], Adversarial Patch [3] and LaVAN [13], all perturb a small region of pixels in an image, but their perturbations are not limited by  $\epsilon$  at the pixel level. Though FAWA is similar to patches, in the text images, watermarks overlay on the text instead of covering the text.

**Generating minimal adversarial perturbations.** FAWA generates adversarial perturbations using the either gradient-based or optimization-based methods.

Gradient-based methods add perturbations generated from gradient against input pixels. *FGSM* [8] is a  $L_\infty$ -norm one-step attack. It is efficient but only provides a coarse approximation of the optimal perturbations. *BIM* [16] takes multiple smaller steps and the resulting image is clipped by the same bound  $\epsilon$ . Thus BIM produces superior results to FGSM. *MI-FGSM* [6] extends BIM with momentum. MI-FGSM can not only stabilize update directions but also escape from poor local maxima during iterations, and thus generates more transferable adversarial examples. Considering the efficiency, we adopt MI-FGSM in FAWA.

Optimization-based methods directly solve the box-constrained optimization problem to minimize the  $L_p$ -norm distance between the original and adversarial image, while yielding the targeted result. *Box-constraint L-BFGS* [22] seeks adversarial examples with L-BFGS. Though L-BFGS constructs subtle perturbations, it is inefficient. Instead of applying cross-entropy as loss function, *C&W  $L_2$  attack* [4] constructs a new loss function and solves it with gradient descent. *OCR attack* [21] is the only available work of attacking sequence-based OCRs as far as we know, which generates adversarial examples using the CTC loss for sequential labeling tasks. In FAWA, we use the same setting as OCR attack.

**Perturbations with other optimization goals.** Besides minimizing the perturbation level, many works make efforts to produce smooth and natural perturbations. Laplace attack [10] smooths perturbations relying on Laplacian smoothing. HAAM [11] creates edge-free perturbations using harmonic functions for disguising natural shadows or lighting. However, the smoothing and disguising are for photos but not for text images. Instead of manipulating pixel values directly, [27] produces perceptually realistic perturbations with spatial transformation. Though avoiding background pollution, it cannot guarantee the readability of text when the attack needs large deformation. [23] also utilizes the watermark idea but performs attacks only by scaling and rotating plain watermarks. Without adding pixel-level perturbations, it does not offer a high attack success rate.

### 3 Fast Adversarial Watermark Attack

In this section, we introduce *fast adversarial watermark attacks* (FAWA). FAWA consists of three steps. 1) We automatically determine a good position to add the initial watermark in the text image so that we can confine the perturbation generation to that region. 2) We generate the watermark-style perturbations with either the *gradient-based method* or *optimization-based method*. 3) Optionally, we convert gray watermarks into full-color ones to improve the text readability.

#### 3.1 Preliminaries

**Problem definition of adversarial image generation.** Given a text image  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$  for any  $x_i \in [0, 1]$ , where  $n$  is the number of pixels in the image, our goal is to generate an adversarial example  $\mathbf{x}'$  with minimum  $L_p$ -norm perturbations  $\|\mathbf{x}' - \mathbf{x}\|_p$  against the white-box model  $f$  with intent to trick model  $f$  into outputting the targeted result  $\mathbf{t}$ ,  $f(\mathbf{x}') = \mathbf{t}$ . Formally, the problem of adversarial image generation is  $\min_{\mathbf{x}'} \|\mathbf{x}' - \mathbf{x}\|_p$  s.t.  $f(\mathbf{x}') = \mathbf{t}$  and  $\mathbf{x}' \in [0, 1]^n$ . Also, we define the *CTC loss function* with respect to the image  $\mathbf{x}$  as  $\ell_{\text{CTC}}(\mathbf{x}, \mathbf{t})$  and the target labels  $\mathbf{t} = [t_1, t_2, \dots, t_N]$  for  $t_i \in \mathcal{T}$ , where  $\mathcal{T}$  is the character set.

**Saliency map.** The saliency map [20] is a versatile tool that not only provides us valuable information to cause the targeted misclassification in the threat model but also assists us in intuitively explaining some attack behaviors as a visualization tool. The saliency map indicates the output sensitivity, relevant to the adversarial targets, to the input features. In the saliency map, a larger value indicates a higher likelihood of fooling the model  $f$  to output the target  $\mathbf{t}$  instead of the ground-truth  $f(\mathbf{x})$ . We can construct a saliency map of an image using the forward derivative with respect to each input component in the text image.

**$L_p$ -norms.**  $L_p$ -norms are commonly-used metrics to measure the perceptual similarity between the clean image  $\mathbf{x}$  and the adversarial image  $\mathbf{x}'$ , denoted by  $\|\mathbf{x} - \mathbf{x}'\|_p = (\sum_{i=1}^n |\mathbf{x} - \mathbf{x}'|_i^p)^{\frac{1}{p}}$ ,  $p = 2, \infty$ . In gradient-based methods, we apply  $L_p$ -norms to prune the saliency map for generating  $L_p$ -norm perturbations. In optimization-based methods,  $L_p$ -norms are usually as an optimization term in

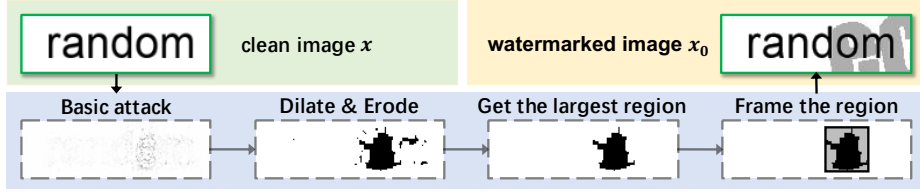


Fig. 1. Find the position of watermarks.

the objective function. Particularly,  $L_2$ -norm is greatly useful to enhance the visual quality.  $L_\infty$ -norm is to measure the maximum variation of perturbations.

