# Vulnerability Analysis of Transformer-based Optical Character Recognition to Adversarial Attacks

Lucas Beerens[*]        Desmond J. Higham[†]

November 30, 2023

## Abstract

Recent advancements in Optical Character Recognition (OCR) have been driven by transformer-based models. OCR systems are critical in numerous high-stakes domains, yet their vulnerability to adversarial attack remains largely uncharted territory, raising concerns about security and compliance with emerging AI regulations. In this work we present a novel framework to assess the resilience of Transformer-based OCR (TrOCR) models. We develop and assess algorithms for both targeted and untargeted attacks. For the untargeted case, we measure the Character Error Rate (CER), while for the targeted case we use the success ratio. We find that TrOCR is highly vulnerable to untargeted attacks and somewhat less vulnerable to targeted attacks. On a benchmark handwriting data set, untargeted attacks can cause a CER of more than 1 without being noticeable to the eye. With a similar perturbation size, targeted attacks can lead to success rates of around 25%—here we attacked single tokens, requiring TrOCR to output the tenth most likely token from a large vocabulary.

[*]School of Mathematics and The Maxwell Institute for Mathematical Sciences, University of Edinburgh, EH8 9BT, UK

[†]School of Mathematics and The Maxwell Institute for Mathematical Sciences, University of Edinburgh, EH8 9BT, UK

A robot may not injure a human being

$+\Delta x$

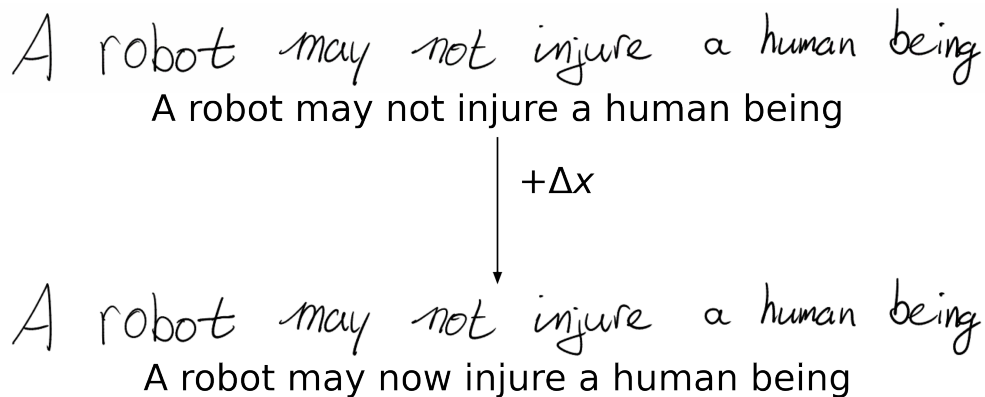A robot may now injure a human being

Figure 1: Adversarial attack created by tailoring a Carlini & Wagner attack [9] to the TrOCR setting. The full algorithm is described in section 3.5. Here a targeted version of the attack is used—we aim for the optical character recognition system to create a specific output on the perturbed sentence. The original sentence (upper), handwritten by the first author and scanned, was correctly recognized by TrOCR as 'A robot may not injure a human being'. After the imperceptible adversarial perturbation is added (lower), TrOCR incorrectly outputs 'A robot may now injure a human being'. This was the output target supplied to the attack algorithm. The original sentence is inspired by the first law of robotics by Isaac Asimov [3].

# 1    Introduction

Optical Character Recognition (OCR) is the conversion of images of printed or handwritten text into machine-interpretable text. This has many applications, ranging from document digitization and automated data entry to facilitating accessibility for visually impaired individuals. OCR has gone through several decades of development [27]. In the last decade, the resurgence of neural networks has led to significant improvements in OCR capabilities. Both text detection and text recognition were able to take advantage of the advances in convolutional neural networks (CNNs). In particular, text recognition models have been formulated as encoder-decoder systems, where the encoder uses CNNs, and the decoder uses recurrent neural networks (RNNs) [39].

In recent years, the widespread adoption of transformers [44], originally

designed for natural language processing (NLP), has yielded remarkable advances in various domains, including NLP [17, 54, 1], Computer Vision (CV) [18, 20], and speech processing [21]. Utilization of transformers in NLP and CV has now been extended to OCR. First, this was done while keeping CNNs in the backbone [13]. Later, a model without CNN backbone was created called Transformer-based Optical Character Recognition (TrOCR) [22].

However, amidst these achievements, a critical concern remains. Many AI systems can be manipulated by mischievous, malicious or criminal third parties. Adversarial attacks, which, for example, introduce imperceptible perturbations into clean images to deceive deep learning classification models, have drawn substantial attention in the past decade [41, 15, 23]. A cat-and-mouse game between attack and defence strategies is continuing, with Carlini [8] observing that

> "Historically, the vast majority of adversarial defenses published at top-tier conferences . . . are quickly broken."

Given that recent TrOCR technology offers benefits in many high-risk and safety-critical areas, including legal, financial and healthcare applications, it is crucial to understand its vulnerability to adversarial attack. This issue, which currently remains unaddressed, motivates our work.

At a policy-making level, concerns about AI in general have led to a call for regulation to ensure the safety and transparency of AI. An example of such regulation is the proposed EU AI act, which was first released in April 2021 [11] and amended in June 2023 [32]. Article 15 – paragraph 4 – subparagraph 1 of the amended proposal states:

> "High-risk AI systems shall be resilient as regards to attempts by unauthorised third parties to alter their use, behaviour, outputs or performance by exploiting the system vulnerabilities."

Highlighting one of many application domains, we note that Annex III – paragraph 1 – point 3 – point b of this same proposal mentions that the following class of AI systems is deemed high risk:

> "AI systems intended to be used for the purpose of assessing students in educational and vocational training institutions and for assessing participants in tests commonly required for admission to those institutions."

OCR systems can be used as part of an AI pipeline designed for student assessment [34]. Therefore, in this and other high-risk domains, understanding and addressing their weaknesses is imperative before deployment.

Our main contributions in this work are as follows.

- We create a novel framework to assess the resilience of Transformer-based Optical Character Recognition.

- We tailor strategies to TrOCR by building on established image classification attacks, resulting in attacks such as the one in fig. 1.

- We assess the vulnerability of TrOCR to both targeted and untargeted attacks on the IAM handwriting dataset by comparing multiple attack algorithms.

# 2 Related Work

## 2.1 Adversarial Attacks on CNNs

The first adversarial attacks were carried out on fully connected neural networks and CNNs in the context of image classification [41, 15]. Subsequently, these ideas have been extended to other areas, including NLP, and other architectures within CV [46, 53, 48].

Adversarial attacks can be classified as white-box or black-box. The first adversarial attacks fell into the white-box category: they had access to information about the neural network, including gradients with respect to input components. Examples of widely-used white-box methods include L-BFGS [41], FGSM [15], and DeepFool [28]. These techniques are readily generalized to black-box attacks, where less information about the targeted neural network is available [30]. A common approach in black-box attacks is to use gradient-based white-box attack algorithms optimized for transferability, replacing exact gradients with approximations. Two such types of attack are transfer-based and query-based. In the transfer-based case, the gradient of a surrogate model is used [30]. In the query-based case, the gradient is estimated by querying the model with a number of inputs [10, 4]. Hence, white-box methods are relevant to the black-box context.

