

Attacking Optical Character Recognition (OCR) Systems with Adversarial Watermarks

Lu Chen¹ and Wei Xu¹

¹Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing, China
lchen17@mails.tsinghua.edu.cn
weixu@tsinghua.edu.cn

Abstract

Optical character recognition (OCR) is widely applied in real applications serving as a key preprocessing tool. The adoption of deep neural network (DNN) in OCR results in the vulnerability against adversarial examples which are crafted to mislead the output of the threat model. Different from vanilla colorful images, images of printed text have clear backgrounds usually. However, adversarial examples generated by most of the existing adversarial attacks are unnatural and pollute the background severely. To address this issue, we propose a watermark attack method to produce natural distortion that is in the disguise of watermarks and evade human eyes' detection. Experimental results show that watermark attacks can yield a set of natural adversarial examples attached with watermarks and attain similar attack performance to the state-of-the-art methods in different attack scenarios.

Introduction

Optical Character Recognition (OCR) is a widely adopted application for converting printed or handwritten images to text, which becomes a critical preprocessing component in text analysis pipelines, such as document retrieval and summarization. OCR has been significantly improved in recent years thanks to the wide adoption of the deep neural network (DNN), and thus deployed in many critical applications where OCR's quality is vital. For example, photo-based ID recognition depends on OCR's quality to automatically structure information into databases, and automatic trading sometimes relies on OCR to read certain news articles for determining the sentiment of news.

Unfortunately, OCR also inherits all counter-intuitive security problems of the DNNs. Especially, the OCR model is also vulnerable to *adversarial examples*, which are crafted by making human-imperceptible perturbations on original images with the intent of misleading the model. The wide adoption of OCR in real pipelines gives more incentives for adversaries to game the OCR, such as causing fake ID information, incorrect readings of metrics or instructions, etc. Figure 2 and 3 in the evaluation section illustrate two real-world examples with attacking the ID number and financial

report number. This paper provides a preliminary study on the possibility of OCR attacks.

Many prior works (Nguyen, Yosinski, and Clune 2015; Goodfellow, Shlens, and Szegedy 2014; Papernot et al. 2016a; Szegedy et al. 2013) have shown that changing the prediction of DNNs is practicable by applying carefully-designed perturbations (usually background noise) to original images, in traditional image classification tasks. Recent projects, Adversarial Patch (Brown et al. 2017) and LaVAN (Karmon, Zoran, and Goldberg 2018), introduced the adversarial patch attack, which puts visible perturbations confined to a small region or location in the image.

However, these methods are not directly applicable to OCR attacks for the following three reasons:

First, the input image to OCR is on a white paper with a spotless background. Thus any perturbation added by existing attacks will appear so obvious to human readers that it will cause suspicion.

Second, in complex languages like Chinese, there are many characters (e.g., the dataset we use contains 5,989 unique characters). If an adversary wants to perform a *targeted* attack, i.e., changing one character to another specific one (target) in a sentence and meanwhile resulting in semantically meaningful recognition results, it requires a large number of perturbations that are too obvious to hide.

Third, instead of classifying characters individually, the modern OCR model is an end-to-end neural network, inputting a variable-sized image and outputting sequences of labels. In other words, it works on feeding images line by line. It is usually called the *sequential labeling* task, which is relatively harder to be attacked than the image classification task. It is insufficient to just add perturbations to a single character. Instead, the perturbations are required to span multiple characters. Also, as the OCR model is end-to-end, the internal feature representations rely on nearby characters (contexts). Thus the perturbations of attacking a single character are designed given its contexts.

In this preliminary study, we propose a new attack method, `WATERMARK` attack, against modern OCR models. Watermarks are images or patterns commonly used in documents as background to identify things, such as marking a document proprietary, secret, urgent, or even simply as

decoration. Similar to watermarks, in Asian countries, documents often contain stamps to verify their authority. Human eyes are so used to these watermarks and ignore them. In this paper, we generate natural watermark-style perturbations. That is, we limit all perturbations within a small region of the image, i.e., a watermark. Given the bound, we minimize the perturbation level. In comparison, classic adversarial examples spread noise all over the image. Our approach is similar to the patch-based attacks (Brown et al. 2017; Karmon, Zoran, and Goldberg 2018). Different from that patches absolutely cover part of the images, watermarks do not hinder text’s readability and thus look more natural. (Heng, Zhou, and Jiang 2018) generated disguising perturbations as shadows, exposure problems or color problems. And (Hanwei Zhang 2019) generated smooth noise by Laplacian smoothing. But none of these solve the clear background challenge for OCR.

We focus on *white-box, targeted attack* in this paper. That is, we assume adversaries have perfect knowledge of the DNN architecture and parameters (white-box model) and aim to generate specific recognition results (targeted attack). Given that many real OCR softwares are based on similar open-source OCR DNN models, we believe the white-box model, in addition to being a starting point, also has real-world applicability.

As a consequence, the `WATERMARK` attack is an adversarial attack on the OCR model. The `WATERMARK` attack attaches natural watermark-style noise, tricks the OCR model into outputting specific recognition results, and preserves the readability of adversarial images at the same time. To some extent, the `WATERMARK` attack solves the clear background problem.

As an evaluation, we performed the `WATERMARK` attack on a state-of-the-art open-source OCR model using DenseNet + CTC neural network architectures for the Chinese language. We used a dataset with 3.64 million images and 5,989 unique characters. With 158 pairs of original-target, we show that the `WATERMARK` attack can generate quite human-eye friendly adversarial samples with a high probability of success. Some of `WATERMARK` adversarial examples even work on Tesseract OCR (Google 2019) in a black-box manner.

Even more, we applied our model to real-world scenarios. In Figure 3, we employed the `WATERMARK` method to a page of an annual report of a listed Chinese company and changed the semantics, in the meantime, the image looks natural to human readers.