### 3.2 Finding the Position of Watermarks

To automatically determine a good position of watermarks, as shown in Fig. 1, we perform the following steps. 1) We produce adversarial perturbations of the basic attack (i.e., Grad-Basic or Opt-Basic in Section 3.3) in the text image. 2) We binarize such adversarial image, and get its perturbed regions  $\mathbf{r} = [r_1, r_2, \dots, r_n]^T$ , where  $r_i = 1$  if  $x_i \neq x'_i$  and 0 otherwise. 3) In order to find relatively complete perturbed regions  $\mathbf{r}'$ , we apply a combination of erosion  $\oplus$  and dilation  $\ominus$  operations twice in the perturbed regions  $\mathbf{r}$ ,  $\mathbf{r}' = ((\mathbf{r} \oplus \mathcal{K}) \ominus \mathcal{K})^2$ , where we set the kernel  $\mathcal{K} = \mathbf{1}_{3 \times 3}$ . 4) After sorting the perturbed regions in  $\mathbf{r}'$  by their area, we obtain the largest perturbed region in  $\mathbf{r}'$ . We can find that our target texts locate in the same position as the largest found region. 5) Last, in the text image, we place a watermark big enough to cover the found position, and obtain an initial watermarked image  $\mathbf{x}_0$  and a binary watermark mask  $\Omega_w$  with the same shape of  $\mathbf{x}$ , where  $\Omega_{w,i} = 1$  if the position  $i$  is inside the watermark and 0 otherwise.

### 3.3 Generating Adversarial Watermarks

After finding the position of watermarks, we need to generate the perturbations within the watermarks to mislead the OCR models to output target texts. Integrated with two popular methods, we use the following methods to attack.

**Gradient-based Watermark Attack (Grad-WM).** Considering the high efficiency of MI-FGSM [6], we apply it as our gradient-based method (we refer to as *Grad-Basic*). Each iterative update of the Grad-Basic is to 1) get the saliency map normalized by  $L_1$ -norm with the cross-entropy loss, 2) adjust the update direction in the saliency map and update the momentum, and 3) update a new  $L_p$ -norm adversarial image  $\mathbf{x}'_{i+1}$  bounded by  $\epsilon$  using the updated momentum.

There are two main differences between Grad-WM and Grad-Basic. 1) Different from off-the-shelf MI-FGSM, where it applies the cross-entropy as loss function in the image classification tasks, in our Grad-WM, we use the CTC loss function to compute the saliency map. CTC loss function fits better because it is widely used in OCR models to handle the sequential labeling tasks. 2) In Grad-WM, to hide perturbations in the watermarks, we confine the perturbations

**Algorithm 1** Gradient-based Watermark Attack

**Input:** A text image  $\mathbf{x}$ , OCR model  $f$  with CTC loss  $\ell_{\text{CTC}}$ , targeted text  $\mathbf{t}$ ,  $\epsilon$ -bounded perturbation, attack step size  $\alpha$ , # of maximum iterations  $I$ , decay factor  $\mu$ .

**Output:** An adversarial example  $\mathbf{x}'$  with  $\|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon$  or attack failure  $\perp$ .

1: Initialization:  $\mathbf{g}_0 = \mathbf{0}$ ;  $\mathbf{x}'_0 = \mathbf{x}_0$

2: **for all** each iteration  $i = 0$  to  $I - 1$  **do**

3:   Input  $\mathbf{x}'_i$  to  $f$  and obtain the saliency map  $\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})$

4:   Update  $\mathbf{g}_{i+1}$  by accumulating  $\mathbf{g}_i$  in the saliency map normalized by L<sub>1</sub>-norm as

$$\mathbf{g}_{i+1} = \mu \cdot \mathbf{g}_i + \frac{\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})}{\|\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})\|_1} \quad (1)$$

5:   Update  $\mathbf{x}'_{i+1}$  by applying watermark-bounded L<sub>p</sub>-norm perturbations as

$$\mathbf{x}'_{i+1} = \mathbf{x}'_i + \text{clip}_{\epsilon}(\alpha \cdot (\Omega_w \odot \frac{\mathbf{g}_{i+1}}{\|\mathbf{g}_{i+1}\|_p})) \quad (2)$$

6:   **if**  $f(\mathbf{x}'_{i+1}) == \mathbf{t}$  **then**

7:     **return**  $\mathbf{x}'_{i+1}$

8:   **end if**

9: **end for**

10: **return** attack failure  $\perp$

within the boundary of the watermark  $\Omega_w$ , rather than spreading perturbations over the entire image in the result of the background pollution like Grad-Basic.

Algorithm 1 summarizes the Grad-WM generation procedure. Notably, input images are the initial watermarked images  $\mathbf{x}_0$  created in Section 3.2. Then the saliency map is produced with the CTC loss function relevant to the adversarial targets  $\mathbf{t}$ . Through element-wisely multiplying the watermark mask  $\Omega_w$  with the updated L<sub>p</sub>-norm  $\mathbf{g}_{i+1}$ , we get rid of the perturbations outside the watermark to maintain the clean background. We gain a significant visual improvement in the Grad-WM than the Grad-Basic. Besides, for further improving the attack efficiency, we determine whether to stop the attacks in every few iterations.