Many strategies have been employed in the development of white-box techniques. These can be broadly categorized as gradient optimization, constrained optimization, and generative models. Within these categories, a

4

targeted attack aims for a specific, new output class, whereas a non-targeted attack aims for any change of output class.

Since attacks typically aim for the largest output change for a given (small) size of input change, the norm used to measure perturbations is a key ingredient. Most attacks use $\ell_0$, $\ell_2$, or $\ell_\infty$ norms [38]. When the objective is to hide a perturbation in the image, these norms can be interpreted as proxies for the visibility of the attack. There are also other forms of disguise, such as attacking speech recognition systems using psychoacoustic hiding [36] or attacking handwritten digit classification by hiding the perturbation in the ink [5]. In this work, we will be using the $\ell_2$ norm to measure perturbation size.

## 2.2 Optical Character Recognition

Attacks on CNN models have focused mainly on classification. Here, we consider the TrOCR model, which is designed for Optical Character Recognition. Two types of input can be considered: handwritten and printed text. Printed text is the simpler case due to the uniformity of dimensions, positions, and fonts. We focus on handwritten text. Handwritten OCR may also be divided into two cases: online and offline. In the former case, the system can use all information that can be captured while the text is written. This means that the system knows the direction and order of the pen strokes while not having any issues with stroke thickness. In the offline case, this information is not available and only the resulting image can be used. We focus on the more widely applicable offline case.

Typically, OCR refers to the entire pipeline from data acquisition to post-processing. This pipeline can be separated into the following steps: acquisition, pre-processing, segmentation, feature extraction, classification, and post-processing. The first three steps of the text processing pipeline can be grouped together as text detection, while the last three can be grouped together as text recognition. We focus on segmented examples from the IAM Handwriting Database [26]. (Figure 1 uses a scan of a sentence that was handwritten by the first author.)

## 2.3 Transformer-based OCR

TrOCR is an end-to-end text recognition system, which takes textline segmented images and performs text recognition [22]. It does not use any con-

volutional layers in its architecture. Instead, it uses the transformer architecture [44]. TrOCR consists of three steps. First the image is resized to a fixed resolution and split into a sequence of image patches. These are used as input to a pre-trained Vision Transformer model, which encodes this sequence of patches [14]. This is decoded by a pre-trained transformer language model, which generates wordpiece tokens based on the image and the context generated before. TrOCR was pre-trained with synthetic data and fine-tuned with human-labeled datasets. Different combinations of encoder and decoder architectures are used to create versions of different sizes. $\text{TrOCR}_{\text{SMALL}}$ uses $\text{DeiT}_{\text{SMALL}}$ for the decoder [42] and MiniLM for the encoder [45], giving a total of 62M parameters, compared with 334M and 558M parameters, respectively, for $\text{TrOCR}_{\text{BASE}}$ and $\text{TrOCR}_{\text{LARGE}}$. We focus on $\text{TrOCR}_{\text{SMALL}}$ which allows us to conduct more extensive experiments; this choice is supported by empirical [31, 50, 35] and theoretical [43] results showing that smaller models are more resilient to attack, in general.

## 2.4   Adversarial attacks on transformer models

The rise of transformers has led to investigations into their vulnerability to adversarial attacks. In particular, Visual Transformers (ViTs) have been examined, with results suggesting that they are generally more robust than CNNs [16, 37, 2]. However, the attacks originally designed for CNN systems remain somewhat effective against ViTs and some recent work has focused on improving transferability [47, 25, 52].

Transformers have also been used in NLP to build language models such as BERT [12] and MiniLM [45]. Such models can be used to create systems for a variety of applications, such as machine translation, sentiment analysis, and summarization. Models that generate tokens are the most relevant to our work, as they are broadly similar to the TrOCR decoder. Adversarial attacks have been found that generate effective adversarial examples for these types of models [29].

Since TrOCR has image inputs, like the ViTs, but has generative text output, like language transformer models, it is substantially different from both cases. Adversarial attacks on TrOCR have been investigated in the context of adding diacritics to printed text [7]. This is markedly different from our work since it specifically uses printed text, and the perturbations are within the ASCII character set. To the best of our knowledge, adversarial attacks on models such as TrOCR have not been investigated before in the

6

context of small perturbation to input images.

# 3   Method

In section 3.1 we present a framework for dealing with TrOCR mathematically. Then we introduce adversarial attacks against TrOCR based on various existing methods in sections 3.2 to 3.5. Finally, we discuss the implementation in section 3.6. section 4 presents experimental results and section 5 gives conclusions.

## 3.1   TrOCR

In creating adversarial attacks against TrOCR, we will be using gradients of functions applied to the TrOCR output. Therefore, it is important to describe this output carefully. Images are represented by a real-valued vector $x$. In the inference process, images are first passed through a preprocessor to give them a similar shape and distribution. Then, they are passed through an encoder that encodes image patches into vectors. Let us call the output of these first two steps $E(x)$. The decoder will iteratively generate tokens from this. It starts with a special beginning-of-sentence (bos) token. In every iteration, the output is the current sequence of tokens shifted to the left (so without the bos token), with the new token added at the end. Once the end-of-sequence (eos) token has been generated or the maximum amount of tokens has been reached, the process stops. The output is a sequence of logit vectors. The elements in the vectors correspond to vocabulary tokens. Applying the softmax function would result in probabilities, but we do not need to consider such a step. The element with the maximum value determines the token assigned. The inputs and outputs have the form

$$(E(x), ((\text{bos}), t_1, \ldots, t_k)) \mapsto (\mathbf{t}_1, \ldots, \mathbf{t}_{k+1}), \tag{1}$$

where for all $i$ in $\{1, \ldots k + 1\}$ we have

$$\mathbf{t}_i = \begin{pmatrix} l_{i,1} \\ \vdots \\ l_{i,v} \end{pmatrix}, \tag{2}$$

and $t_i = \arg\max_j \{l_{i,j}\}$. Here, $v$ is the size of the vocabulary. Vector $\mathbf{t}_i$ does not change in subsequent iterations. This means that the logits for the

7

generation of all tokens can be obtained by taking the final output. The final logit matrix is denoted by $F(x)$ with entries $F_j^i(x)$, where $i$ is the position of the token and $j$ is the token class in the vocabulary.

## 3.2 FGSM

The Fast Gradient Sign Method (FGSM) [15] is an early and widely used adversarial attack strategy. Let us first consider the untargeted version of FGSM. The gradient of the loss function with respect to the correct label points in the direction where it will increase the most, locally. FGSM looks at the sign of this gradient for every pixel and takes a step in that direction. Increasing the loss means moving away from the correct answer. The attack is computed as

$$\Delta x = \epsilon \operatorname{sign}\left(\nabla_x \mathcal{L}(x, L_{\text{out}})\right), \tag{3}$$

where $\mathcal{L}$ denotes the loss function used in training with respect to the output labels $L_{\text{out}} = (t_0, \ldots, t_k)$ and $\epsilon$ denotes some allowed perturbation size. For small $\epsilon$, the perturbation $\Delta x$ in (3) causes the largest change under the constraint $\|\Delta x\|_\infty \leq \epsilon$.