The contributions of this paper include: 1) We propose the `WATERMARK` adversarial attack to generate natural-to-human watermark-style perturbations, targeting DNN-based OCR systems. We also demonstrated a method to hide perturbations that human eyes are accustomed to, in a watermark-bounded region. 2) Using difficult OCR cases (Chinese), we demonstrated the success rate of `WATERMARK` attacks comparing to existing ones.

Background and Related Work

Optical Character Recognition (OCR)

Generally, the OCR pipeline, as shown in Figure 1, begins with line segmentation, which includes page layout analysis for locating the position of each line, de-skewing the image, and segmenting the input image into line images. After pre-processing line images, such as rescaling and normalizing, such images are fed into the recognition model, which outputs recognition results.

There are two types of OCR models. 1) Character-based models are the classic way (Smith 2007). Such a recognition model segments the image into per-character sub-images and classifies each sub-image into the most likely recognition result. Obviously, its performance heavily relies on the character segmentation. 2) End-to-end models are a segmentation-free approach. It recognizes entire sequences of characters in a variable-sized image. (Bengio et al. 1995; Espana-Boquera et al. 2011) adopted sequential models. (Breuel et al. 2013; Wang et al. 2012) utilize DNNs as the feature extractor for the end-to-end sequential recognition. Sequential DNN models (Graves et al. 2006) introduced a segmentation-free approach, connectionist temporal classification(CTC), which allows variable-sized input images and output results.

In end-to-end models, *sequence labeling* is a task that assigns a sequence of discrete labels to variable-sized sequential input data. In our case, the input is a variable-size image \mathbf{x} and the output is a sequence of characters $\mathbf{t} = [t_1, t_2, \dots, t_N]$, $t_i \in \mathcal{T}$ from predefined character set \mathcal{T} .

Connectionist Temporal Classification (CTC). CTC is an alignment-free method for training DNNs on the sequential labeling task, which provides a kind of loss enabling us to recognize sequences without explicit segmentation while training DNNs. Therefore, many state-of-the-art OCR models use CTC as the model’s loss function. Given the input image \mathbf{x} , let $f(\mathbf{x}) = \mathbf{y} = [y_1, y_2, \dots, y_M]$ be the sequence of model f ’s outputs, where $M \geq N$ and $y_i \in [0, 1]^{|\mathcal{T}|}$ is the probability distribution over the character set \mathcal{T} in observing label i .

CTC requires calculating the likelihood $p(\mathbf{t}|\mathbf{x})$, which is barely directly measured from the model’s probability distribution $f(\mathbf{x})$ and the target sequence \mathbf{t} . To settle this, CTC uses a valid alignment $\mathbf{a} = [a_1, a_2, \dots, a_M]$ of \mathbf{t} , $a_i \in \mathcal{T} \cup \{\text{blank}\}$, where the target sequence \mathbf{t} can be obtained by removing all blanks and sequential duplicate characters (e.g. both $[a, -, a, b, -]$ and $[-, a, a, -, -, a, b, b]$ are valid alignments of $[a, a, b]$). The likelihood $p(\mathbf{t}|\mathbf{x})$ is to sum up the probability of all possible valid alignments denoted as A .

$$p(\mathbf{t}|\mathbf{x}) = \sum_{\mathbf{a} \in A} \prod_{i=1}^M p(a_i|\mathbf{x}) = \sum_{\mathbf{a} \in A} \prod_{i=1}^M (y_i)_{a_i} \quad (1)$$

The negative log-likelihood of $p(\mathbf{t}|\mathbf{x})$ is the CTC loss function $\ell_{\text{CTC}}(f(\mathbf{x}), \mathbf{t})$.

$$\ell_{\text{CTC}}(f(\mathbf{x}), \mathbf{t}) = -\log p(\mathbf{t}|\mathbf{x}) \quad (2)$$

To obtain the most probable output sequence $\arg \max_t p(t|\mathbf{x})$, a greedy path decoding algorithm can select the most probable alignment at each step. However, the greedy algorithm does not guarantee to find the most probable labeling. A better method, beam search decoding, simultaneously keeps a certain number of the most probable alignments at each step and chooses the most probable output in the top-alignment list.

Attacking DNN-based computer vision tasks

Where to add perturbations? Attacking DNN models is a popular topic in both computer vision and security fields. Many projects focus on finding small L_p -bounded perturbations, hoping that the bound L_p will keep the perturbations visually imperceptible. FGSM (Goodfellow, Shlens, and Szegedy 2014), L-BFGS (Szegedy et al. 2013), DeepFool (Moosavi-Dezfooli, Fawzi, and Frossard 2016), Carlini L_2, L_∞ (Carlini and Wagner 2017), PGD (Madry et al. 2017) and EAD (Chen et al. 2018) all performed modifications at the pixel level by a small amount bounded by ϵ .

Other attacks such as JSMA (Papernot et al. 2016a), Carlini L_0 (Carlini and Wagner 2017), Adversarial Patch (Brown et al. 2017) and LaVAN (Karmon, Zoran, and Goldberg 2018), perturb a small region of pixels in an image but the pixel-level perturbations are not bounded by L_p .

As we have mentioned, neither approach can hide perturbations from the normal human vision in OCR tasks, as a document with enough readability usually has a spotless background and vivid text, which is greatly different from natural RGB photos.

How to generate perturbations? There are two types of methods to generate perturbations.

1) Gradient-based attack is to add perturbations generated by gradient against input pixels. Formally, we can describe the general problem as: For a L_∞ -bounded adversary, we compute an adversarial example \mathbf{x}' given an original image \mathbf{x} and target labels t where perturbations' bound $\epsilon > 0$ is tiny enough to be indistinguishable to human observers.