**Optimization-based Watermark Attack (Opt-WM).** We employ OCR attack [21] as our basic optimization-based method (Opt-Basic). Opt-Basic solves the following optimization problem:  $\min_{\mathbf{x}'} c \cdot \ell_{\text{CTC}}(\mathbf{x}', \mathbf{t}) + \|\mathbf{x}' - \mathbf{x}\|_2^2$  s.t.  $\mathbf{x}' \in [0, 1]^n$ , where  $c$  is a hyper-parameter. To improve the visual quality of adversarial images, Opt-Basic adopts L<sub>2</sub>-norm for penalizing perturbations. To eliminate the box-constraint  $[0, 1]$  of  $\mathbf{x}'$ , it reformulates the problem using the change of variables [4] as  $\min_{\boldsymbol{\omega}} c \cdot \ell_{\text{CTC}}(\frac{\tanh(\boldsymbol{\omega})+1}{2}, \mathbf{t}) + \|\frac{\tanh(\boldsymbol{\omega})+1}{2} - \mathbf{x}\|_2^2$ , which optimizes a new variable  $\boldsymbol{\omega}$  instead of the box-constrained  $\mathbf{x}'$ . The fact  $-1 \leq \tanh(\cdot) \leq 1$  implies that  $\frac{\tanh(\cdot)+1}{2}$  satisfies the box-constraint  $[0, 1]$  automatically. Intuitively, we treat the variable  $\boldsymbol{\omega}$  as the adversarial example  $\mathbf{x}'$  without the box-constraint.

Similar to Grad-WM, we first get the initial watermarked input image  $\mathbf{x}_0$  and the watermark mask  $\Omega_w$ . To confine perturbations in watermarks more conveniently, we separate the perturbation term from  $\mathbf{w}$ , and rewrite the original

$\mathbf{w}$  to  $\mathbf{w} + \mathbf{x}_0$  that represents the adversarial image of combining the perturbation term  $\mathbf{w}$  and the initial watermarked image term  $\mathbf{x}_0$ . The objective function changes to  $\min_{\mathbf{w}} c \cdot \ell_{\text{CTC}}(\frac{\tanh(\mathbf{w} + \mathbf{x}_0) + 1}{2}, \mathbf{t}) + \|\frac{\tanh(\mathbf{w} + \mathbf{x}_0) + 1}{2} - \mathbf{x}_0\|_2^2$ . To constrain the manipulation region in the watermark fashion, we perform the element-wisely multiplication between the perturbation variable  $\mathbf{w}$  and the watermark mask  $\Omega_w$ . Formally, we reformulate the optimization problem as

$$\min_{\mathbf{w}} c \cdot \ell_{\text{CTC}}(\frac{\tanh(\Omega_w \odot \mathbf{w} + \mathbf{x}) + 1}{2}, \mathbf{t}) + \|\frac{\tanh(\Omega_w \odot \mathbf{w} + \mathbf{x}) + 1}{2} - \mathbf{x}\|_2^2. \quad (3)$$

We adopt the Adam optimizer [14] for both Opt-Basic and Opt-WM to seek the watermark-style adversarial images. We utilize the binary search to adapt the tradeoff hyper-parameter  $c$  between the loss function and the  $L_2$ -norm distance.

### 3.4 Improving Efficiency

To improve attack efficiency, we employ batch attack and early stopping by increasing the number of parallel attacks and reducing redundant attack iterations.

**Batch attack.** Because OCR models only support batch image processing with the same size, when we attack a large number of variable-size text images, it will be time-consuming to attack them one by one. To improve attack parallelism (i.e., attack multiple images simultaneously), we 1) resize the variable-size images into the same height, 2) pad them into the maximum width among these images, 3) put the same-size images into a matrix. After the preprocessing, we could perform batch adversarial attacks in a single matrix to improve attack efficiency.

**Early stopping.** For avoiding redundant attack iterations, we insert the early-stop mechanism in the attack iterations. We evaluate attack success rate (ASR) every few iterations. We stop attack iteration once ASR achieves 100%, indicating that all adversarial images are misidentified as the targeted texts by the threat model. To a certain extent, early stopping reduces the attack efficiency as it takes time to check attack status during iterations. This is a basic tradeoff between the cycle of ASR evaluation and the maximum attack iteration setting.

### 3.5 Improving Readability

Sometimes the gray watermarks still hinder the readability of the text as it reduces the contrast. To achieve better readability of the text, we employ the following two effective strategies: adding text masks and full-color watermark.

**Adding Text Masks.** When generating the initial watermarked image  $\mathbf{x}_0$  in Section 3.2, to prevent the text from being obscured by the watermarks, we add in watermarks outside the text. To achieve this, we binarize the text images  $\mathbf{x}$  by a threshold  $\tau$  to get text masks  $\Omega_t$  with the same size of  $\mathbf{x}$ , where  $\Omega_{t,i} = 1$  if  $x_i > \tau$  and 0 otherwise. Then we get the initial watermarked images  $\mathbf{x}_0 = \mathbf{x} \odot (1 - \Omega_w \odot \overline{\Omega_t}) + \beta \cdot \Omega_w \odot \overline{\Omega_t}$ , where  $\beta$  is the grayscale value of initial watermarks.