FGSM can be used in a targeted manner by instead computing the loss with respect to some target labels $L_{\text{target}}$. The targeted attack can be written as

$$\Delta x = -\epsilon \operatorname{sign}\left(\nabla_x \mathcal{L}(x, L_{\text{target}})\right). \tag{4}$$

Both versions can be applied to TrOCR in a straightforward manner. As a loss function, we take the sum of the cross-entropy losses of the tokens, which was used in training. Then we can use backpropagation to compute the gradient with respect to the pixel values.

## 3.3 DeepFool

The DeepFool algorithm [28] was developed as an extension to FGSM. It is an untargeted attack algorithm. The key idea is that, at each iteration, the decision boundaries are approximated by means of linearization. This creates a polyhedron. Then, the closest boundary is chosen and a perturbation to that boundary is made. This means that in each step we aim for the class to change to one other specific class, characterized by the smallest alteration. Iteration stops once the class changes. If we start with an image $x_0$ classified

8

as $\hat{k}$, then for each iteration we first compute the following values for all classes:

$$w_k \leftarrow \nabla F_k(x_i) - \nabla F_{\hat{k}}, \tag{5}$$

$$f_k \leftarrow F_k(x_i) - F_{\hat{k}}(x_i). \tag{6}$$

Then the optimal class to target is

$$\overline{k} \leftarrow \arg\min_{k \neq \hat{k}} \frac{|f_k|}{\|w_k\|_2}, \tag{7}$$

and the resulting update is

$$x_{i+1} \leftarrow x_i + \frac{|f_{\overline{k}}|}{\|w_{\overline{k}}\|_2^2} w_{\overline{k}}. \tag{8}$$

In applying DeepFool to TrOCR, we need to consider decision boundaries for all tokens. We aim for all tokens to be mislabeled. Therefore, we could compute updates for every token consecutively. Because the number of possible tokens in the vocabulary is high, a maximum number of decision boundaries to check is chosen, using the classes with the highest logits. The overall method is given in algorithm 1.

In practice, it was found that a single iteration is sufficient. Almost always, the closest decision boundary was that of the second most likely class. Therefore in testing we use topAmount $= 1$ and maxIt $= 1$ in algorithm 1.

## 3.4   Backward Error

Adversarial attacks based on backward error analysis were first introduced in [6] and extended in [5]. The idea is that at each iteration, the neural network is linearized around the current image, and a linear least-squares constrained optimization problem is solved to find the smallest perturbation leading to the desired misclassification. This optimization problem may then be solved by a state of the art solver; we use OSQP [40].

To extend this to TrOCR, we may change the misclassification constraint to hold for all tokens that do not align with the targets. Instead of comparing the target label with all the options in the vocabulary, we will only compare with the current label. This is done to limit the amount of backpropagation. Finally, we also remove the constraint keeping the pixel values in $[0, 1]$, to be more closely comparable with the other algorithms. The method is given in algorithm 2.

9

---

**Algorithm 1** DeepFool for TrOCR

---

1: **Input:** $F, x, \text{topAmount}, \text{maxIt}$
2: **Output:** $\Delta x$
3: $\text{out} \leftarrow F(x)$
4: $\text{orgLabels}, \text{labels} \leftarrow \text{Labels}(\text{out})$
5: $\text{len} \leftarrow \text{Len}(\text{orgLabels})$
6: **initialize** $x_1 \leftarrow x$ **and** $i \leftarrow 1$
7: **while** $\text{Any}_j \left( \text{orgLabels}_j = \text{labels}_j \right)$ **and** $i \leq \text{maxIt}$ **do**
8:     **for** $t \in \{1, \ldots, \text{len}\}$ **with** $\text{orgLabels}_t = \text{labels}_t$ **do**
9:         $\hat{k} \leftarrow \text{orgLabels}_t$
10:         $\text{topLabels} \leftarrow \text{argTopK}_k \left( \text{out}_k^t, \text{topAmount} + 1 \right)$
11:         **for** $k \in \text{topLabels}$ **with** $k \neq \text{orgLabels}_t$ **do**
12:             $w_k^t \leftarrow \nabla F_k^t(x_i) - \nabla F_{\hat{k}}^t(x_i)$
13:             $f_k^t \leftarrow F_k^t(x_i) - F_{\hat{k}}^t(x_i)$
14:         **end for**
15:         $\overline{k}_t \leftarrow \arg\min_{k \neq \hat{k}} \frac{|f_k^t|}{\|w_k^t\|_2}$
16:         $r_i^t \leftarrow \frac{|f_{\overline{k}}|}{\|w_{\overline{k}}\|_2^2} w_{\overline{k}}$
17:     **end for**
18:     $u_i \leftarrow \sum_t r_i^t$
19:     $x_{i+1} \leftarrow x_i + u_i$
20:     $\text{out} \leftarrow F(x_{i+1})$
21:     $\text{labels} \leftarrow \text{Labels}(\text{out})$
22:     $i \leftarrow i + 1$
23: **end while**
24: $\Delta x = \sum_i u_i$
25: **return** $\Delta x$

---

**Algorithm 2** Backward Error attack for TrOCR
_____

1: **Input:** $F, x, \alpha$, target, iterations
2: **Output:** $\Delta x$
3: out $\leftarrow F(x)$
4: labels $\leftarrow$ Labels(out)
5: $x_1 \leftarrow x$
6: $\Delta x \leftarrow 0$
7: diff $\leftarrow \{t : \text{target}_t \neq \text{labels}_t\}$
8: **for** $i = 1$ **to** iterations **do**
9:     $\delta x \leftarrow$ Variable
10:     $y_{\text{orig}} \leftarrow$ Variable
11:     $y_{\text{targ}} \leftarrow$ Variable
12:     objective $\leftarrow$ Objective$(\|\Delta x + \delta x\|_2)$
13:     constr $\leftarrow$ Constraint$(y_{\text{orig}} \leq y_{\text{targ}})$
14:     **for** $j = 1$ **to** Len(diff) **do**
15:         $t \leftarrow \text{diff}_j$
16:         orig $\leftarrow \text{labels}_t$
17:         targ $\leftarrow \text{target}_t$
18:         constr.add$\left(y_{\text{orig}}^j = \text{out}_{\text{orig}}^t + \nabla_x F_{\text{orig}}^t(x_i) \cdot \delta x\right)$
19:         constr.add$\left(y_{\text{targ}}^j = \text{out}_{\text{targ}}^t + \nabla_x F_{\text{targ}}^t(x_i) \cdot \delta x\right)$
20:     **end for**
21:     $\delta x, y_{\text{orig}}, y_{\text{targ}} \leftarrow$ Solve(objective, constr)
22:     $\Delta x \leftarrow \Delta x + \alpha \delta x$
23:     $x_{i+1} \leftarrow x_i + \alpha \delta x$
24:     out $\leftarrow F(x_{i+1})$
25:     labels $\leftarrow$ Labels(out)
26: **end for**
27: **return** $\Delta x$
_____