Fast Gradient Sign Method (FGSM) (Goodfellow, Shlens, and Szegedy 2014) is a one-step attack that obtains the adversarial image \mathbf{x}' as $\mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(\mathbf{x}', t))$. The original image \mathbf{x} takes a gradient sign step with step size ϵ in the direction that increases the probability of the target label t . It is efficient, but it only provides a coarse approximation of the optimal perturbations.

Basic Iterative Method (BIM) (Kurakin, Goodfellow, and Bengio 2016) takes multiple smaller steps α and the result image is clipped by the same bound ϵ : $\mathbf{x}'_i = \mathbf{x}'_{i-1} + \text{clip}_\epsilon(\alpha \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(\mathbf{x}'_{i-1}, t)))$, where \mathbf{x}'_i is an adversarial example yielded at step i . BIM produces superior results to FGSM.

Momentum Iterative Method (MIM) (Dong et al. 2018) extends BIM with a momentum item. MIM can not only stabilize update directions but also escape from poor local maxima during the iteration. Thus, it generates more transferable adversarial examples. Each iterative update is to adjust the update direction and generate new adversarial image \mathbf{x}'_i using the momentum item \mathbf{g}_i , as following $\mathbf{g}_i =$

$$\mu \cdot \mathbf{g}_{i-1} + \frac{\nabla_{\mathbf{x}} \ell(\mathbf{x}'_{i-1}, t)}{\|\nabla_{\mathbf{x}} \ell(\mathbf{x}'_{i-1}, t)\|_p}, \mathbf{x}'_i = \mathbf{x}'_{i-1} + \text{clip}_\epsilon(\alpha \cdot \text{sign}(\mathbf{g}_i)),$$

where μ is the decay factor.

2) Optimization-based attack directly solves the optimization problem of minimizing the L_p distance between the original example \mathbf{x} and the adversarial example \mathbf{x}' and yielding the incorrect classification.

$$\begin{aligned} & \text{minimize} && \|\mathbf{x}' - \mathbf{x}\|_p \\ & \text{subject to} && f(\mathbf{x}') = t \text{ and } \mathbf{x}' \in [\mathbf{x}_{\min}, \mathbf{x}_{\max}]^{|\mathbf{x}|} \end{aligned}$$

Box-constraint L-BFGS (Szegedy et al. 2013) finds the adversarial examples by solving the box-constraint problem, minimize $\lambda \|\mathbf{x}' - \mathbf{x}\|_p + \ell(\mathbf{x}', t)$ subject to $\mathbf{x}' \in [\mathbf{x}_{\min}, \mathbf{x}_{\max}]^{|\mathbf{x}|}$, where $\ell(\mathbf{x}', t)$ is the cross-entropy loss between the logit output and the target label t . Although the perturbations generated by L-BFGS are much less than the gradient-based attack, L-BFGS has far low efficiency.

C&W (Carlini and Wagner 2017) is a L_p -oriented attack that can successfully break undefended and defensively distilled DNNs. Given the logit Z of the model f , other than applying cross-entropy as the loss function, C&W attack designed a new loss function $\ell(\mathbf{x}', t) = \max(\max\{Z(\mathbf{x}')_i : i \neq t\} - Z(\mathbf{x}')_t, -\kappa)$ solved by gradient descent, where κ is the confidence of misclassification.

Defense methods against these attacks

People have proposed many practical defense methods against adversarial examples. *Adversarial training* (Tramèr et al. 2017) improves the robustness of DNNs by injecting label-corrected adversarial examples into the training procedure and training a robust model that has a resistance to perturbations generated by gradient-based methods. *Defensive distillation* (Papernot et al. 2016b) defends against adversarial perturbations using the distillation techniques (Hinton, Vinyals, and Dean 2015) to retrain the same network with class probabilities predicted by the original network. There are also methods focusing on detecting adversarial samples (Xu, Evans, and Qi 2017; Lu, Issaranon, and Forsyth 2017; Grosse et al. 2017; Feinman et al. 2017).

Methodology

Preliminaries. We assume that attackers have full knowledge of the threat model, such as the model architecture and parameter values. Given an input image $\mathbf{x} \in [\mathbf{x}_{\min}, \mathbf{x}_{\max}]^{|\mathbf{x}|}$, an adversarial image $\mathbf{x}' \in [\mathbf{x}_{\min}, \mathbf{x}_{\max}]^{|\mathbf{x}|}$ where $|\mathbf{x}|$ is the number of pixels in the original image, OCR model f , the adversarial example's prediction result is $f(\mathbf{x}')$. Given target label t , loss function of the target model f with respect to the input image \mathbf{x} is $\ell(\mathbf{x}, t)$. Besides, we assume $f(\mathbf{x})$ to be differentiable almost everywhere (sub-gradient may be used at discontinuities). Because the gradient-descent approach is applicable to any DNNs with a differentiable discriminant function.

Distance Metric. We define the distance metric L_p to quantify the similarity between the original image \mathbf{x} and the corresponding adversarial image \mathbf{x}' . Such a distance metric reflects the cost of manipulations. L_p -norm is a widely-used

distance metric defined as $\|x - x'\|_p = \sqrt[p]{\sum_{i=1}^d |x_i - x'_i|^p}$ where d -dimensional vector $\mathbf{x} = [x_1, \dots, x_d]$ for any $p \in \{0, 1, 2, \infty\}$. L_1 accounts for the total variation in the perturbations and serves as a popular convex surrogate function of the L_0 that measures the number of modified pixels. L_2 measures the standard Euclidean distance, which is usually used to improve the visual quality. L_∞ measures the maximum change of the perturbations between \mathbf{x} and \mathbf{x}' .

