# Counterfeit Answers: Adversarial Forgery against OCR-Free Document Visual Question Answering

Marco Pintore*, Maura Pintor*, Dimosthenis Karatzas‡ and Battista Biggio*

*University of Cagliari, Italy
†CINI, Italy
‡Computer Vision Center, UAB, Spain

*Abstract*—Document Visual Question Answering (DocVQA) enables end-to-end reasoning grounded on information present in a document input. While recent models have shown impressive capabilities, they remain vulnerable to adversarial attacks. In this work, we introduce a novel attack scenario that aims to forge document content in a visually imperceptible yet semantically targeted manner, allowing an adversary to induce specific or generally incorrect answers from a DocVQA model. We develop specialized attack algorithms that can produce adversarially forged documents tailored to different attackers' goals, ranging from targeted misinformation to systematic model failure scenarios. We demonstrate the effectiveness of our approach against two end-to-end state-of-the-art models: Pix2Struct, a vision-language transformer that jointly processes image and text through sequence-to-sequence modeling, and Donut, a transformer-based model that directly extracts text and answers questions from document images. Our findings highlight critical vulnerabilities in current DocVQA systems and call for the development of more robust defenses. We release our open source code at https://anonymous.4open.science/r/adv_docVQA-E7C5.

*Index Terms*—Adversarial Machine Learning, Document Understanding, LMM Security, Document and Text Processing

Figure 1: **Example DocVQA task on a synthetic invoice.** The model must read and reason over structured document text. In normal operation the model correctly answers questions on the unaltered document (left). By applying a simple adversarial patch perturbation (right), an adversary can force the model to answer a preselected (incorrect) response; for instance, the perturbed document shows $0.00, potentially causing monetary loss.

## 1. Introduction

Document Analysis research has demonstrated incredible advancements thanks to the adoption of Machine Learning (ML). Document Visual Question Answering (DocVQA) is a task that focuses on training ML models to answer questions posed on document images. Moreover, DocVQA has become a standard benchmark to evaluate current Large Multimodal Models (LMMs) [1]. While Optical Character Recognition (OCR) has been employed to extract textual information from document images for downstream tasks, the use of OCR-free techniques has been proposed as an alternative method to obtain highly efficient DocVQA models. The main characteristic of this strategy is that it allows end-to-end training of Deep Learning (DL) models, and overcomes limitations of the OCR-based systems such as the high computational cost, the error propagation to the downstream tasks, and the inflexibility to languages and structure of the document [2], [3]. Typically, OCR-free models receive as input an image and a question, and generate an answer in natural language. Instead of relying on explicitly extracted text, they learn to jointly model visual and layout information, leveraging both extracted language and visual cues such as text regions, formatting, and figures to infer the answer.

Despite being convenient and efficient, DL is known to suffer from adversarial perturbations, i.e. subtle manipulations of the input data that target the sensitivity of the DL models to trigger undesired behaviors [4], [5]. The extent to which end-to-end DocVQA systems are vulnerable to adversarial perturbations has not yet been studied. We believe that adversarial perturbations on DocVQA could open the door to a new age of *document forgery*, a criminal act involving the creation of false documents to deceive or defraud. With the advent of protocols such as the Agents Payments Protocol (AP2) [1], which empower AI agents to initiate financial transactions, such vulnerabilities in automatic document processing may translate directly into tangible monetary consequences. Figure 1 demonstrates with a simple example the potential risk caused by an attack on automatic document processing: a small adversarial perturbation on the bottom-right corner is able to steer

---

1. https://cloud.google.com/blog/products/ai-machine-learning/announcing-agents-to-payments-ap2-protocol

the model's output to the wrong answer, which in a fully automated pipeline could result in an AI agent authorizing a payment with the wrong amount.

In this paper, we present the first threat model for analyzing the adversarial robustness of DocVQA systems, with a particular focus on OCR-free architectures. Our contributions are threefold: (i) we define and formalize the threat model that leverages visually inconspicuous yet disruptive perturbations to manipulate model outputs toward different adversarial goals; (ii) we introduce attack algorithms specifically designed to target OCR-free DocVQA models, enabling the creation of adversarially forged documents; and (iii) we empirically demonstrate the effect of these attacks on two state-of-the-art DocVQA models, namely Pix2Struct [2] and Donut [3].

## 2. Adversarial Attacks against DocVQA

This section details our framework for analyzing the adversarial robustness of OCR-free DocVQA models. We first provide an overview of these models and define the threat model we adopt for this work. We then formalize the attack, and detail the different objectives of several attack scenarios. We conclude the section by discussing the specific challenges in attacking DocVQA models.