## 3.5 Carlini and Wagner

In contrast to FGSM, where the training loss function is used, the Carlini and Wagner (C&W) attack [9] is based on the idea of minimizing a specially created loss function. It can be used in targeted and untargeted situations. When targeting a label $\hat{k}$, the new loss function is defined as

$$\mathcal{L}(\Delta x) = \|\Delta x\|_2^2 + c \, f\left(x + \Delta x, \hat{k}\right), \tag{9}$$

where the first term on the RHS controls the $\ell_2$ norm of the perturbation and the second term encourages the desired classification. The constant $c > 0$ is a parameter that balances the two requirements. The function $f$ is chosen so that $f(x + \Delta x, \hat{k}) \leq 0$ if and only if $\arg\max_k F_k(x + \Delta x) = \hat{k}$. In particular, it is chosen to be

$$f\left(x', \hat{k}\right) = \left(\max_{k \neq \hat{k}} \left(F_k(x')\right) - F_{\hat{k}}(x')\right)^+, \tag{10}$$

where $(\cdot)^+ = \max\{0, \cdot\}$. In that case $\hat{k}$ is the target class. In the untargeted setting, we take $\hat{k}$ as the predicted class and instead define

$$f\left(x', \hat{k}\right) = \left(F_{\hat{k}}(x') - \max_{k \neq \hat{k}} \left(F_k(x')\right)\right)^+. \tag{11}$$

Then we have $f(x + \Delta x, \hat{k}) \leq 0$ if and only if $\arg\max_k F_k(x + \Delta x) \neq \hat{k}$.

Pixels are represented by a change of variables, keeping them in the allowed range $[0, 1]$. We represent perturbed image $x + \Delta x$ with values $w$ such that

$$x + \Delta x = \frac{1}{2}\left(\tanh(w) + 1\right). \tag{12}$$

Now the optimization problem becomes

$$\underset{w}{\arg\min} \left(\begin{array}{c} \left\|\frac{1}{2}\left(\tanh(w) + 1\right) - x\right\|_2^2 \\ + c \, f\left(\frac{1}{2}\left(\tanh(w) + 1\right), \hat{k}\right) \end{array}\right), \tag{13}$$

with the choice of $f$ depending on whether the attack is targeted. For each iteration, the loss function gradient is calculated and a step is taken based on the Adam optimization scheme [19].

With a suitable adjustment to $f$, we can use this approach against TrOCR. Now we have multiple tokens, so in the targeted case the function becomes

$$f\left(x', L\right) = \sum_t \left( \max_{k \neq L^t} \left( F_k^t(x') \right) - F_{L^t}^t(x') \right)^+, \tag{14}$$

where $L$ are the output labels. In the untargeted case the function will become

$$f\left(x', L\right) = \sum_t \left( F_{L^t}^t(x') - \max_{k \neq L^t} \left( F_k^t(x') \right) \right)^+, \tag{15}$$

where $L$ are the target labels.

To limit the number of iterations, a maximum is set. Additionally, when the value of the loss function has not improved for five steps, the optimization process is terminated, and the perturbation with the lowest loss so far will be used as the output. In practice, when the targeted version was used, it was found that the loss tends to increase again after reaching the desired target. Therefore, we terminate immediately after reaching the target.

## 3.6 Implementation

Some comments on the implementation are in order. All testing is done using PyTorch [33]. The TrOCR model is accessed through the Transformers package by Hugging Face [49]. The image preprocessor included with TrOCR does not support backpropagation. To overcome this, we wrote an equivalent function that allows backpropagation. Similarly, the function used for inference does not support backpropagation, so we performed the inference steps manually.

# 4 Experiments

## 4.1 Evaluation Metrics

The success of attacks in targeted and untargeted settings may be measured in several ways. In an untargeted setting, we simply wish to quantify the effect on performance. A standard benchmark for TrOCR is the IAM handwriting database test set [26]. This testing set consists of 1861 images of text written by 128 different writers. In line with a previous study of TrOCR

[22], we will score performance using the Character Error Rate (CER), with a smaller CER indicating better model performance. The attacks aim to make the CER higher. This measure is based on the number of substitutions, deletions, and insertions needed to recover the correct answer. Here, we define the correct answer to be the true label given by the dataset. We will report this performance score in terms of the perturbation size.

In the targeted setting, we will look at the percentage of the attacked images that are classified with the target labels.

## 4.2 Experimental Setup

In the untargeted setting, we consider FGSM, DeepFool, and C&W attacks. The resulting attacks are then evaluated at the same relative scales; that is, for any attack $\Delta x$ on an image $x$, we evaluate

$$F\left(x + \epsilon\frac{\|x\|_2}{\|\Delta x\|_2}\Delta x\right),\tag{16}$$

with values $\epsilon \in \{n/10^4 : n = 0, 1, \ldots, 99\}$. Then we find the labels with the highest logit scores and decode the resulting sequence into words. These can then be compared with the original decoded labels using the CER. This will give us graphs of the CER on those perturbation sizes for every algorithm.

FGSM does not require any parameters to be set. DeepFool is used with a single iteration and with a single class other than the original. C&W uses a learning rate and weight decay for Adam of 0.002 and $10^{-5}$ respectively. These values were experimentally fine-tuned from the default values. The value of $c$ is 0.05, based on the recommendation in [9] to choose a low $c$ that leads to misclassification. As initialization, we use a perturbation that changes every pixel by $\eta = 0.00002$ in a random direction. This is then clipped, if necessary, to stay within the range $[0, 1]$. To accommodate the change of variables, we then change every pixel at 0 or 1 to $\eta$ and $1 - \eta$ respectively. The recommendation in [9] that these initial perturbations should be about the size of the eventual perturbation inspired this choice of $\eta$. A maximum number of 30 iterations is used. This is done because in the untargeted case the algorithm often easily finds perturbations leading to the desired misclassification, and then spends more iterations optimizing the perturbation size.

In the targeted case, we compare FGSM, Backward Error, and C&W attacks. For every image, we create a perturbation in which we aim to change

14

a random token to the token in the vocabulary with the tenth highest logit score. We regard this as a fair comparison because it changes tokens which offer the same typical difficulty.

Here we again evaluate the perturbations as in eq. (16), but now for $\epsilon \in \{n/10^5 : n = 0, 1, \ldots, 99, 100, 110, \ldots, 390\}$. We look at smaller perturbations than in the untargeted case, because generally the perturbations made by the algorithms did not become more successful when scaled to be larger. Having finer resolution allows us to more accurately find all successful attacks. More $\epsilon$ are added for a better comparison with the C&W attack. We will view the perturbations created by this attack separately and evaluate them in the size in which they are outputted by the algorithm. This is because scaling these leads to unsuccessful attacks.