WATERMARK attack to CTC-based OCR

In this paper, we propose the MIM-based WATERMARK attack on the CTC-based OCR model to generate adversarial examples. In this section, we will first introduce how to integrate watermarks into MIM (Dong et al. 2018), which induces the MIM-based WATERMARK attack method (WM) to generate adversarial examples satisfying the L_∞ -norm restriction in the targeted and non-targeted attack fashion. We then present several variants of WM to L_∞ -norm bound. The generation pipeline of the WATERMARK adversarial attack is illustrated in Figure 1. Table 1 shows adversarial examples generated from each method.

MIM-based WATERMARK attack (WM). Watermark widely occurs in a mass of documents and files. Making use of the popularity of the watermark in the documents, we apply the idea of the watermark to decorate perturbations as the watermark; that is, we restrict the manipulation region on a specific predefined watermark-shape region Ω .

To generate a targeted L_∞ -bounded adversarial example \mathbf{x}' , we start with an initial image $\mathbf{x}'_0 = \mathbf{x}$ given an original image \mathbf{x} . WM seeks the adversarial example by solving the constrained optimization problem

$$\begin{aligned} \arg \min_{\mathbf{x}'} \ell_{\text{CTC}}(\mathbf{x}', \mathbf{t}) \\ \text{s.t. } \|\mathbf{x}' - \mathbf{x}\|_\infty \leq \epsilon \text{ and } (1 - \Omega) \odot (\mathbf{x}' - \mathbf{x}) = \mathbf{0}, \end{aligned}$$

where ϵ is the size of adversarial perturbations and $\Omega = \{o \in \Omega | o = 1 \text{ inside the watermark, otherwise, } o = 0, \Omega \in \{0, 1\}^{|\mathbf{x}'|\}$. We summarized WM in Algorithm 1. At each attacking iteration i , the attacker first feeds the adversarial example \mathbf{x}'_i to the OCR model f and obtain the gradient $\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})$ through back-propagation. Then, for the purpose of stabilizing update directions and escaping from poor local maxima, update momentum item \mathbf{g}_{i+1} by accumulating the velocity vector in the gradient direction as Equation 3 shown in Algorithm 1. Last, update new adversarial example \mathbf{x}'_{i+1} by applying the Ω -restricted sign gradient with small step size α , and clip the intermediate perturbations to ensure them in the ϵ -ball of \mathbf{x}'_0 as Equation 4. The attacking iteration proceeds until the attack is successful or reaches the maximum iterations I .

Variants of WM attack. To present a more natural appearance of watermark-like perturbations, we design three variants of WM attack.

WM_{init}. Watermark region Ω element-wisely multiplies the sign gradient $\text{sign}(\mathbf{g}_{i+1})$ which attackers only operate pixels inside the watermark region Ω . As shown in the WM column of Table 1, the perturbations of WM adversarial examples

Algorithm 1 MIM-based WATERMARK example generation

Input: A clean image \mathbf{x} , OCR model f with CTC loss ℓ_{CTC} , ground-truth text $f(\mathbf{x})$, targeted text \mathbf{t} , watermark modification region Ω , ϵ -ball perturbation, # of iterations I , decay factor μ
Output: An adversarial example \mathbf{x}' with $\|\mathbf{x}' - \mathbf{x}\|_\infty \leq \epsilon$ or attack failure \perp

- 1: Initialization: $\alpha = \epsilon/I$; $\mathbf{g}_0 = \mathbf{0}$; $\mathbf{x}'_0 = \mathbf{x}$
- 2: **for all** each iteration $i = 0$ to $I - 1$ **do**
- 3: Input \mathbf{x}'_i to f and obtain the gradient $\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})$
- 4: Update \mathbf{g}_{i+1} by accumulating the velocity vector in the gradient direction as

$$\mathbf{g}_{i+1} = \mu \cdot \mathbf{g}_i + \frac{\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})}{\|\nabla_{\mathbf{x}} \ell_{\text{CTC}}(\mathbf{x}'_i, \mathbf{t})\|_p} \quad (3)$$

- 5: Update \mathbf{x}'_{i+1} by applying watermark-bound sign gradient as

$$\mathbf{x}'_{i+1} = \mathbf{x}'_i + \text{clip}_\epsilon(\alpha \cdot \Omega \odot \text{sign}(\mathbf{g}_{i+1})) \quad (4)$$

- 6: **if** $f(\mathbf{x}'_{i+1}) == \mathbf{t}$ **then**
 - 7: **return** $\mathbf{x}' = \mathbf{x}'_{i+1}$
 - 8: **end if**
 - 9: **end for**
 - 10: **return** failure \perp
-

are not dense enough to construct a complete watermark and be a natural watermark. Thus, for filling in the blanks of the watermark region, we start from an initial watermark-pasted image by attaching a watermark to the original image, $\mathbf{x}'_0 = \mathbf{x} + \lambda \cdot (\mathbf{x} > \tau) \odot \Omega$, where λ is the grayscale value of the pasted watermark and $(\mathbf{x} > \tau)$ denotes the position except the text.

WM_{neg}. The sign of the gradient, $\text{sign}(\cdot)$, can be -1 or +1 based on the direction of gradient descent. When the gradient is positive, $\text{sign}(\cdot) = 1$, the pixel value will increase, that is, the pixel looks whiter (the maximum of the grayscale value mean the whitest color or else blackest color). Otherwise, the pixel value decreases, and the pixel looks blacker. Obviously, the pixels in the text region become whiter resulting in the fuzzy text, and it's meaningless to whiten the clear white background. Thus, we only need to keep the negative gradient and leave the positive gradient behind. We generated WM noise but only kept the negative gradient during attacking iteration. After adding the new constraint, the update step of new adversarial example, Equation 4, becomes

$$\mathbf{x}'_{i+1} = \mathbf{x}'_i + \text{clip}_\epsilon(\alpha \cdot \Omega \odot \min(0, \text{sign}(\mathbf{g}_{i+1}))). \quad (5)$$

WM_{edge}. A different way to add perturbations is to confine watermark region around the text edges, pretending to be defectives in printing.