**Full-Color Watermark.** Modern OCR systems always preprocess colored images into grayscale ones before recognition. Compared to grayscale watermarks,

colored watermarks have better visual quality on the black texts. To convert grayscale watermarks into RGB ones, we only manipulate the pixel  $x_i$  inside the watermark, where  $\Omega_{w,i} = 1$ . Given a grayscale value  $Gray$  in the grayscale watermark, we preset  $R$  value and  $B$  value to certain values. Next we make the color conversion with the transform equation:  $Gray = R*0.299 + G*0.587 + B*0.114$ , to calculate the left  $G$  value. Last we get RGB values from a grayscale value.

## 4 Experiments

### 4.1 Experiment Setup

**OCR model.** We choose to attack Calamari-OCR<sup>1</sup> [26]. The OCR model has two convolutional layers, two pooling layers, and the following LSTM layer. We use the off-the-shelf English model as our threat model, trained with CTC loss.

**Generate text images.** After processing the corpus in the IMDB dataset, we generate a dataset with 97 paragraph, 1479 sentence and 1092 word images using the Text Recognition Data Generator<sup>2</sup>. We use five fonts to generate these images: Courier, Georgia, Helvetica, Times and Arial. We set the font size to 32 pixels. We verify the Calamari-OCR can achieve 100% accuracy on these images.

**Generating letter-level attacks targets.** We choose a target word that is 1) a valid word, and 2) with edit distance 1 from the original. We evaluate replacement, insertion and deletion. More similar the replacement target is to the original, the easier the attack is. For example, replacing letter **t** with letter **f** is easier than replacing **t** with letter **j**. We use the *logit* value from the output of the last hidden layer, as the similarity metric between a pair of letters. Given the logit value, we assign replacement attacks into easy, random and hard case.

**Generating word-level attacks targets.** In this task, we replace the entire word in word, sentence and paragraph images. Different from letter-level attacks, we randomly choose a word in the English dictionary which has the same length as the original. However, we don’t constrain the edit distance between them.

**Attack implementation and settings.** We implemented Grad-WM and Opt-WM attacks<sup>3</sup>. We normalize all input images to  $[0, 1]$  and set  $\beta$ , the grayscale value of initial watermarks, to 0.682. Given color transform equation, fixing  $R = 255$  and  $B = 0$ , RGB’s upper bound (255, 255, 0) is equal to  $Gray = 0.882$ . After adding  $\epsilon$ -bounded perturbations ( $\epsilon = 0.2$ ), watermarks are still less than 0.882, the upper bound for valid conversion. We use the watermark style: the word “ecml” of Impact font with 78-pixel font size and 15-degree rotation.

For the gradient-based methods, we use the implementation of MI-FGSM in CleverHans Python library<sup>4</sup>. We set maximum iterations  $I = 2000$ , batch size = 100 and  $\epsilon = 0.2$ . For  $L_2$ -norm,  $\alpha = 0.05$ ; for  $L_\infty$ -norm,  $\alpha = 0.01$ . The momentum

<sup>1</sup> <https://github.com/Calamari-OCR/calamari>

<sup>2</sup> <https://github.com/Belval/TextRecognitionDataGenerator>

<sup>3</sup> <https://github.com/strongman1995/Fast-Adversarial-Watermark-Attack-on-OCR>

<sup>4</sup> <https://github.com/tensorflow/cleverhans>



	original $\mathbf{x}$	WM <sub>0</sub> $\mathbf{x}'_0$	Grad-Basic	Grad-WM	Opt-Basic	Opt-WM	Color-WM	target
letter	year	year	year	year	year	year	year	hear
	here	here	here	here	here	here	here	there
	short	short	short	short	short	short	short	sort
word	tennis	tennis	tennis	tennis	tennis	tennis	tennis	amazon
	entry	entry	entry	entry	entry	entry	entry	waken
	hoy	hoy	hoy	hoy	hoy	hoy	hoy	jow

**Table 1.** Adversarial examples of letter-level and word-level attacks. Colored watermark examples are converted from Grad-WM examples.

decay factor  $\mu$  is 1.0. For the optimization-based methods, we use the Adam optimizer [14] for 1000 steps with mini-batch size 100 and learning rate 0.01. We choose  $c = 10$  as the tradeoff between the targeted loss and perturbation level.

**Evaluation metrics.** We evaluate attack capability from the following aspects.
















- 1) Perturbation level.** We quantify the perturbation level with three metrics:  $MSE = \frac{1}{n} \|\mathbf{x} - \mathbf{x}'\|_2^2$ , evaluates the difference between two images.  $PSNR = 10 \log \left( \frac{D^2}{MSE} \right)$ , where  $D$  is the range of pixel intensities.  $SSIM$  [25] captures structural information and measures the similarity between two images. For these metrics, we take the average of all images in the following results. Smaller MSE, larger PSNR and SSIM closer to 1 indicate less perturbations.
- 2) Success rate.** Targeted attack success rate,  $ASR = \frac{\#(f(\mathbf{x}')=\mathbf{t})}{\#(\mathbf{x})}$ , is the proportion of adversarial images of fooling OCR models to output targeted text.
- 3) Attack efficiency.**  $I_{avg}$  is average iterations of images to reach 100% ASR.

## 4.2 Letter-level Attack Performance

Table 1 illustrates generated adversarial examples. We find that the perturbations of basic attacks are quite obvious, especially when we want to replace the entire word, while the watermarks help hide the perturbations from human eyes.

To intuitively analyze the effects of watermarks on our attacks, in Fig. 2, we first illustrate the saliency map that highlights the influence of each pixel on the target output  $\mathbf{t}$ . In replacement cases, we observe that the white background needs more perturbations, both positive (+) and negative (-), while the gray background requires 72% less perturbations due to reduced contrast. Watermarks approximate the gray background and look more natural. Specifically, the attack of the white background also adds significant a-shaped negative perturbations to weaken the letter **a**, while the other two require negative perturbations less than 50% due to lower contrast. In addition, the attacks add perturbations to neighboring letters, as the sequence-based OCR models also consider them.