### 2.1. End-to-end DocVQA Systems

An end-to-end DL-based Document Visual Question Answering (DocVQA) model $f$ takes as input an image $\mathbf{x}$ and a natural-language question about its content $\mathbf{q}$ whose ground-truth (GT) answer is $\mathbf{y} = (y_1, \ldots, y_T)$. The model estimates a conditional distribution and factorizes the answer generation through auto-regression as

$$\hat{\mathbf{y}} \sim P_{\boldsymbol{\theta}}(\cdot \mid \phi(\mathbf{x}), \mathbf{q}) = \prod_{t=1}^{\hat{T}} P_{\boldsymbol{\theta}}\left(\hat{y}_t \mid \hat{y}_{<t}, f_{\text{enc}}(\phi(\mathbf{x}), \mathbf{q})\right),$$
(1)

where $\phi$ denotes image preprocessing (e.g., resizing, normalization), and $f_{\text{enc}}(\cdot)$ is the multimodal encoder. Note that $\hat{T}$ denotes the length of the generated string, that might be different than the GT answer length $T$. The encoder $f_{\text{enc}}$ extracts a multimodal latent representation from the preprocessed image:

$$\mathbf{H} = f_{\text{enc}}(\phi(\mathbf{x}), \mathbf{q}) \in \mathbb{R}^{N \times d},$$
(2)

where $\mathbf{H}$ is a sequence of $N$ feature vectors. The encoder may encode the question jointly with the image (as Pix2Struct [2]), or accept the question separately as prompt tokens and concatenate it with an encoding of the image (as Donut [3]). Both behaviors are captured by $f_{\text{enc}}$. The decoder produces, at each step $t$, unnormalized scores (logits) over the vocabulary of size $\mathcal{V}$, conditioned on the encoder output and previously produced tokens:

$$\mathbf{z}_t = f_{\text{dec}}(\hat{y}_{<t}, \mathbf{H}) \in \mathbb{R}^{\mathcal{V}}.$$
(3)

The output is then obtained by applying a softmax:

$$P_{\boldsymbol{\theta}}(\hat{y}_t \mid \hat{y}_{<t}, \phi(\mathbf{x}), \mathbf{q}) = \text{softmax}(\mathbf{z}_t).$$
(4)

The final token at step $t$ is taken as the token with the highest probability, i.e.

$$\hat{y}_t = \arg\max_{v \in \mathcal{V}} \text{softmax}(\mathbf{z}_t)[v]$$
(5)

where $\text{softmax}(\mathbf{z}_t)[v]$ denotes the probability assigned to the token $v$.

**Model training.** The goal of these models is to produce an answer such that $\hat{\mathbf{y}} = \mathbf{y}$. This is done by fine-tuning a pretrained model on a dataset

$$\mathcal{D}^{\text{tr}} = \left\{ \left( \mathbf{x}_i, \{(\mathbf{q}_{i,j}, \mathbf{y}_{i,j})\}_{j=1}^{M_i} \right) \right\}_{i=1}^{N},$$

where each document $\mathbf{x}_i$ is associated with $M_i$ question-answer pairs (QA pairs). The model parameters $\boldsymbol{\theta}$ are found by minimizing the token-level negative log-likelihood:

$$L_{\boldsymbol{\theta}}(\phi(\mathbf{x}), \mathbf{q}, \mathbf{y}) = - \sum_{(\mathbf{x}, \mathbf{q}, \mathbf{y}) \in \mathcal{D}^{\text{tr}}} \sum_{t=1}^{T} \log(\text{softmax}(\mathbf{z}_t)[y_t]),$$
(6)

where $(\text{softmax}(\mathbf{z}_t)[y_t])$ is the probability assigned to the correct token $y_t$.

**Model evaluation.** In the DocVQA task, exact matches between the correct answer and the intended response cannot be used directly, as minor variations in formatting, punctuation, or tokenization would cause valid variations of the correct answers to be scored as wrong. For this reason, metrics such as the Average Normalized Levenshtein Similarity (ANLS) are used to account for minor variations in text answers while capturing semantic fidelity [6]. The ANLS metric provides a robust assessment that smoothly captures minor by penalizing minor deviations, ensuring that semantically correct answers are still appropriately recognized. It is defined as

$$\text{ANLS} = \frac{1}{\sum_{i=1}^{N} M_i} \sum_{i=1}^{N} \sum_{j=1}^{M_i} s_{i,j},$$
(7)

where $\sum_{i=1}^{N} M_i$ is the number of question-answer pairs, and $s_{i,j}$ is the similarity score for the $j$-th QA pair of the $i$-th document, computed as:

$$s_{i,j} = \begin{cases} \text{NLS}(\hat{\mathbf{y}}_{i,j}, \mathbf{y}_{i,j}) & \text{if } \text{NLS}(\hat{\mathbf{y}}_{i,j}, \mathbf{y}_{i,j}) \geq \tau \\ 0 & \text{otherwise} \end{cases}$$

Here, $\text{NLS}(\hat{\mathbf{y}}, \mathbf{y})$ is the Normalized Levenshtein Similarity (which ranges from 0 to 1, where 1 is the most similar) between the predicted answer $\hat{\mathbf{y}}$ and the GT answer $\mathbf{y}$, and $\tau$ is a threshold that distinguishes between answers that are correctly identified but improperly recognized versus those that are fundamentally incorrect.

### 2.2. Adversarial Documents

We now detail how we leverage gradient-based attacks to perturb the inputs of DocVQA systems to elicit *manipulated* answers. Following the adversarial machine learning literature [7], we first outline the assumptions underlying the threat model of this attack.