Again, we decode the labels from the outputs into words. However, this time we will not look at the average CER with respect to the original labels. Instead, we compute what percentage has a CER of 0 with respect to the target labels for some perturbation up to size $\epsilon$. That is, for every $n_0 \in \{0, 1, \ldots, 99, 100, 110, \ldots, 390\}$ we compute

$$\frac{1}{I} \sum_i \mathbf{1} \left( \exists_{n \leq n_0} : \mathrm{CER}_{i,n} = 0 \right), \tag{17}$$

where $I$ is the total number of images, $\mathbf{1}$ indicates whether a statement is true and

$$\mathrm{CER}_{i,n} = \mathrm{CER} \left( \mathrm{pertText}_i(n/10^5), \mathrm{targetText}_i \right), \tag{18}$$

is the CER between target text of image $i$ and the decoded text output from image $i$ with a perturbation of $\epsilon = n/10^5$.

The targeted FGSM algorithm does not require parameters. The Backward Error attack, denoted BE, will use 5 iterations with a step size of $\alpha = 0.5$. C&W uses the same Adam learning rate and weight decay as in the untargeted case. To reach the desired classifications, it is now necessary to increase $c$ to 15 and $\eta$ to 0.002. The maximum number of iterations is increased to 50. Although not reached in most cases, we found that this slightly higher upper bound was occasionally useful in allowing the iteration to reach the target.

## 4.3   Results

The mean CER values for the different untargeted attacks are shown in fig. 2. We notice that C&W performs best for small $\epsilon$. This can be explained by the
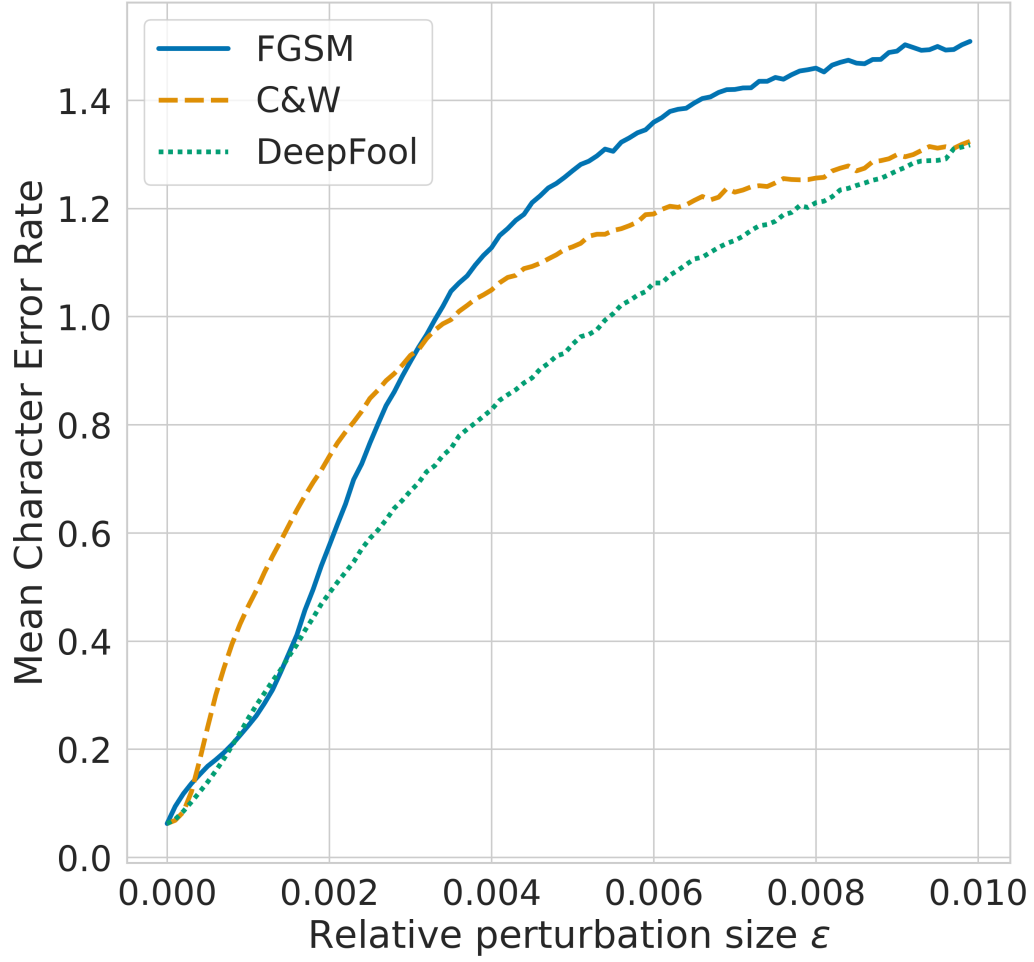
15

Figure 2: Comparison of mean Character Error Rate for adversarial attacks of different sizes generated by untargeted attacks FGSM, C&W, and DeepFool on IAM handwriting database.

fact that it focuses on changing the token labels for the smallest perturbation. However, this token change might not change all characters in the words since two different tokens can contain overlapping characters. Also, FGSM is focusing on the cross-entropy loss, which looks at all the logits. Increasing this loss will initially not always lead to misclassifications of tokens. Further, it could be trying to increase the loss of one token more than that of the others.

For larger perturbation sizes, FGSM becomes better than C&W. This can be explained by C&W being satisfied with any token change and optimizing the perturbation size for that change. Meanwhile, FGSM aims to increase the cross-entropy loss, which means that the logit for the actual class should decrease as much as possible. This will also lead to similar tokens not getting very high values. Also note that FGSM only overtakes when the CER is higher than 1. For such a high CER, insertions are made to make the sentences longer.

In this context, DeepFool performs the worst of the three algorithms for perturbations of all sizes. For larger perturbations, it suffers from the attack only considering a token misclassification, just like C&W. Since DeepFool tries to make the smallest step to the second-best class based on a linearization, it might be too small a step for multiple tokens when we consider small perturbations. It does not focus well on any of the two aspects that make the other two algorithms work well in their own respect.

To get a sense of the difference in scale of the perturbation sizes and the CER scores, we show an example of a C&W attack in fig. 3. We show the same attack rescaled for different values of $\epsilon$. The correct sentence is '" The Thetans, " he said, " are presumably here to'. For $\epsilon = 0.001$, we see that the words in the output sentence look fairly similar to the words in the original image. For instance, 'Thetans' becomes 'Rhetau' and 'are' becomes 'one'. For the larger $\epsilon$ values in the figure, the output text no longer matches the original image. It is interesting that the output for $\epsilon = 0.01$ consists of words that have some kind of mutual coherence. For $\epsilon = 0.1$, TrOCR states part of the alphabet and gets caught in a loop where it outputs 'k /' repeatedly. The CER generally increases as $\epsilon$ increases, but decreases for $\epsilon = 0.1$. This can be explained by the fact that shorter single-letter tokens need fewer insertions to create, which is the main driver of CER for large $\epsilon$ values. The image perturbations are essentially invisible for all but the largest $\epsilon$.