Table 2 and Fig. 3 show the quantitative analysis of attack performance and use different targets in letter-level attacks. In Fig. 3, a sharper slope indicates

	clean	adversarial	saliency map	MSE	saliency map+	MSE+	saliency map-	MSE-	target
replacement	parts	parts		55.82		34.77		21.05	<u>port</u>
	parts	parts		15.47		6.92		8.55	<u>port</u>
	parts	parts		24.97		14.13		10.84	<u>port</u>
+	parts	parts		14.73		7.10		7.63	part <u>is</u>
-	parts	parts		4.31		2.42		1.88	pars

**Fig. 2.** Saliency map visualization. First three lines are replacement. They have clean, gray and watermark backgrounds, respectively. Last two lines are insertion(+) and deletion(-). We fetch the positive part of saliency map to generate saliency map+ and the negative part to generate saliency map-. MSE, MSE+, MSE- represent the perturbation level in the saliency map, saliency map+ and saliency map-, respectively.

a higher efficiency, that reaches a higher ASR in a fixed number of iterations. ASRs of all attacks are 100%. We observe that: 1) MSEs of watermark attacks are only 40% of basic attacks' on average, confirming to the intuitive analysis above. 2) Watermark attacks only require 78% fewer  $I_{avg}$  on average, and have around 3 to 8 times sharper slope than basic attacks, showing the significant improvement of attack efficiency. 3) Not surprisingly, hard cases require both more perturbations (higher MSE) and more iterations (larger  $I_{avg}$ ) with lower efficiency (flatter slope), due to lack of similarity between two letters. 4) In addition, deletions are easier than insertions and replacements. This is because OCR models are sensitive to perturbations, and classify fuzzy letters into blank tokens that will be ignored by CTC in the output. Intuitively, in the saliency map of Fig. 2, perturbations of deletion are much slighter than other cases. Thus few perturbations can achieve deletion. 5) Courier font is easier to attack because it is thinner than other fonts. Thus it requires less perturbations. 6) Comparing gradient-based and optimization-based methods, the perturbation level is higher in gradient-based methods, no matter if they have watermarks or not. Because perturbations of optimization-based methods are not constrained by  $\epsilon$ .

### 4.3 Word-Level Attack Performance

We can still achieve 100% ASR in the word-level attacks, but both MSE and  $I_{avg}$  of word-level attacks are significantly higher than those of letter-level attacks.

We perform word-level attacks in word, sentence and paragraph images. Table 3 and Fig. 3(c) summarize results in word-level attacks. Due to limited space, we only show results of gradient-based methods. The optimization-based methods have similar insights. In word-level attacks, we have similar observations as letter-level attacks. 1) Watermark attacks spend 50% less  $I_{avg}$  with higher efficiency (sharper slope) generating adversarial images than basic attacks. 2) We observe 56% lower MSE with watermarks averagely. 3) It is easier to attack

**Table 2.** Letter-level attacks using Grad-Basic, Grad-WM, Opt-Basic, Opt-WM attacks. Last line is the target output. We denote each font with their first letter.

		Gradient-based					Optimization-based					
		replacement			insertion	deletion	replacement			insertion	deletion	
		easy	random	hard			easy	random	hard			
Basic	C	10.5 59	14.0 74	17.0 70	11.6 50	3.2 21	25.4 266	30.3 313	36.7 321	25.4 309	13.6 43	
	G	27.4 43	32.8 99	37.3 104	22.1 83	17.3 55	52.0 292	59.4 318	67.5 328	41.6 337	45.0 169	
	H	27.0 51	33.6 113	38.6 113	23.0 70	16.7 43	52.1 301	60.2 328	68.2 340	47.1 321	45.0 178	
	T	26.4 62	31.5 85	35.8 109	20.3 98	17.2 68	49.9 294	56.1 324	61.6 345	41.7 314	44.3 172	
	A	29.8 51	36.7 73	42.5 66	24.3 88	19.2 59	56.3 304	65.3 327	73.8 341	48.3 324	51.0 176	
		parts	parts	parts	parts	parts	parts	parts	parts	parts	parts	
WM	C	2.8 30	3.6 18	4.3 27	3.6 21	0.7 8	16.7 116	20.1 96	20.4 95	31.1 29	3.2 13	
	G	7.8 15	8.9 33	9.8 30	5.1 39	3.5 21	31.6 30	35.1 32	38.3 37	21.7 12	16.2 9	
	H	8.4 9	10.0 52	11.2 52	6.3 23	3.7 19	33.3 31	37.0 42	38.8 53	25.1 13	16.5 9	
	T	7.3 15	8.3 20	9.3 34	4.5 7	3.4 21	30.3 22	33.9 26	35.9 36	19.2 11	15.4 8	
	A	9.4 13	11.1 14	12.7 25	6.2 33	4.4 20	37.2 30	40.4 45	43.6 50	25.4 16	19.4 10	
		parts	parts	parts	parts	parts	parts	parts	parts	parts	parts	
output		pants	pacts	pasts	partis	pars	pants	pacts	pasts	partis	pars	

Courier font than other fonts for less MSE and  $I_{avg}$ . 4) Due to a larger number of pixels in sentence and paragraph images, their MSEs are much lower than MSE of word images, but the absolute number of affected pixels stays the same.