**Attacker's Goal.** The overall objective of the attacker is to perturb the input document in order to manipulate the model's outputs. The attacker's objective can be: (i) steer the model toward a specific target answer for a single question (targeted, single-QA manipulation), (ii) control the answers to multiple questions on the same document (targeted, multi-QA manipulation), or (iii) broadly prevent the model from producing correct answers for a document

(untargeted, single- or multi-QA manipulation). We see the first two as an *evasion* attack (integrity violation at test time) and the latter as a *Denial of Answer*[2] attack (availability violation at test time). After presenting the overall threat model and general formulation, we will formalize these cases, which we call *scenarios*, in terms of their objectives.

**Attacker's knowledge.** We consider a white-box attack scenario in which the adversary has complete knowledge of the target DocVQA model, including not only its architecture, parameters, and gradients, but also the full system (including the preprocessing steps). In this setting, the attacker possesses both the original target document $\mathbf{x}$ and knows the specific correct answer (target answer) $\mathbf{y}^{\star}$ they wish to avoid (elicit) from the model. Furthermore, the attacker knows the question $\mathbf{q}$ that will be asked on the document, e.g., a system prompt used by an automatic information extraction system (see Figure 1). This represents a strong level of access, where the adversary can leverage full model transparency to craft highly effective adversarial perturbations. The white-box assumption enables the use of gradient-based optimization techniques to systematically modify the document visual content (i.e. the input image) in ways that are visually inconspicuous to humans yet semantically manipulate the model's outputs. While this scenario may seem overly restrictive, it provides important insights into the fundamental vulnerabilities of DocVQA systems and establishes an upper bound on attack effectiveness that can inform the development of robust defenses. Notably, some of these assumptions can be relaxed in practice, e.g. by averaging perturbations over multiple QA pairs to reduce the dependence on knowing the exact ones. We leave this as a future development.

**Attacker's capability.** In our formulation, the attacker can modify, for a single document $\mathbf{x}$, its pixel values. This means that the attacker has only access to the test data, i.e. the attack happens at test time. The attacker can modify the pixels and their colors independently, as long as the perturbation remains small or contained. This is often enforced with an $\ell_p$ constraint centered on the sample $\mathbf{x}$, or equivalently, containing the $\ell_p$-norm of the perturbation. Otherwise, especially for the case of adversarial documents, the perturbation can also be restricted to a region of contiguous pixels, like in the patch attacks [8]. Notably, while the traditional patch attacks usually need to generalize over different sizes, rotations and dimensions of the patch, in this case, these can be decided by the attacker and can be very accurately enforced [9]. We assume that the question $\mathbf{q}$ is fixed, as the case of a system prompt for automatic document processing. Therefore, the attacker can only modify the document, whereas has no control of the question $\mathbf{q}$.

**Attack strategy.** We will first outline the attack formulation for a single QA pair (single-QA attack), then we will derive the variations covering the other more complex scenarios.

As we select a single sample $\mathbf{x}$ to craft this attack, we will omit the index $i$ unless otherwise specified. We define an adversarial example as a perturbed input $\mathbf{x}' = \mathbf{x} + \boldsymbol{\delta}$, where $\boldsymbol{\delta}$ is constrained in magnitude to be imperceptible (e.g., under some $\ell_p$ norm). In this work, we consider

the $\ell_{\infty}$ norm, and we apply the perturbation to the entire input (which we refer to as the *full-document* setting) and to a patch-based perturbation model (*patch* setting), i.e., perturbations applied only to a selected contiguous region of the document image, such as stamps, watermarks, or logos [8]. The general objective of the attack is to find a perturbation that causes the model to output an incorrect (or targeted) answer, despite the image being visually similar to the untainted one. The attacker has control only over the input image $\mathbf{x}$, whereas the question $\mathbf{q}$ cannot be changed (as it is, for instance, a system prompt provided by an automatic processing system). Importantly, our approach performs end-to-end manipulation in the input space, meaning we directly perturb the raw document image $\mathbf{x}$ rather than the intermediate representations after preprocessing or the embeddings. This ensures that the resulting adversarial example $\mathbf{x} + \boldsymbol{\delta}$ is a valid image that can be stored, transmitted, and processed by any DocVQA system, making the attack practically deployable.