WM_{edge} is similar with WM. We define the watermark as the region of the text edge, which can be obtained by image erosion in morphological image processing. The erosion operation \ominus erodes the original image using the specified structuring element that determines the shape of a pixel neighborhood over which the minimum is taken. In experiments, we use a 3×3 rectangular structuring element as a kernel, $\mathcal{K} := \mathbf{1}^{3 \times 3}$. We take the bolder text region after erosion as

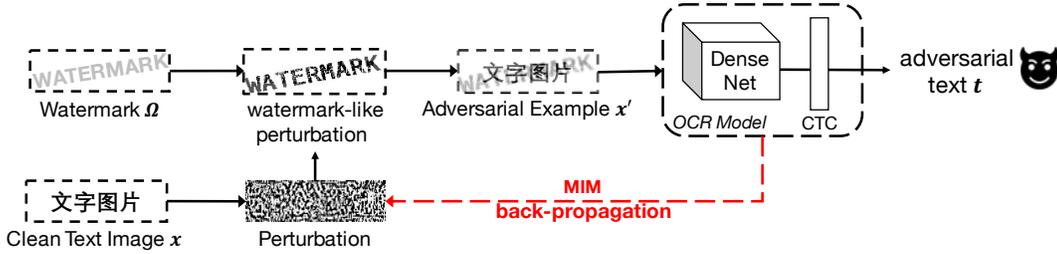


Figure 1: The pipeline of the WATERMARK attack. We generate noise using MIM with CTC loss function back propagating the targeted DenseNet and then mask the noise outside the watermark region. We iterating the procedure above until an iterations threshold.

	original	MIM	WM	WM _{init}	WM _{neg}	WM _{edge}	OCR output	English translation
substitution	靠左行驶	靠左行驶	靠左行驶 ^K	靠左行驶 ^K	靠左行驶 ^K	靠左行驶	靠右行驶	drive left → drive right
	9月1号	9月1号	9月1号 ^K	9月1号 ^K	9月1号 ^K	9月1号	9月9号	1 Sep. → 9 Sep.
	我叫小方	我叫小方	我叫小方 ^K	我叫小方 ^K	我叫小方 ^K	我叫小方	我叫孙方	I am Xiao Fang → I am Sun Fang
	hire	hire	hire ^K	hire ^K	hire ^K	hire	fire	hire → fire
	一点下课	一点下课	一点下课 ^K	一点下课 ^K	一点下课 ^K	一点下课	二点下课	class is over at 1 → class is over at 2
-	一点下课	一点下课	一点下课 ^K	一点下课 ^K	一点下课 ^K	一点下课	一下课	class is over at 1 → once class is over
+	一点下课	一点下课	一点下课 ^K	一点下课 ^K	一点下课 ^K	一点下课	一点不下课	class is over at 1 → class is not over at 1

Table 1: Adversarial examples with different attacks. The last two rows show text deletion and insertion. Other rows show text substitution.

the watermark. Thus, the text-edge shape watermark Ω_{edge} is defined as $\Omega_{\text{edge}} = \{o \in \Omega_{\text{edge}} | o = \begin{cases} 1, & \text{if } \mathbf{x} \ominus \mathcal{K} > \tau \\ 0, & \text{if } \mathbf{x} \ominus \mathcal{K} \leq \tau \end{cases}$.

Experiments

In this section, we generate adversarial examples on the CTC-based OCR model. We compared the performance of the basic MIM, WM, and its variants.

Setup

Threat model. We performed WM attack on the DenseNet + CTC neural network¹ which is trained in the Chinese text image dataset. DenseNet (Huang et al. 2017) is one of powerful DNNs for visual recognition, which can capture complex features. Thus, we utilize DenseNet as the feature extractor and CTC (Graves et al. 2006) as the loss function. In the test phase, the DenseNet+CTC OCR recognition model achieved 98.3% accuracy on the validation dataset that involves 36400 images. The Chinese text image dataset has 3.64 million images that are generated by altering fonts, scale, grayscale, blur, sketch based on Chinese news articles. The character set has 5989 unique characters, including Chinese and English characters, digits, punctuations.

¹https://github.com/YCG09/chinese_ocr

Attack setting. The attack setting is applied among all experiments. Our experiment setup is based on MIM’s framework. We use the implementation of MIM in CleverHans package². We use the attack setting which runs $I = 1000$ iterations at the most. We utilize an early stopping criterion based on the attacking result at each iteration. The L_{∞} -norm perturbations ϵ is bounded by 0.2. The pixel value of the image ranges from 0 (black) to 1 (white). For the initial watermark in WM_{init}, we set the grayscale value λ of the watermark to 0.3, and we put the watermark in the center of the image by default. The watermarks are set to the font size 30.

Choose attacked candidates. To satisfy the semantic fluency in our OCR attack, we choose 691 pairs of antonym characters with high similarity of character shape, $\text{Sim}(c, c') > 0.6$ ³, so that the adversarial attack only requires adversarial perturbations as less as possible to fool the OCR system. Then we match selected antonym character pairs in the corpus of People’s Daily in 1988 and choose

²<https://github.com/tensorflow/cleverhans>

³Given two characters c and c' , character similarity is defined as $\text{Sim}(c, c') = w_1 \text{Stroke}(c, c') + w_2 \text{Sijiao}(c, c') + w_3 \text{Edit}(c, c')$ where $\text{Stroke}(c, c')$ denotes the absolute value of strokes between c and c' . $\text{Sijiao}(c, c')$ is the Levenshtein distance of sijiao between c and c' , which is an encoding approach to fast retrieving Chinese characters. $\text{Edit}(c, c')$ is edit distance of character images between c and c' . The weights w_1, w_2 and w_3 is chosen as 0.33, 0.33 and 0.34, respectively.