#### 4.4 Effects of Hyper-parameter Settings

To better understand these attacks, we evaluate the effects of hyper-parameter settings, such as tradeoff  $c$ ,  $L_p$ -norms, font-weight, position of watermarks, etc.

**Effects of plain watermarks.** First, we evaluate the effects of plain watermarks, i.e., initial watermarks without perturbations. The accuracy of Courier, Georgia, Helvetica, Times, Arial drop to 0.066, 0.768, 0.531, 0.753 and 0.715, respectively, showing that it is quite trivial to launch *untargeted attacks* on OCR. Also, it further confirms that thinner font, Courier, is highly sensitive to plain watermarks. In samples of incorrect recognition, we count the percentage of output text shorter than the ground-truth text, Courier(0.65), Georgia(0.12), Helvetica(0.08), Times(0.09), Arial(0.14). It's easy to induce that thinner font is more sensitive to perturbations resulting in losing letters like the deletion case.

**Tradeoff between efficiency and visual quality.** In both gradient-based and optimization-based methods, we need to set tradeoff parameters to balance between efficiency ( $I_{avg}$ ) and visual quality (MSE). In Fig. 4, we plot the change of  $I_{avg}$  and MSE along with distinct tradeoff settings. In gradient-based methods, we can reduce perturbation level with smaller step size  $\alpha$ , at the cost of attack efficiency (raised  $I_{avg}$ ). In optimization-based methods, the main quality-efficiency tradeoff parameter  $c$  balances the targeted loss and the  $L_p$ -norm distance. We find that a smaller  $c$  makes the optimized objective function pay more attention to reduce the perturbation level with larger  $I_{avg}$ . Thus we can use the binary search for  $c$  to find adversarial examples with a lower perturbation level under

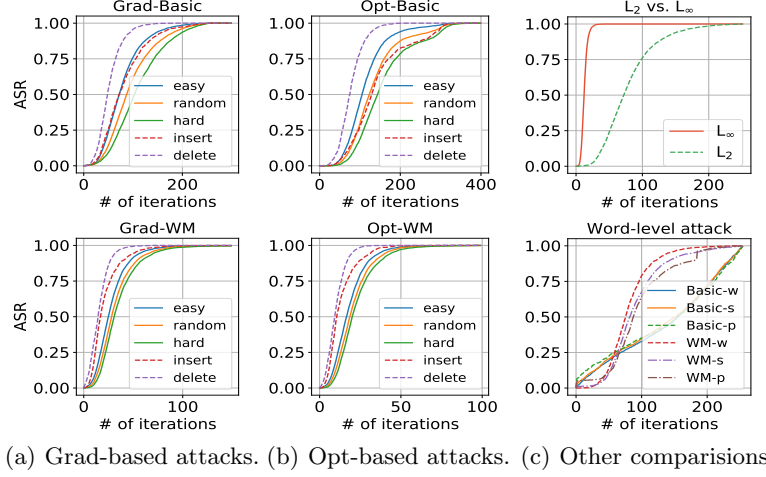
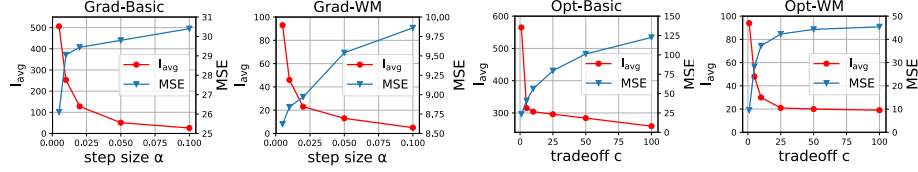


Fig. 3. Attack efficiency in word images with Arial font.

Fig. 4. In easy case of Arial word images, Grad-Basic and Grad-WM with different step size  $\alpha = 0.005, 0.010, 0.020, 0.050, 0.100$ . In easy case of Arial word images, Opt-Basic and Opt-WM with different tradeoff  $c = 1, 5, 10, 25, 50, 100$ .

the premise of maintaining perfect ASR. In addition, in various hyperparameter settings, the performance of watermark attacks is better than the basic attacks.

**Setting  $L_p$ -norms.**  $L_\infty$ -norm for measuring the maximum variation of perturbations has the same perturbed value at each perturbed pixel, and thus significantly narrows down operation space. Even more, observing  $L_\infty$ -norm examples in Table 4,  $L_\infty$ -norm causes severe background noise.  $L_2$ -norm yields perturbed values varying in perturbed pixels, and thus offers better flexibility to perform stronger adversarial attacks. Although  $L_2$ -norm is lower efficient than  $L_\infty$ -norm in Fig. 3(c),  $L_2$ -norm’s image quality metrics, MSE, PSNR and SSIM, in Table 4, all are better than those of  $L_\infty$ -norm. Intuitively,  $L_2$ -norm examples also have better visual quality. Therefore, we choose the  $L_2$ -norm in our experiments.

**Bold fonts settings.** We also investigate the bold version of the five fonts. Table 4 shows most bold fonts require slightly more MSE and  $I_{\text{avg}}$  to attack successfully, except for Courier requiring doubling MSE. This is not surprising, as bold fonts contain more useful pixels per letter. So they need more perturbations.