We formalize the objective of the single-QA attack as:

$$\begin{aligned}
\arg\min_{\boldsymbol{\delta}} \quad & \gamma L_{\boldsymbol{\theta}}(\phi(\mathbf{x} + \boldsymbol{\delta}), \mathbf{q}, \mathbf{y}^{\star}) \\
\text{s.t.} \quad & \|\boldsymbol{\delta}\|_{\infty} \leq \epsilon \\
& \mathbf{x}_{\text{lb}} \leq \mathbf{x} + \boldsymbol{\delta} \leq \mathbf{x}_{\text{ub}},
\end{aligned} \tag{8}$$

where $\gamma$ is set to $-1$ (negative) for untargeted attacks (maximizing loss w.r.t. the GT answer $\mathbf{y}^{\star} = \mathbf{y}$) and $+1$ (positive) for targeted attacks (minimizing loss w.r.t. the target answer $\mathbf{y}^{\star} = \mathbf{y}^{t}$). The loss $L$ enforces the output to get close to the target. Within our setup, we employ two distinct loss functions: the first is the standard loss used for fine-tuning the model to the DocVQA tasks (eq. 6), while the second is a custom loss introduced in eq. (10), which we designed to improve the attack's efficacy against the Donut model. The norm constraint $\|\boldsymbol{\delta}\|_{\infty}$ limits the magnitude of the perturbation in the original input space, before applying any preprocessing $\phi$. The last constraint defines $\mathbf{x}_{\text{lb}}$ and $\mathbf{x}_{\text{ub}}$ as pixel-wise lower and upper bounds (typically $[0, 255]$ for images). Again, this is enforced before preprocessing. This same constraint also covers patch-based attacks. In these cases, the mask is naturally encoded by setting $\mathbf{x}_{\text{lb}} = \mathbf{x}_{\text{ub}} = \mathbf{x}$ for pixels outside the target patch region, effectively constraining $\boldsymbol{\delta}$ to be zero in those areas while allowing perturbations within the patch. Additionally, for some models, such as Pix2Struct [2], the question is directly rendered onto the image by the preprocessing function $\phi$. In this case, the perturbation should be applied only on the part of the document that does not contain the question (as it is practically unavailable to the attacker). $\phi$ A suitable strategy would be to impose a patch-like constraint or apply a spatial mask so that gradients only update pixels outside the question overlay. In our attack, we instead reimplement the preprocessing pipeline end-to-end and apply the perturbation before the question is rendered. This ensures that the perturbation is introduced solely on the document content while still allowing gradients to flow through the differentiable rendering and scaling operations performed by $\phi$.

Once the optimization problem in Eq. (8) is defined, a standard and practical solver is Projected Gradient Descent (PGD) [10]. Starting from an initial perturbation $\boldsymbol{\delta}^{(0)}$ (e.g.,

---

2. as it shares similarities with the Denial of Service, as the adversary effectively denies the system's intended service.

zero or a small random initialization), PGD performs the iterative updates

$$\boldsymbol{\delta}^{(k+1)} = \Pi_{\mathcal{B}}\Big(\boldsymbol{\delta}^{(k)} - \alpha\,\nabla_{\boldsymbol{\delta}}\Big[\gamma L_{\boldsymbol{\theta}}\Big(\phi(\mathbf{x}+\boldsymbol{\delta}^{(k)}),\mathbf{q},\mathbf{y}\Big)\Big]\Big),$$

where $\alpha$ is a step size, $\nabla_{\boldsymbol{\delta}}$ denotes the gradient of the loss w.r.t. the input perturbation (computed by backpropagating through $\phi$ and the model), and $\Pi_{\mathcal{B}}$ is the projection operator onto the feasible set

$$\mathcal{B} = \{\boldsymbol{\delta} : \|\boldsymbol{\delta}\| \leq \epsilon,\ \mathbf{x}_{\mathrm{lb}} \leq \mathbf{x}+\boldsymbol{\delta} \leq \mathbf{x}_{\mathrm{ub}}\}.$$

For patch attacks, the projection additionally enforces zeros outside the patch (equivalently achieved by applying a binary mask to the update).

Since document images are typically quantized to discrete pixel levels, the resulting adversarial image is quantized after each iteration to preserve realistic values:

$$\boldsymbol{\delta} = \mathrm{Quantize}(\mathbf{x}+\boldsymbol{\delta}) - \mathbf{x}.$$

This ensures that each iteration yields a perturbation consistent with the discrete image domain. In practice, variants of PGD such as momentum-accelerated PGD [11] or adaptive optimizers (e.g., Adam) may be used in place of the plain gradient step to improve convergence.

**Loss design.** During initial experiments, we observed that the standard fine-tuning loss of Eq. (6), stops improving without achieving the desired target answers for all models considered. Therefore, we design a new loss function intended to amplify the contribution of target tokens and reduce the influence of tokens that compete most strongly with them. This approach significantly increases the success rate of attacks in those models where standard loss does not suffice for the attack to be successful.

Let $\mathbf{y}^{\star} = (y_1, \ldots, y_T)$ be the target sequence and $\mathbf{z}_t \in \mathbb{R}^{|\mathcal{V}|}$ the logit vector produced by the decoder at position $t$, calculated by conditioning the perturbed input $\mathbf{x}'$ and the prompt (see Sect. 2.1). Following the notation from the previous sections, we define the target logit token as $\mathbf{z}_t[y_t]$, and denote the highest logit among all tokens generated by the model at that position as $z_t^{\mathrm{top}}$:

$$z_t^{\mathrm{top}} = \max_k \mathbf{z}_t[k].$$

We then define the token-level logit loss as

$$L_{\boldsymbol{\theta}}^t = \begin{cases} z_t^{\mathrm{top}} - \mathbf{z}_t[y_t], & \text{if } \arg\max_k \mathbf{z}_t[k] \neq y_t, \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