For further illustration, in fig. 4, we show the absolute values of the perturbation in fig. 3 of size $\epsilon = 0.01$ against a white background. We see
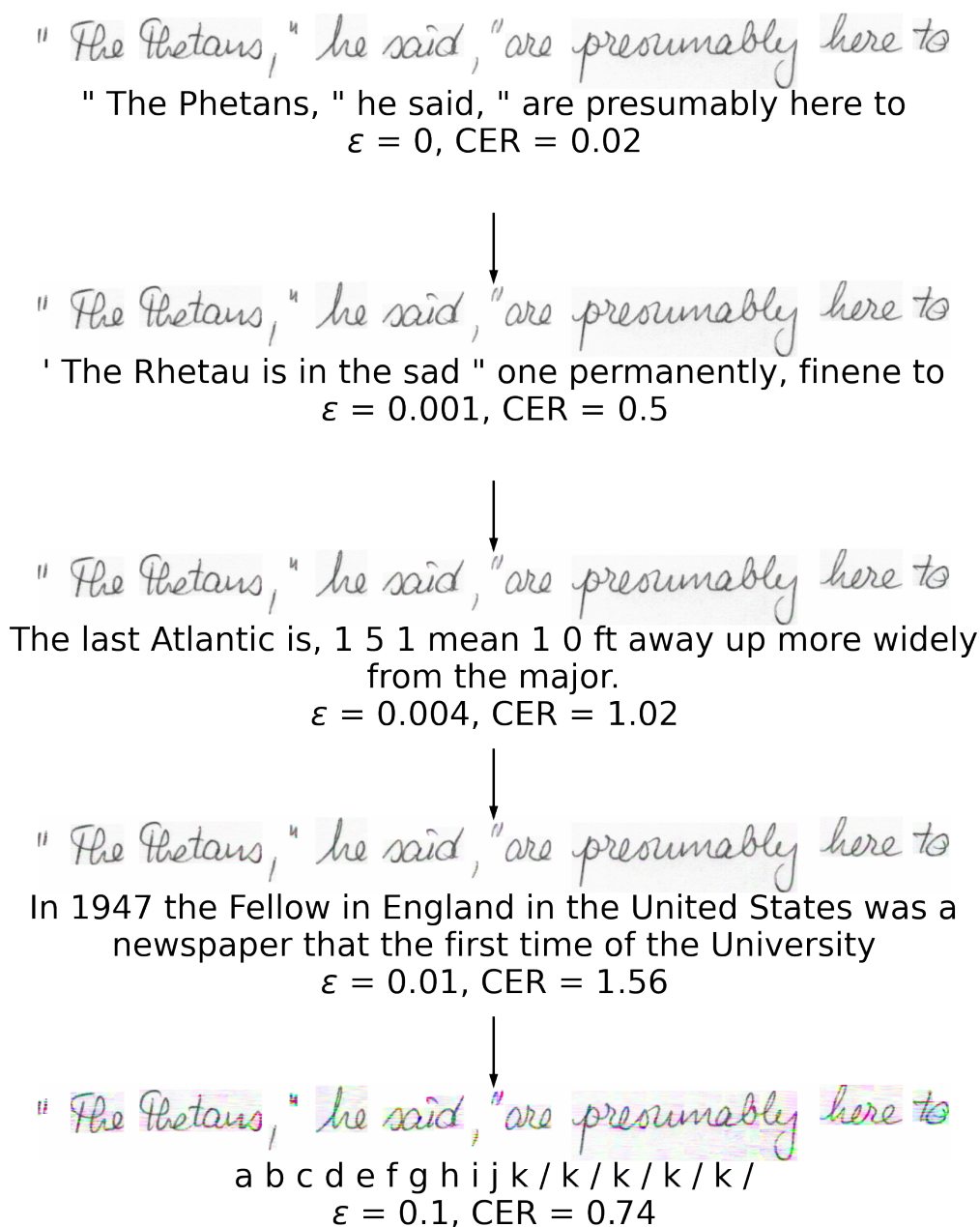
" The Phetans, " he said, " are presumably here to
$\varepsilon = 0$, CER = 0.02



' The Rhetau is in the sad " one permanently, finene to
$\varepsilon = 0.001$, CER = 0.5



The last Atlantic is, 1 5 1 mean 1 0 ft away up more widely
from the major.
$\varepsilon = 0.004$, CER = 1.02



In 1947 the Fellow in England in the United States was a
newspaper that the first time of the University
$\varepsilon = 0.01$, CER = 1.56



a b c d e f g h i j k / k / k / k / k /
$\varepsilon = 0.1$, CER = 0.74

Figure 3: Perturbed images created by C&W algorithm. From top to bottom, we start with the unperturbed image and gradually increase the relative norm of the perturbation. Each image evaluated with TrOCR. The output sentences are mentioned below the images. The correct text is '" The Thetans, " he said, " are presumably here to'.

Figure 4: Almost imperceptible element-wise absolute value of the perturbation from fig. 3 with $\epsilon = 0.01$.

only a vague indication of the perturbation; it is almost imperceptible. We also note that a CER of 1 was already reached with smaller perturbations.

From these tests, it can be concluded that TrOCR is highly vulnerable to untargeted attacks. The output can be completely changed to nonsensical text with a perturbation that is not visible to the human eye.

For the targeted case, fig. 5 shows that C&W needs significantly larger perturbation sizes to successfully reach the target labels, but it can successfully perturb 25.7% of the images. This is significantly higher than FGSM and BE, which are able to successfully perturb only 8.5% and 8.4% respectively. FGSM performs worse than BE, but quickly catches up as $\epsilon$ increases.

The fact that C&W has larger perturbation sizes can be attributed to the initial perturbation. A larger starting perturbation is needed to increase the number of successful attacks, but the algorithm is not very good at decreasing the size of the perturbation. Successful attacks are done using perturbations of sizes close to $\epsilon = 0.003$. In fig. 3 it can be seen that this size of perturbation is not noticeable. Therefore, it is preferable to have higher success rates for the cost of the perturbations being this size.

In these tests, TrOCR is seen to be somewhat vulnerable to targeted attacks, but less so than it is to untargeted attacks. Note that in the targeted case, the size of the vocabulary complicates the choice of objective—we chose to compute results for the tenth most likely class as a representative example, but it would be of interest to consider other options.

## 5   Conclusion

In this study, we devised the first range of attack strategies for the TrOCR model and conducted a comprehensive evaluation.