**Table 3.** Word-level attacks in word, sentence and paragraph images.

		word image				sentence image				paragraph image			
		MSE	PSNR	SSIM	$I_{avg}$	MSE	PSNR	SSIM	$I_{avg}$	MSE	PSNR	SSIM	$I_{avg}$
Grad-Basic	Courier	50.6	31.1	0.944	235	5.8	40.5	0.993	153	7.0	39.7	0.993	113
	Georgia	118.6	27.4	0.900	326	14.6	36.5	0.988	239	21.3	34.9	0.986	203
	Helvetica	124.0	27.2	0.894	254	14.2	36.6	0.988	238	22.5	34.6	0.984	233
	Times	114.4	27.5	0.904	291	13.2	36.9	0.989	201	17.4	35.7	0.989	164
	Arial	133.9	26.9	0.888	222	15.8	36.1	0.987	273	23.2	34.5	0.984	242
	example	parts				This one did exactly that.							
Grad-WM	Courier	12.9	37.0	0.993	51	3.4	42.8	0.999	60	8.7	38.7	0.998	57
	Georgia	35.4	32.6	0.985	131	5.2	41.0	0.999	90	9.2	38.5	0.998	81
	Helvetica	40.0	32.1	0.984	124	5.4	40.8	0.999	91	11.8	37.4	0.998	96
	Times	34.8	32.7	0.985	160	4.5	41.6	0.999	79	6.7	39.9	0.999	74
	Arial	44.6	31.6	0.982	138	6.1	40.3	0.999	98	12.0	37.3	0.998	98
	example	parts				This one did exactly that.							
target output		taupe				Tale one did exactly that.							

**Table 4.** Comparison of  $L_2$ -norm and  $L_\infty$ -norm in Grad-Basic. Examples' target output are "parts". Fonts are denoted by their initial letters. Bold fonts marked as 'b'.

	$L_2$						$L_2$					$L_\infty$				
	MSE	PSNR	SSIM	$I_{avg}$	example		MSE	PSNR	SSIM	$I_{avg}$	example	MSE	PSNR	SSIM	$I_{avg}$	example
Cb	23.6	34.4	0.977	69	<b>parts</b>	C	10.5	37.9	0.988	59	<b>parts</b>	72.6	29.5	0.771	11	<b>parts</b>
Gb	25.9	34.0	0.975	59	<b>parts</b>	G	27.4	33.7	0.974	43	<b>parts</b>	163.6	26	0.645	11	<b>parts</b>
Hb	32.8	33.0	0.970	75	<b>parts</b>	H	27.0	33.8	0.975	51	<b>parts</b>	159.4	26.1	0.658	10	<b>parts</b>
Tb	26.3	33.9	0.975	66	<b>parts</b>	T	26.4	33.9	0.975	62	<b>parts</b>	156.6	26.2	0.653	11	<b>parts</b>
Ab	34.2	32.8	0.968	69	<b>parts</b>	A	29.8	33.4	0.972	51	<b>parts</b>	169.3	25.8	0.656	11	<b>parts</b>

**Attacking sequence-based OCRs.** As OCR recognizes entire sequences rather than individual characters, we demonstrate that we can replace a letter even the watermark does not overlap with the letter, indicating we will not add perturbations over it. In this experiment, we shift the watermark about 10 pixels right to the target letter when generating initial watermarked images. Table 5 summarizes results that the attacks require 18 times MSE and 34 times  $I_{avg}$ , and ASRs drop to less than 50% except for the simplest Courier case. It confirms the fact that in sequential-based OCRs, the influence around the letter is not as important as that overlaying the letter. Also, it reveals the necessity of finding the position of watermarks accurately to perform strong adversarial attacks.

#### 4.5 Extensions and Applications of Watermark Attacks

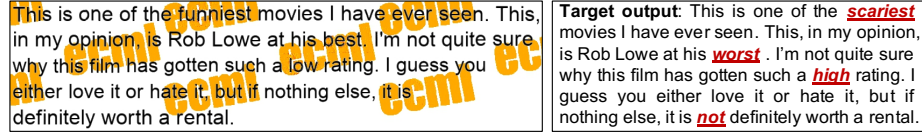
**Full-color watermarks.** Sometimes adding gray watermarks still hurts human readability. We use the fact that modern OCR first transforms colored images into gray ones before recognition. Colored watermarks significantly improve overall readability when mixed with black texts. Fig. 5 shows a colored-watermark

**Table 5.** Non-overlapping Grad-WM in easy case of word images. Add a shift about 10 pixels right to the target letter.

	MSE	I <sub>avg</sub>	ASR	example	output
Courier	48.5	186	0.83	love	hove
Georgia	156.6	546	0.36	move	move
Helvetica	146.8	474	0.45	parts	pants
Times	145.9	515	0.43	hoy	boy
Arial	153.5	463	0.44	broad	bread

**Table 6.** Protection mechanism. Acc and Acc\* are the prediction accuracy with and without protection, respectively.

	MSE	I <sub>avg</sub>	Acc	Acc*	example
Courier	0.6	5	1.0	0.066	part
Georgia	0.2	1	1.0	0.768	part
Helvetica	0.2	1	1.0	0.531	part
Times	0.2	1	1.0	0.753	part
Arial	0.2	1	1.0	0.715	part

**Fig. 5.** Full-color watermarks on a paragraph image. MSE/PSNR/SSIM: 9.36/38.42/0.998.

example of altering a positive movie review into a negative one by replacing and inserting words in it. Note that not all watermarks are malicious (e.g., the watermark on “I have ever seen” of Fig. 5 does not include adversarial perturbations). We evenly distribute watermarks over the paragraph image to make it look more similar to a naturally watermarked paragraph image in the real-world scenario.