The total loss on the entire target sequence is obtained by summing the contributions of each token:

$$L_{\boldsymbol{\theta}} = \sum_{t=1}^{T} L_{\boldsymbol{\theta}}^t. \quad (10)$$

In this way, the loss penalizes only target tokens that are not currently already the maximum logit. This guides the optimization towards bringing out the target tokens without changing those that are already correct. The function is differentiable with respect to the logits $\mathbf{z}_t$, allowing backpropagation to the input $\mathbf{x}'$ to update the perturbation. **Attack Scenarios.** Here, we specialize the above formulation to cover different goals of the attacker. Starting from Eq. (8), we formulate several attack variants, each differing in its constraints and optimization targets.

*Targeted, single-answer manipulation.* The attacker selects a single QA pair $(\mathbf{x}, \mathbf{q}, \mathbf{y})$ and a target answer $\mathbf{y}^{\star}$, optimizing the objective in Eq.(8) with a positive sign (i.e., $\gamma = 1$, minimizing the loss $L_{\boldsymbol{\theta}}$ on the target answer). *Targeted, multi-answer manipulation.* The attacker selects a set of QA pairs on the sample $\mathbf{x}$, where for each question there is a specific target answer. Thus, given $\{(\mathbf{x}, \mathbf{q}_j, \mathbf{y}_j^{\star})\}_{j=1}^{M}$, the objective jointly minimizes the sum of losses across all QA pairs:

$$\arg\min_{\boldsymbol{\delta}} \quad \sum_{j=1}^{M} L_{\boldsymbol{\theta}}(\phi(\mathbf{x}+\boldsymbol{\delta}), \mathbf{q}_j, \mathbf{y}_j^{\star}) \quad (11)$$

*Denial of Answer.* The attacker seeks to degrade performance globally by inducing wrong answers to one or all questions on a sample $\mathbf{x}$. This can be obtained by modifying slightly the objective in eq. (11), using the negative sign of the loss $L$, (i.e. $\gamma = -1$) so that it maximizes the loss to the correct answers, and sets $\mathbf{y}_j^{\star} = \mathbf{y}_j$, $j = 1, \ldots M$ for inducing a total denial of service in which the model outputs wrong answers to all questions about the given sample. Ideally, by averaging over several questions, the attack could also generalize to unseen questions. However, in our results, we found that this does not seem to happen in practice (see Figure 5 in our experimental results).

**Specific challenges.** Implementing custom adversarial attacks for DocVQA models presents several challenges. First, unlike standard vision models trained on datasets such as ImageNet, where preprocessing is relatively uniform and limited to simple operations (e.g., resizing, cropping, normalization), DocVQA systems often employ model-specific preprocessing pipelines. These may include a combination of lossy and highly specialized transformations such as compression, adaptive rescaling, document layout adjustments, or text-enhancement operations. Consequently, the preprocessing function $\phi$ available to an adversary is only an approximation of the true, model-specific pipeline. While one option in our work is to reverse-engineer this pipeline (as will be discussed in Sect. 3 for each target model), alternative strategies such as gradient-based approximation of the transformation function are also possible [12]. Second, optimizing over the full generation-aware objective, which requires backpropagation through the autoregressive decoding process, is computationally demanding.

## 3. Experiments

In this section, we present the experiments carried out to evaluate our attack methods. We start by describing the experimental setup in Sect. 3.1, including the dataset used, the models, and the attack settings. In Sect. 3.2, we describe the metrics that we use to evaluate the effectiveness of the attacks. In Sect. 3.3 we show the results obtained, while Sect. 3.4 illustrates some qualitative examples of forged documents and the model's responses.

### 3.1. Experimental Setup

**Dataset.** We conduct our experiments in the PFL-DocVQA dataset [13]. This dataset contains real documents related to invoices, in which each invoice is associated with a

question and multiple answers. It is originally designed to test existing privacy techniques on multi-modal DocVQA scenarios. In total, it contains $336,842$ question-answer pairs on $117,661$ pages, resulting in $37,669$ documents from $6,574$ different providers. Although the authors provide a Blue Team/Red Team split to separate the data between training and privacy attack evaluation, we merged them to obtain a single unified set. To build our evaluation set, we extracted $N = 1000$ unique samples from the merged data, where each sample is composed of a single image $\mathbf{x}$ and exactly $M = 5$ associated QA pairs.

**Models.** We consider two state-of-the-art DocVQA models: Pix2Struct-Base [2] and Donut [3], which propose end-to-end architectures designed for OCR-free document understanding. We use the publicly available checkpoint from HuggingFace [3].

*Pix2Struct-Base.* This model is a 282 M parameters image encoder-text decoder pretrained to perform layout parsing from image-text pairs. We use the version fine-tuned for the DocVQA task. This version uses a special preprocessing to render the question as a header at the top of the original image, providing both the question and image jointly with the visual modality, relying exclusively on visual inputs. Moreover, the preprocessing of Pix2Struct uses a special rescaling that extracts the maximum number of fixed-size patches from the document without aspect-ratio distortion. Finally, the document image is normalized with its own mean and standard deviation. This model achieves strong performance on the DocVQA benchmark [1], reaching an ANLS score of 72.1%, and on the Infographic VQA benchmark [14], with an ANLS score of 38.2%.