158 sentences containing selected characters which do not cause syntactic errors after substituting the corresponding antonym character. Last, we generate the line image containing chosen sentences.

Evaluation metrics. To quantify perturbations of adversarial images x' compared with benign images x , we measure perturbations from MSE, PSNR and SSIM. *Mean-squared error* (MSE) denotes the difference between adversarial images and original images, calculated by $MSE = \frac{1}{|x|} \|x - x'\|_2$. *Peak-signal-to-noise ratio* (PSNR) is a ratio of maximum possible power of a signal and power of distortion, calculated by $PSNR = 10 \log \left(\frac{D^2}{MSE} \right)$ where D denotes the dynamic range of pixel intensities, e.g., for an 8 bits/pixel image we have $D = 255$. *Structural similarity index* (SSIM) attempts to model the structural information of an image from luminance, contrast and structure respectively.

To evaluate the efficiency of adversarial attacks, we calculate *attack success rate* (ASR) by $ASR = \frac{\#(f(x')=t) + \#(f(x) \neq f(x'))}{\#(x)}$ that is the fraction of adversarial images that were successful in fooling the DNN model, *targeted attack success rate* (ASR*) of adversarial attacks calculated as $ASR^* = \frac{\#(f(x')=t)}{\#(x)}$, the average time to generate adversarial perturbations from the clean images.

Comparison of attacks on single character altering

We compare different methods of altering a single character. Table 1 shows some successful adversarial examples that different attack methods generated. Our intuitions are: 1) MIM generates human-perceptible and unnatural noise on account of the dirty background, which distributes all over the image, and harms the image structure similarity and image quality. 2) WM and its variants retain the noise in the watermark region bringing in a more clear background and reasonable perturbations. 3) The watermark-fashion perturbations of WM are relatively light, and do not look like a real watermark. 4) WM_{init} and WM_{neg} look more real with darker and more complete watermark's shape. 5) The perturbations of WM_{edge} are around the edge of the text, which makes the text looks bolder and similar to printing/scanning defects.

Intuitively, we can see that WM family of attacks generate better visual quality (in terms of looking natural) images if the attack is successful.

We evaluate the attack performance of altering a single character using the corpus discussed above. In Table 2, we report the metrics above (MSE/PSNR/SSIM, ASR*/ASR), as well as the average time required to generate an attack sample. Our observations are: 1) From Table 2, compared to MIM, we can observe that WM, WM_{neg} and WM_{edge} obtained lower MSE, higher PSNR and higher SSIM, indicating that the noise level is indeed lower on a successful attack. 2) Due to the lower noise level, the attack success rates (ASR* and ASR) of WM and its variants are also lower than MIM's. We believe there are several reasons why they are lower. First, in this preliminary study, we always choose a fixed shape and location of the watermark that is at the center of the original image. The fixed location severely limits what the adversary can do. As our future work, we will allow multiple shapes

of the watermark (e.g., different texts, logos of the watermark), and different locations. 3) WM_{neg} casts away the positive gradient noise which possesses a certain attack ability. Hence, WM_{neg} behaved worse than WM. However, WM_{neg} does generate more natural examples visually. 4) WM_{edge} is a special case of WM which restricts the watermark region to the shape of the text edge, which has no location problem of the watermarks like other WM-series methods. It achieved good ASR and preserved the naturalness of perturbations. It's obvious to find that retaining all gradient noise is better than only keeping the negative gradient noise. 5) WM_0 is to attach an initial watermark to the original image, which can evaluate the impact of the watermark originally. After attaching an initial watermark, 37.9% and 17.1% images are misclassified by the DenseNet+CTC model and Tesseract OCR, respectively. Thus watermark owns attacking properties intrinsically. 6) The time for producing each adversarial example is similar and within a reasonable range. Thus, a practical strategy is to combine different attack methods to improve ASR.

Attack transferability to blackbox OCR

We want to see if adversarial examples can mislead other (black-box) models, or have commonly called *transferability* (Liu et al. 2016; Papernot, McDaniel, and Goodfellow 2016; Sharma and Chen 2017; Papernot et al. 2017).

We adopt the widely-used latest version Tesseract OCR (Google 2019) as a black-box model to perform adversarial attacks. We fed the adversarial samples, which are generated by attack methods above in the Densenet+CTC model, into the off-the-shelf Tesseract OCR, and evaluated recognition results (ASR*/ASR) shown in the last two columns of Table 2.

We find that all attacks produce transferable adversarial examples in terms of ASR. It may be due to the reason that the noise indeed perturbs the intrinsic features of a character sequence for different models, or because Tesseract OCR cannot handle noise. However, ASR* reduces significantly because perturbations are still trained on a different model.

Real-world examples

Figure 2 and 3 show two real-world examples of WM. In Figure 2, using watermarks, we successfully altered the license number recognition results. Figure 3 shows an example of a paragraph of an annual financial report of a public company. By adding the AICS watermark, we altered all the revenue numbers in the recognition results.

Defense against these attacks

We evaluate the robustness of these attacks against common defense methods that preprocess the input images, and Table 3 summarizes the results.