In untargeted settings, we observed the vulnerability of TrOCR to adversarial attacks, particularly highlighting the efficacy of FGSM and C&W attacks. C&W demonstrated superior performance for small perturbation
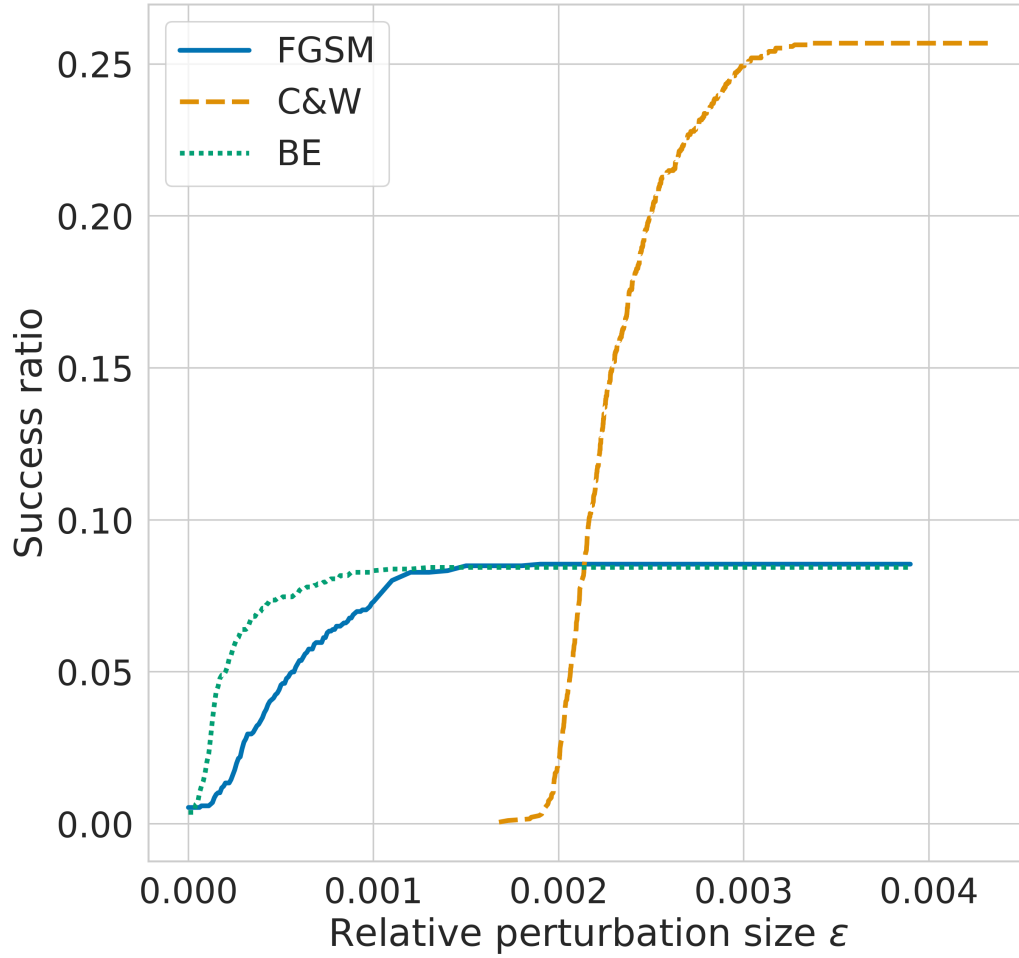
Figure 5: Success ratio for targeted attacks created by FGSM, C&W, and BE on IAM handwriting database test set. The target is to change a random token position to the tenth most probable token from the vocabulary.

sizes, emphasizing its ability to focus on token changes. However, as perturbation sizes increased, FGSM outperformed C&W, showcasing the nuances in attack strategies and their impact on TrOCR performance.

The results further revealed the susceptibility of TrOCR to targeted attacks, where we specifically targeted the tenth most likely token within the vocabulary. C&W, although requiring larger perturbation sizes, achieved a success rate of approximately 25%, outperforming the other algorithms.

TrOCR uniquely combines the advantages of CV and NLP models to create a powerful OCR process. However, our study demonstrates that TrOCR also inherits the vulnerabilities of the components that make up the overall computational pipeline. This raises immediate security concerns, especially in high-risk applications such as finance, law and education, that must be understood and addressed. Our results are also highly relevant for the current activity around regulation of AI. Clearly, for any new regulations to be meaningful and realistic, they must be informed by results about current algorithmic limitations.

Finally, we note that since our attack strategies build on established methodologies, there is potential to adapt existing defense strategies [51], including adversarial training [24]. We therefore hope that this work motivates further research into the development of OCR systems that are both powerful and resilient.

# Funding

# Data Statement

Code for the experiments presented here will be made available upon publication.

# References

[1] F. A. Acheampong, H. Nunoo-Mensah, and W. Chen, *Trans-*

*former models for text-based emotion detection: a review of bert-based approaches*, Artificial Intelligence Review, (2021), pp. 1–41.

[2] A. ALDAHDOOH, W. HAMIDOUCHE, AND O. DEFORGES, *Reveal of vision transformers robustness against adversarial attacks*, arXiv preprint arXiv:2106.03734, (2021).

[3] I. ASIMOV, *I, Robot*, Doubleday science fiction, Doubleday, 1950.

[4] Y. BAI, Y. WANG, Y. ZENG, Y. JIANG, AND S.-T. XIA, *Query efficient black-box adversarial attack on deep neural networks*, Pattern Recognition, 133 (2023), p. 109037.

[5] L. BEERENS AND D. J. HIGHAM, *Adversarial ink: Componentwise backward error attacks on deep learning*, IMA Journal of Applied Mathematics, (2023).

[6] T. BEUZEVILLE, P. BOUDIER, A. BUTTARI, S. GRATTON, T. MARY, AND S. PRALET, *Adversarial attacks via backward error analysis.* hal-03296180, version 3, Dec. 2021.

[7] N. BOUCHER, J. BLESSING, I. SHUMAILOV, R. ANDERSON, AND N. PAPERNOT, *When vision fails: Text attacks against ViT and OCR*, arXiv preprint arXiv:2306.07033, (2023).

[8] N. CARLINI, *A LLM assisted exploitation of AI-Guardian*, arXiv preprint arXiv:2307.15008, (2023).

[9] N. CARLINI AND D. WAGNER, *Towards evaluating the robustness of neural networks*, in 2017 ieee symposium on security and privacy (sp), Ieee, 2017, pp. 39–57.

[10] P.-Y. CHEN, H. ZHANG, Y. SHARMA, J. YI, AND C.-J. HSIEH, *Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models*, in Proceedings of the 10th ACM workshop on artificial intelligence and security, 2017, pp. 15–26.

[11] E. COMISSION, *Proposal for a regulation of the european parliament and of the council laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts*, (2021).

[12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, *Bert: Pre-training of deep bidirectional transformers for language understanding*, arXiv preprint arXiv:1810.04805, (2018).

[13] D. H. Diaz, S. Qin, R. Ingle, Y. Fujii, and A. Bissacco, *Rethinking text line recognition models*, arXiv preprint arXiv:2104.07787, (2021).

[14] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, *An image is worth 16x16 words: Transformers for image recognition at scale*, arXiv preprint arXiv:2010.11929, (2020).

[15] I. J. Goodfellow, J. Shlens, and C. Szegedy, *Explaining and harnessing adversarial examples*, in 3rd International Conference on Learning Representations, San Diego, CA, Y. Bengio and Y. LeCun, eds., 2015.

[16] J. Heo, S. Seo, and P. Kang, *Exploring the differences in adversarial robustness between vit-and cnn-based models using novel metrics*, Computer Vision and Image Understanding, 235 (2023), p. 103800.