*Donut.* This model is a 176 M parameters transformer architecture pre-trained with a pseudo-OCR task (i.e., it learns to extract text from the input image). For the version fine-tuned for the DocVQA task, the question is provided as a starting prompt to the decoder using a fixed template with special tokens, i.e. `<s_docvqa><s_question>`q`</s_question>` `<s_answer>`. The embedding extracted from the image and the question is tokenized and prepended to the decoder input before passing the image, enabling the model to autoregressively generate the answer text. The image preprocessing applied first resizes the document, adds padding to match the model's expected resolution, and then the pixel values are normalized to the ImageNet default mean $(0.485, 0.456, 0.406)$ and standard deviation $(0.229, 0.224, 0.225)$. This model achieves an ANLS score of 67.5% on the DocVQA benchmark.

**Baseline Performance on the PFL-DocVQA subset.** We evaluate Pix2Struct-Base and Donut on our custom PFL-DocVQA subset. Using the publicly available DocVQA checkpoints, Pix2Struct-Base reaches an ANLS score of 51.27%, while Donut reaches 41.14%. These results will be reported in the experimental plots in Sect. 3.3, providing baseline references for our experiments, which we refer to as the ANLS-baseline.

**Model-specific preprocessing.** We summarize the preprocessing steps for the two target models in Table 1. Adversarial perturbations must ultimately be applied in the input image space, i.e., before any model-specific preprocessing, to remain valid inputs for the system. This

3. https://huggingface.co/models

TABLE 1: Summary of the preprocessing steps applied by Pix2Struct and Donut.

| Step | Pix2Struct [2] | Donut [3] |
|---|---|---|
| Question handling | Rendered onto image | Passed to prompt decoder |
| Resizing | Patch extraction | Resize and padding |
| Normalization | Sample-wise | Fixed (ImageNet) |

means that, after the attack, the document can be saved as a file and then re-loaded again and processed directly without errors or loss of attack functionality.

To create attacks that are end-to-end differentiable, for each model, we reverse-engineered the preprocessing function $\phi$. We will now detail how we reverse-engineered the preprocessing steps in order to let the gradient needed to craft the perturbation flow through a fully-differentiable computational graph back to the input.

*Question handling.* First, we need to make sure the perturbation is only applied to the document, and not to the question or to the parts that are not in control of the attacker. Pix2Struct renders the question text directly onto the top of the input image as a header, meaning the question becomes part of the visual input processed by the encoder. However, the attacker cannot modify the header part, as this is applied directly by the model during inference. To ensure this, we apply a mask to the header part when backpropagating through the model, so that these gradients are zeroed during optimization and won't contribute to the overall perturbation. In contrast, Donut treats the question as a separate text prompt passed to the decoder, therefore the entire image can be perturbed without requiring any specific measure.

*Resizing.* Pix2Struct employs aspect-ratio preserving scaling followed by patch-based extraction, which avoids distorting the original document layout. Donut, however, resizes images to a fixed canvas size with padding to maintain aspect ratio. Scaling preprocessing relies on standard interpolation routines (e.g. bilinear resampling in PyTorch), which are lossy operations that slightly alter pixel values. To ensure gradient flow through this step, we reimplemented the preprocessing pipeline using fully differentiable operations, allowing backpropagation through the rescaling and normalization stages.

*Normalization.* The normalization strategies simply apply mean centering and variance scaling. Pix2Struct applies per-image standardization using the sample-wise mean and standard deviation. Donut uses fixed ImageNet statistics for normalization. We handle backpropagation through these normalizations again by implementing them in fully differentiable operations.

Thus, backpropagating through the preprocessing required reimplementing the entire pipeline within the PyTorch framework. With this adaptation, the pipeline maintains end-to-end differentiability. This allows gradients to flow from the model's loss back through all preprocessing operations to the raw input image, enabling gradient-based optimization of adversarial perturbations. We illustrate the difference between the standard non-differentiable preprocessing and our differentiable implementation in Figure 2.

We care to specify that we apply these modifications to the pipeline only for crafting the attacks. Once saved, the
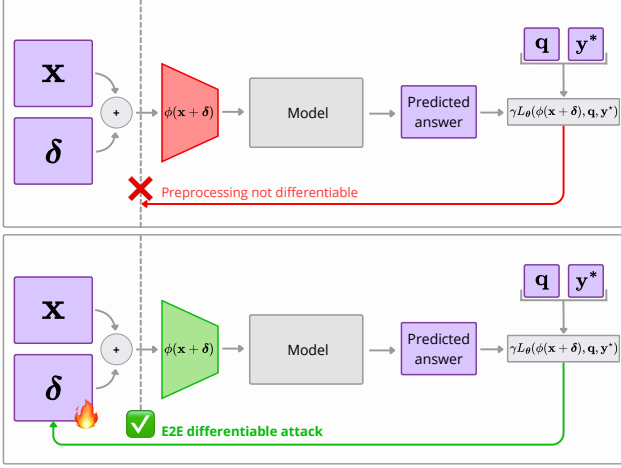
Figure 2: **End-to-end attack.** The non-differentiable pre-processing (*top*), which breaks the computational graph before $\delta$, and our end-to-end differentiable attack (*bottom*), obtained by reverse engineering the preprocessing $\phi$.

manipulated images are tested against a newly-instantiated, unmodified version of the HuggingFace models, thereby reflecting realistic attack conditions in which the adversary interacts with the model as publicly deployed.