Noise removing methods with local smoothing. Local smoothing makes use of nearby pixels to smooth each pixel. We use three local smoothing methods from OpenCV (OpenCV b), *average blur*, *median blur* and *gaussian blur*. We observe that 1) median blur is particularly effective in removing sparsely-occurring black and white pixels in an

	DenseNet+CTC					Tesseract OCR		
	MSE	PSNR	SSIM	ASR*	ASR	Time (s)	ASR*	ASR
MIM	0.0102	32.27	84.92	92.4	93.7	20.8	19.6	84.2
WM	0.0020	34.70	96.03	60.8	61.4	17.5	19.0	88.0
WM_{neg}	0.0023	34.26	94.11	52.5	53.2	20.7	19.6	88.6
WM_{init}	0.0094	34.64	89.02	55.1	71.5	19.9	0.0	100.0
WM_{edge}	0.0058	34.43	93.87	87.3	88.0	13.5	19.6	83.5
WM_o	0.0034	30.92	93.37		37.9			17.1

Table 2: Comparison of varying adversarial attacks on DenseNet+CTC and Tesseract OCR. See text for the metric definitions.

		ASR*/ASR	MIM	WM	WM _{neg}	WM _{init}	WM _{edge}	example
no deformation			92.4/93.7	60.8/61.4	52.5/53.2	55.1/71.5	87.3/88.0	靠左行驶
Deformation	AvgBlur@2x2		44.30/62.7	47.47/51.27	46.84/54.43	18.99/79.11	48.73/55.06	靠左行驶
	MedianBlur@3x3		2.53/99.37	1.27/98.73	1.90/95.57	0.00/100.0	2.53/96.84	靠左行驶
	GaussianBlur@3x3,0		51.27/58.86	48.10/52.53	46.20/55.70	27.85/72.78	52.53/55.06	靠左行驶
	Salt&Pepper@2%		4.43/99.37	3.16/98.10	2.53/97.47	0.00/100.0	3.80/98.73	靠左行驶
	Compress@20		38.61/68.99	46.20/55.06	47.47/58.86	22.78/76.58	56.33/63.92	靠左行驶
	inpainting@2			0.00/100.0	0.00/100.0	0.00/100.0	0.00/100.0	靠左行驶

Table 3: Comparison of ASR*/ASR (%) of adversarial examples under different preprocessing / defense methods.



Figure 2: An adversarial attack example in driver license recognition. The OCR output a licenses number of *NAL12505717* while it is actually *NHL12506717*.

image while preserving edges of the text well. 2) Median blur with kernel size 3 blurs texts so much that OCR algorithms no longer work, although it reduces ASR*. 3) Average smoothing with kernel size 2 and Gaussian smoothing with kernel size 3 have similar performance. Although MIM has a high ASR, it seems more sensitive to various deformations than WM and WM_{edge}.

Salt&pepper noise is a common way against adversarial examples. We find that it is particularly effective in decreasing ASR* (to 0), but the result is an increase of ASR to almost 100%. It indicates that salt&pepper noise harms the general image quality too much in exchange for reducing adversarial perturbations.

Image compression. We show that adversarial examples can survive lossy image compression and decompression, applying the standard JPEG compression method, with the quality parameter is set to 20. We can point out that WATERMARK attack has more chance of surviving under the compression process.

Watermark removing techniques. Inpainting is a commonly used method to remove (real) watermarks. We use the

Input Image:

2018年度公司实现营业收入 70,599.25 万元, 同比减少 43,695.78 万元, 减幅 38.23%, 其中: 酒类收入 31,594.15 万元, 同比减少 3.54%, 略有下降; 房地产收入 3,651.21 万元, 同比减少 64.42%, 主要是本期受到全岛限购政策影响和存量房较少以致收入减少; 贸易收入 28,140.19 万元, 同比减少 57.99%; 饮料收入 6,483.03 万元, 同比增加 3.56%, 略有增长。

OCR Output:

2018年度公司实现营业收入 70,599.95 万元, 同比减少 43,695.78 万元, 减幅 38.23%, 其中: 酒类收入 81,594.15 万元, 同比减少 3.54%, 略有下降; 房地产收入 3,651.21 万元, 同比减少 74.42%, 主要是本期受到全岛限购政策影响和存量房较少以致收入减少; 贸易收入 98,140.19 万元, 同比减少 57.99%; 饮料收入 6,483.03 万元, 同比增加 3.60%, 略有增长。

Figure 3: Attack on a listed Chinese company’s annual report. All the revenue numbers are altered in the OCR result.

inpainting method (Telea 2004) implemented in OpenCV (OpenCV a), which required the mask of watermarks as a priori and tried to recover the watermark region according to surrounding pixels. While inpainting eliminated the watermark, because the text region overlapped with the watermark region, the inpainting method removed too many useful pixels, that is text pixels, causing OCR to fail completely. We observe that the text even lost readability to human eyes.

Conclusion and Future Work

Generating adversarial examples for OCR systems is different from normal CV tasks. We propose a method that successfully hides perturbations from human eyes while making them effective in the modern sequence-based OCR, by pretending perturbations as a watermark, or printing defects. We show that even with preliminary implementations, our perturbations can be still effective, transferable, and deceiv-

ing to human eyes.

There are many future directions. For example, allowing different watermark shapes and locations, as well as on longer sequences. Also, it is interesting to add semantic-based (language model) attacks even further to improve the attack effectiveness. Also, the adversarial attack calls for better defense methods other than traditional image transformations.