[17] K. S. Kalyan, A. Rajasekharan, and S. Sangeetha, *Ammus: A survey of transformer-based pretrained models in natural language processing*, arXiv preprint arXiv:2108.05542, (2021).

[18] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, *Transformers in vision: A survey*, ACM Comput. Surv., 54 (2022), pp. 1–41.

[19] D. P. Kingma and J. Ba, *Adam: A method for stochastic optimization*, arXiv preprint arXiv:1412.6980, (2014).

[20] J. Lahoud, J. Cao, F. S. Khan, H. Cholakkal, R. M. Anwer, S. Khan, and M.-H. Yang, *3D vision with transformers: a survey*, arXiv preprint arXiv:2208.04309, (2022).

[21] S. Latif, A. Zaidi, H. Cuayahuitl, F. Shamshad, M. Shoukat, and J. Qadir, *Transformers in speech processing: A survey*, arXiv preprint arXiv:2303.11607, (2023).

[22] M. Li, T. Lv, J. Chen, L. Cui, Y. Lu, D. Florencio, C. Zhang, Z. Li, and F. Wei, *TrOCR: Transformer-based optical character recognition with pre-trained models*, in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 37, 2023, pp. 13094–13102.

[23] H. Liang, E. He, Y. Zhao, Z. Jia, and H. Li, *Adversarial attack and defense: A survey*, Electronics, 11 (2022), p. 1283.

[24] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, *Towards deep learning models resistant to adversarial attacks*, in 6th International Conference on Learning Representations, Vancouver, BC, OpenReview.net, 2018.

[25] K. Mahmood, R. Mahmood, and M. Van Dijk, *On the robustness of vision transformers to adversarial examples*, in Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 7838–7847.

[26] U.-V. Marti and H. Bunke, *The iam-database: an english sentence database for offline handwriting recognition*, International Journal on Document Analysis and Recognition, 5 (2002), pp. 39–46.

[27] J. Memon, M. Sami, R. A. Khan, and M. Uddin, *Handwritten optical character recognition (OCR): A comprehensive systematic literature review (SLR)*, IEEE Access, 8 (2020), pp. 142642–142668.

[28] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, *Deepfool: a simple and accurate method to fool deep neural networks*, in Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2016, pp. 2574–2582.

[29] J. X. Morris, E. Lifland, J. Y. Yoo, J. Grigsby, D. Jin, and Y. Qi, *Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp*, arXiv preprint arXiv:2005.05909, (2020).

[30] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, *Practical black-box attacks against machine learning*, in Proceedings of the 2017 ACM on Asia conference on computer and communications security, 2017, pp. 506–519.

[31] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, *Distillation as a defense to adversarial perturbations against deep neural networks*, in IEEE Symposium on Security and Privacy (SP), San Jose, CA, USA, 2016, pp. 582–597.

[32] E. Parliament, *Amendments adopted by the european parliament on 14 june 2023 on the proposal for a regulation of the european parliament and of the council on laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts*, (2023).

[33] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, *PyTorch: An imperative style, high-performance deep learning library*, in Advances in Neural Information Processing Systems 32, Curran Associates, Inc., 2019, pp. 8024–8035.

[34] R. Rajesh and R. Kanimozhi, *Digitized exam paper evaluation*, in 2019 IEEE International Conference on System, Computation, Automation and Networking (ICSCAN), IEEE, 2019, pp. 1–5.

[35] D. Rodriguez, T. Nayak, Y. Chen, R. Krishnan, and Y. Huang, *On the role of deep learning model complexity in adversarial robustness for medical images*, BMC Med Inform. Decis. Mak., 22 (2022).

[36] L. Schönherr, K. Kohls, S. Zeiler, T. Holz, and D. Kolossa, *Adversarial attacks against automatic speech recognition systems via psychoacoustic hiding*, arXiv preprint arXiv:1808.05665, (2018).

[37] R. Shao, Z. Shi, J. Yi, P.-Y. Chen, and C.-J. Hsieh, *On the adversarial robustness of vision transformers*, arXiv preprint arXiv:2103.15670, (2021).

[38] M. Sharif, L. Bauer, and M. K. Reiter, *On the suitability of lp-norms for creating and preventing adversarial examples*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2018, pp. 1605–1613.

[39] B. Shi, X. Bai, and C. Yao, *An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition*, IEEE transactions on pattern analysis and machine intelligence, 39 (2016), pp. 2298–2304.

[40] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, *OSQP: an operator splitting solver for quadratic programs*, Mathematical Programming Computation, 12 (2020), pp. 637–672.

[41] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, *Intriguing properties of neural networks*, arXiv preprint arXiv:1312.6199, (2013).

[42] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, *Training data-efficient image transformers & distillation through attention*, in International conference on machine learning, PMLR, 2021, pp. 10347–10357.

[43] I. Y. Tyukin, D. J. Higham, A. Bastounis, E. Woldegeorgis, and A. N. Gorban, *The feasibility and inevitability of stealth attacks*, arXiv:2106.13997, (2021).

[44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, *Attention is all you need*, in Advances in Neural Information Processing Systems, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds., vol. 30, Curran Associates, Inc., 2017, pp. 5998–6008.

[45] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, *Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers*, Advances in Neural Information Processing Systems, 33 (2020), pp. 5776–5788.

[46] X. Wang, H. Wang, and D. Yang, *Measure and improve robustness in NLP models: A survey*, in NAACL 2022 - 2022 Conference of the North American Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, 2022, pp. 4569–4586.

[47] Z. Wei, J. Chen, M. Goldblum, Z. Wu, T. Goldstein, and Y.-G. Jiang, *Towards transferable adversarial attacks on vision transformers*,

in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 36, 2022, pp. 2668–2676.

[48] Z. Wei, J. Chen, Z. Wu, and Y.-G. Jiang, *Cross-modal transferable adversarial attacks from images to videos*, in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), June 2022, pp. 15064–15073.

[49] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, *Transformers: State-of-the-Art Natural Language Processing*, in Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations, Association for Computational Linguistics, Oct. 2020, pp. 38–45.

[50] K. Y. Xiao, V. Tjeng, N. M. M. Shafiullah, and A. Mądry, *Training for faster adversarial robustness verification via inducing Relu stability*, International Conference on Learning Representations, New Orleans, USA, (2019).

[51] X. Yuan, P. He, Q. Zhu, and X. Li, *Adversarial examples: Attacks and defenses for deep learning*, IEEE transactions on neural networks and learning systems, 30 (2019), pp. 2805–2824.

[52] J. Zhang, Y. Huang, W. Wu, and M. R. Lyu, *Transferable adversarial attacks on vision transformers with token gradient regularization*, in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 16415–16424.

[53] J. Zhang, B. Li, J. Xu, S. Wu, S. Ding, L. Zhang, and C. Wu, *Towards efficient data free black-box adversarial attack*, in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), June 2022, pp. 15115–15125.

[54] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie, and J.-R. Wen, *A survey of large language models*, arXiv preprint arXiv:2303.18223, (2023).