**Attack Setup.** Following the formalization introduced in eq. (8), we create our attacks using PGD with an $\ell_\infty$ constraint. A concise overview of all losses and hyperparameters is given in Table 2, with the detailed configuration described below.

*Pix2Struct attack setup.* For the targeted attack scenario, we use the loss defined in eq. (6). The full-document attack is carried out using $K = 20$ attack steps, with a step size $\alpha = 2$ and the perturbation budget of $\epsilon = 8$. For the patch attack, we restrict the perturbation to a square with size equal to the 15% of the minimum dimension of the document, placed in the bottom-right corner, and we employ $K = 25$, $\alpha = 24$ and $\epsilon = 96$. For the untargeted attack (DoA), we optimize eq. (11) with the same hyperparameters, but setting $K = 1$, showing it is already enough to disrupt the models' correct functionality on all QA pairs.

*Donut attack setup.* We use the loss defined in eq. (10). In the full-document setting we use $K = 100$, $\alpha = 2$, and $\epsilon = 32$. For the patch attack, we use the same patch constraint as above, with hyperparameters set to $K = 100$, $\alpha = 24$ and $\epsilon = 96$. For the untargeted scenario (DoA), we optimize again eq. (11) with $K = 1$.

In all cases, we optimize adversarial examples such that, for each optimization step involving $B \leq M$ questions, the model is increasingly pushed to produce the first $B$ responses from the fixed set: $y^t \in \{$"No Answer", "Unclear", "Retry", "Try later", "I won't tell you"$\}$. We highlight here that the attacker can select any string as a target, including custom amounts (to cause directly financial fraud) or even hidden commands (to execute indirect prompt injections in downstream models [15]). We provide a visualization in Figure 3 to show how, with the given setup, the attack remains difficult to detect.

TABLE 2: Summary of the losses and hyperparameters used for the attacks.

| Model (Setting) | Loss | $\epsilon$ | $\alpha$ | $K$ |
|---|---|---|---|---|
| Pix2Struct (full) | eq. (6) | 8 | 2 | 20 |
| Pix2Struct (patch) | eq. (6) | 96 | 24 | 25 |
| Donut (full) | eq. (10) | 32 | 2 | 100 |
| Donut (patch) | eq. (10) | 96 | 24 | 100 |



Figure 3: **Attack Visualization.** The document altered with a patch in the lower right corner, with $\epsilon = 96$.

## 3.2. Evaluation Metrics

Across all scenarios, we evaluate attacks using three main metrics: (i) Attack Success Rate (ASR), which measures whether the targeted answers are successfully forced (or corrupted, in untargeted attacks); (ii) Collateral Damage (CDMG), which quantifies unintended changes to non-targeted answers within the same document; and (iii) ANLS-based scores, which capture token-level textual similarity and allow us to assess how predictions deviate from ground truth or move toward targeted answers. We measure these metrics as a function of $B$, i.e. the number of QA pairs involved in the optimization objective.

**Attack Success Rate.** For all the attack scenarios, we measure the Attack Success Rate (ASR), defined as the percentage of cases where there is an exact match between the model answer and the one with which the adversarial example was optimized, i.e. $\hat{y} = y^\star$ for targeted attacks and $\hat{y} \neq y$ for the untargeted. As before, we use the index $i$ to indicate the sample, and $j$ to index the QA pairs related to that sample. The ASR for targeted attacks is

then defined as:

$$\text{ASR} = \frac{1}{B} \sum_{i=1}^{N} \left( \prod_{j=1}^{B} \mathbf{1} \left[ \hat{\mathbf{y}}_{i,j} = \mathbf{y}_{i,j}^{\star} \right] \right) \qquad (12)$$

where $B$ is the number of QA pairs targeted by the attack. Note that, for multi-answer objectives, the success counts only if *all* the QA pairs match the objective. For the untargeted attacks, the ASR is computed using $\hat{\mathbf{y}}_{i,j} \neq \mathbf{y}_{i,j}$ in the indicator function. Similarly to the other case, if multiple QA pairs are involved, the success is counted only if *all* the answers are incorrect.