**Attacking Chinese Characters.** In addition to English, we show that the method is applicable to other languages. Fig. 6 shows a Chinese example where we almost altered all important information, and its perturbation level is much larger than that in Fig. 5, because of the complex structure of Chinese characters.

**Using FAWA to enhance the OCR readability of watermarked contents.** Sometimes we want people to notice vital watermarks (e.g., urgent, confidential), but we do not want them to affect OCR’s accuracy. In such case, we generate accuracy-enhancing watermarks by setting the ground-truth as the target. In Table 6, with few I<sub>avg</sub> and MSEs, Acc (accuracy with protection mechanism) increases to 1.0 compared with the low Acc\* (accuracy of initial watermarked images). It works because protective perturbations strengthen the features of the ground-truth and boost the confidence of the original, making it more “similar” to the target (i.e., the ground-truth). As the target is the same as

**Fig. 6.** A Chinese paragraph example. MSE/PSNR/SSIM: 735.34/19.46/0.697.

the original, both MSE and  $I_{\text{avg}}$  stay low. So we can use FAWA as a protection mechanism to produce both human-friendly images and OCR-friendly images.

## 5 Conclusion

DNN-based OCR systems are vulnerable to adversarial examples. On the text images, we hide perturbations in watermarks, making adversarial examples look natural to human eyes. Specifically, we develop FAWA that automatically generate the adversarial watermarks targeting sequence-based OCR models. In extensive experiments, while maintaining perfect attack success rate, FAWA exhibits the outstanding attack capability. For example, FAWA boosts the attack speed up to 8 times, and reduces the perturbation level lower than 40% on average. In word, sentence and paragraph contexts, FAWA works well with letter-level and word-level targets. We further extend our natural watermarked adversarial examples in many scenarios, such as full-color watermarks to increase the text readability, applicability for other languages, protection mechanism for enhancing OCR’s accuracy. We believe human-eye-friendly adversarial samples are applicable in many other scenarios, and we plan to explore them as future work.

## Acknowledgments

This work is supported in part by the National Natural Science Foundation of China (NSFC) Grant 61532001 and the Zhongguancun Haihua Institute for Frontier Information Technology.

## References

1. Bengio, Y., LeCun, Y., Nohl, C., Burges, C.: Lerec: A nn/hmm hybrid for on-line handwriting recognition. *Neural Computation* **7**(6), 1289–1303 (1995)
2. Breuel, T.M., Ul-Hasan, A., Al-Azawi, M.A., Shafait, F.: High-performance ocr for printed english and fraktur using lstm networks. In: 2013 12th International Conference on Document Analysis and Recognition. pp. 683–687. IEEE (2013)
3. Brown, T.B., Mané, D., Roy, A., Abadi, M., Gilmer, J.: Adversarial patch. *arXiv preprint arXiv:1712.09665* (2017)
4. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: 2017 IEEE Symposium on Security and Privacy (SP). pp. 39–57. IEEE (2017)
5. Chen, P.Y., Sharma, Y., Zhang, H., Yi, J., Hsieh, C.J.: Ead: elastic-net attacks to deep neural networks via adversarial examples. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
6. Dong, Y., Liao, F., Pang, T., Su, H., Zhu, J., Hu, X., Li, J.: Boosting adversarial attacks with momentum. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9185–9193 (2018)
7. Espana-Boquera, S., Castro-Bleda, M.J., Gorbe-Moya, J., Zamora-Martinez, F.: Improving offline handwritten text recognition with hybrid hmm/ann models. *IEEE transactions on pattern analysis and machine intelligence* **33**(4), 767–779 (2011)

8. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
9. Graves, A., Fernández, S., Gomez, F., Schmidhuber, J.: Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In: Proceedings of the 23rd international conference on Machine learning. pp. 369–376. ACM (2006)
10. Hanwei Zhang, Yannis Avrithis, T.F.L.A.: Smooth adversarial examples. arXiv preprint arXiv:1903.11862 (2019)
11. Heng, W., Zhou, S., Jiang, T.: Harmonic adversarial attack method. arXiv preprint arXiv:1807.10590 (2018)
12. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8), 1735–1780 (1997)
13. Karmon, D., Zoran, D., Goldberg, Y.: Lavan: Localized and visible adversarial noise. arXiv preprint arXiv:1801.02608 (2018)
14. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
15. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems. pp. 1097–1105 (2012)
16. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial machine learning at scale. arXiv preprint arXiv:1611.01236 (2016)
17. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
18. 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. pp. 2574–2582 (2016)
19. Nguyen, A., Yosinski, J., Clune, J.: Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 427–436 (2015)
20. Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z.B., Swami, A.: The limitations of deep learning in adversarial settings. In: 2016 IEEE European Symposium on Security and Privacy (EuroS&P). pp. 372–387. IEEE (2016)
21. Song, C., Shmatikov, V.: Fooling ocr systems with adversarial text images (2018)
22. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)
23. Wang, G., Chen, X., Xu, C.: Adversarial watermarking to attack deep neural networks. In: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1962–1966. IEEE (2019)
24. Wang, T., Wu, D.J., Coates, A., Ng, A.Y.: End-to-end text recognition with convolutional neural networks. In: Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012). pp. 3304–3308. IEEE (2012)
25. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)
26. Wick, C., Reul, C., Puppe, F.: Calamari-a high-performance tensorflow-based deep learning package for optical character recognition. arXiv preprint arXiv:1807.02004 (2018)
27. Xiao, C., Zhu, J.Y., Li, B., He, W., Liu, M., Song, D.X.: Spatially transformed adversarial examples. ArXiv **abs/1801.02612** (2018)