Acknowledgements. 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

- [Bengio et al. 1995] Bengio, Y.; LeCun, Y.; Nohl, C.; and Burges, C. 1995. Lerec: A nn/hmm hybrid for on-line handwriting recognition. *Neural Computation* 7(6):1289–1303.
- [Breuel et al. 2013] Breuel, T. M.; Ul-Hasan, A.; Al-Azawi, M. A.; and Shafait, F. 2013. High-performance ocr for printed english and fraktur using lstm networks. In *2013 12th International Conference on Document Analysis and Recognition*, 683–687. IEEE.
- [Brown et al. 2017] Brown, T. B.; Mané, D.; Roy, A.; Abadi, M.; and Gilmer, J. 2017. Adversarial patch. *arXiv preprint arXiv:1712.09665*.
- [Carlini and Wagner 2017] Carlini, N., and Wagner, D. 2017. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, 39–57. IEEE.
- [Chen et al. 2018] Chen, P.-Y.; Sharma, Y.; Zhang, H.; Yi, J.; and Hsieh, C.-J. 2018. Ead: elastic-net attacks to deep neural networks via adversarial examples. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [Dong et al. 2018] Dong, Y.; Liao, F.; Pang, T.; Su, H.; Zhu, J.; Hu, X.; and Li, J. 2018. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9185–9193.
- [España-Boquera et al. 2011] España-Boquera, S.; Castro-Bleda, M. J.; Gorbe-Moya, J.; and Zamora-Martinez, F. 2011. Improving offline handwritten text recognition with hybrid hmm/ann models. *IEEE transactions on pattern analysis and machine intelligence* 33(4):767–779.
- [Feinman et al. 2017] Feinman, R.; Curtin, R. R.; Shintre, S.; and Gardner, A. B. 2017. Detecting adversarial samples from artifacts. *arXiv preprint arXiv:1703.00410*.
- [Goodfellow, Shlens, and Szegedy 2014] Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- [Google 2019] Google. 2019. Tesseract.
- [Graves et al. 2006] Graves, A.; Fernández, S.; Gomez, F.; and Schmidhuber, J. 2006. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, 369–376. ACM.
- [Grosse et al. 2017] Grosse, K.; Manoharan, P.; Papernot, N.; Backes, M.; and McDaniel, P. 2017. On the (statistical) detection of adversarial examples. *arXiv preprint arXiv:1702.06280*.
- [Hanwei Zhang 2019] Hanwei Zhang, Yannis Avrithis, T. F. L. A. 2019. Smooth adversarial examples. *arXiv preprint arXiv:1903.11862*.
- [Heng, Zhou, and Jiang 2018] Heng, W.; Zhou, S.; and Jiang, T. 2018. Harmonic adversarial attack method. *arXiv preprint arXiv:1807.10590*.
- [Hinton, Vinyals, and Dean 2015] Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- [Huang et al. 2017] Huang, G.; Liu, Z.; Van Der Maaten, L.; and Weinberger, K. Q. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4700–4708.
- [Karmon, Zoran, and Goldberg 2018] Karmon, D.; Zoran, D.; and Goldberg, Y. 2018. Lavan: Localized and visible adversarial noise. *arXiv preprint arXiv:1801.02608*.
- [Kurakin, Goodfellow, and Bengio 2016] Kurakin, A.; Goodfellow, I.; and Bengio, S. 2016. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*.
- [Liu et al. 2016] Liu, Y.; Chen, X.; Liu, C.; and Song, D. 2016. Delving into transferable adversarial examples and black-box attacks. *arXiv preprint arXiv:1611.02770*.
- [Lu, Issaranon, and Forsyth 2017] Lu, J.; Issaranon, T.; and Forsyth, D. 2017. Safetynet: Detecting and rejecting adversarial examples robustly. In *Proceedings of the IEEE International Conference on Computer Vision*, 446–454.
- [Madry et al. 2017] Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.
- [Moosavi-Dezfooli, Fawzi, and Frossard 2016] Moosavi-Dezfooli, S.-M.; Fawzi, A.; and Frossard, P. 2016. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2574–2582.
- [Nguyen, Yosinski, and Clune 2015] Nguyen, A.; Yosinski, J.; and Clune, J. 2015. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 427–436.
- [OpenCV a] OpenCV. inpainting(cv.inpaint). <https://docs.opencv.org/3.0-beta/modules/photo/doc/inpainting.html>.
- [OpenCV b] OpenCV. Median filter(cv.medianblur). https://docs.opencv.org/master/d4/d13/tutorial_py_filtering.html.
- [Papernot et al. 2016a] Papernot, N.; McDaniel, P.; Jha, S.; Fredrikson, M.; Celik, Z. B.; and Swami, A. 2016a. The limitations of deep learning in adversarial settings. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*, 372–387. IEEE.
- [Papernot et al. 2016b] Papernot, N.; McDaniel, P.; Wu, X.; Jha, S.; and Swami, A. 2016b. Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE Symposium on Security and Privacy (SP)*, 582–597. IEEE.
- [Papernot et al. 2017] Papernot, N.; McDaniel, P.; Goodfellow, I.; Jha, S.; Celik, Z. B.; and Swami, A. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, 506–519. ACM.
- [Papernot, McDaniel, and Goodfellow 2016] Papernot, N.; McDaniel, P.; and Goodfellow, I. 2016. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*.
- [Sharma and Chen 2017] Sharma, Y., and Chen, P.-Y. 2017. Attacking the madry defense model with l_1 -based adversarial examples. *arXiv preprint arXiv:1710.10733*.

- [Smith 2007] Smith, R. 2007. An overview of the tesseract ocr engine. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, volume 2, 629–633. IEEE.
- [Szegedy et al. 2013] Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- [Telea 2004] Telea, A. 2004. An image inpainting technique based on the fast marching method. *Journal of graphics tools* 9(1):23–34.
- [Tramèr et al. 2017] Tramèr, F.; Kurakin, A.; Papernot, N.; Goodfellow, I.; Boneh, D.; and McDaniel, P. 2017. Ensemble adversarial training: Attacks and defenses. *arXiv preprint arXiv:1705.07204*.
- [Wang et al. 2012] Wang, T.; Wu, D. J.; Coates, A.; and Ng, A. Y. 2012. End-to-end text recognition with convolutional neural networks. In *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, 3304–3308. IEEE.
- [Xu, Evans, and Qi 2017] Xu, W.; Evans, D.; and Qi, Y. 2017. Feature squeezing: Detecting adversarial examples in deep neural networks.