**Collateral Damage.** To evaluate how the attack affects the rest of the answers for a given document, i.e. the ones that are not in the optimization objective, we also compute the Collateral Damage (CDMG). We define the CDMG as the percentage of QA pairs not optimized whose prediction $\hat{\mathbf{y}}$ is wrong with respect to the answer $\mathbf{y}$:

$$\text{CDMG} = \frac{1}{CN} \sum_{i=1}^{N} \left( \sum_{j=1}^{C} \mathbf{1} \left[ \hat{\mathbf{y}}_{i,j} \neq \mathbf{y}_{i,j}^{\star} \right] \right) \qquad (13)$$

where $C = M - B$ denotes the set of QA pairs not involved in the optimization of the attack. As our dataset involves documents with $M = 5$ QA pairs, when $B = 5$, i.e. when we optimize all QA pairs for the document, this metric is undefined (and for this reason the curves are shown only for $C \in 1, \ldots, 4$). This metric is measured in the same way for targeted and untargeted attacks.

**ANLS.** For each attack scenario, we also measure the ANLS as defined in eq. (7), with $\tau = 0.5$. Different from the ASR metric – which evaluates adversarial success jointly over all targeted question of a document – the ANLS score is computed independently for each QA pair. Moreover, the ANLS does not account for an exact match, but uses the NLS as described in Sect. 2.1. Each QA pair contributes one similarity score, and the final value is obtained by averaging across the full evaluation set. In our setting, we are interested in three distinct evaluations of this metric. The first one is the average ANLS on the ground truth, which we denote as ANLS-baseline. This is the starting point, i.e. the score without any perturbation. The ANLS-baseline measures the similarity between the model predictions and the original GT answers of the evaluation set (on all the $T$ QA pairs). The second one is ANLS-B, which measures the similarity between the model predictions and the $B$ answers involved in the adversarial optimization. Note that the ANLS-B is expected to increase in the case of the targeted attacks, as the goal is to bring the answers closer to the targets. Conversely, in the case of untargeted attacks, we expect the metric to decrease (lower similarity) with respect to the correct answers. Finally, we measure the ANLS-C, which measures the ANLS on the set of QA pairs kept out of the optimization objective. Thus, this metric should be correlated with the CDMG but captures if the effect of the damage is only disrupting minimally the other answers (thus just enough to fail the exact match as in eq. 13), or, if the perturbation completely changes the original answers.

**Discussion.** For targeted, single-answer attacks, we expect high ASR and low CDMG, as the perturbation is optimized towards manipulating a specific QA pair. On the other hand, for targeted multi-answer manipulation scenarios,

high ASR might be more difficult to achieve, especially in the targeted attack (due to the fact that multiple QA pairs are optimized, and the ASR accounts for exact matches of *all* answers). CDMG may moderately increase because the visual features that influence one of the multiple answers optimized may overlap with those relevant to other questions, and therefore the optimization may inadvertently modify representations shared between multiple answers. Finally, in the denial-of-answer setting, the ASR is less restrictive as it is enough to disrupt the exact match of the answers with the ground truth (i.e. even by a single character). Moreover, in this case, the perturbation is explicitly designed to degrade the overall performance on the document rather than targeting specific answers. Thus, the high CDMG is a *desired* result.

### 3.3. Experimental results

To fully evaluate the effectiveness of our method, for each model and each attack scenario, we show how the attack performance varies as the complexity of the attack objective increases, i.e. as a function of the number of QA pairs involved in the optimization.

**Targeted Single-Answer manipulation** ($B = 1$)**.** Our first goal is to establish the effectiveness of the attack in the base scenario, i.e. the manipulation of a single answer. This case corresponds to the starting point of the curve in Figure 4, i.e. $B = 1$. When perturbing the full document, our method achieves an ASR of nearly 100% on Pix2Struct and just below 80% on Donut. The CDMG remains low for both models, as the remaining answers that are not affected by the attack tend to remain correct. The same trends, with a smoother behavior, can be seen in the ANLS-B (Figure 5), where the metric is high, confirming similarity of the answers to the targets, and ANLS-C, meaning that the remaining answers remain close to the ground truth. In the patch setting, the attack against Pix2Struct reaches an ASR close to 100%, while against Donut it succeeds in less than 20% of the cases. Again, the impact on the other answers is low for both models, and the correlation on the ANLS scores can be observed. *These results show that the attack is highly effective in optimizing a perturbation that induces the model to output specific answers for one selected QA pair.*

**Targeted Multi-Answer manipulation** ($B > 1$)**.** We analyze how the attack behaves when the objective includes multiple QA pairs simultaneously, i.e. $B > 1$. Moving along the x-axis of Figure 4, the attack's effectiveness decreases with the number of QA pairs used in the optimization objective, reflecting a progressively more challenging scenario. For Pix2Struct, the ASR remains high for small $B$ and gradually decreases as more QA pairs are included. The CDMG remains contained, indicating that the remaining QA pairs are barely impacted. The ANLS-B remains relatively high even as $B$ increases, indicating that the attack still achieves to improve similarity with the target answers across multiple QA pairs, but does not achieve exact match in many cases. The ANLS-C remains high even with the addition of more questions in the optimization, confirming that they are not affected and the document remains intact except for the targeted QA pairs. The effect of the perturbation on the remaining asnwers is visible through the increase of the CDMG and
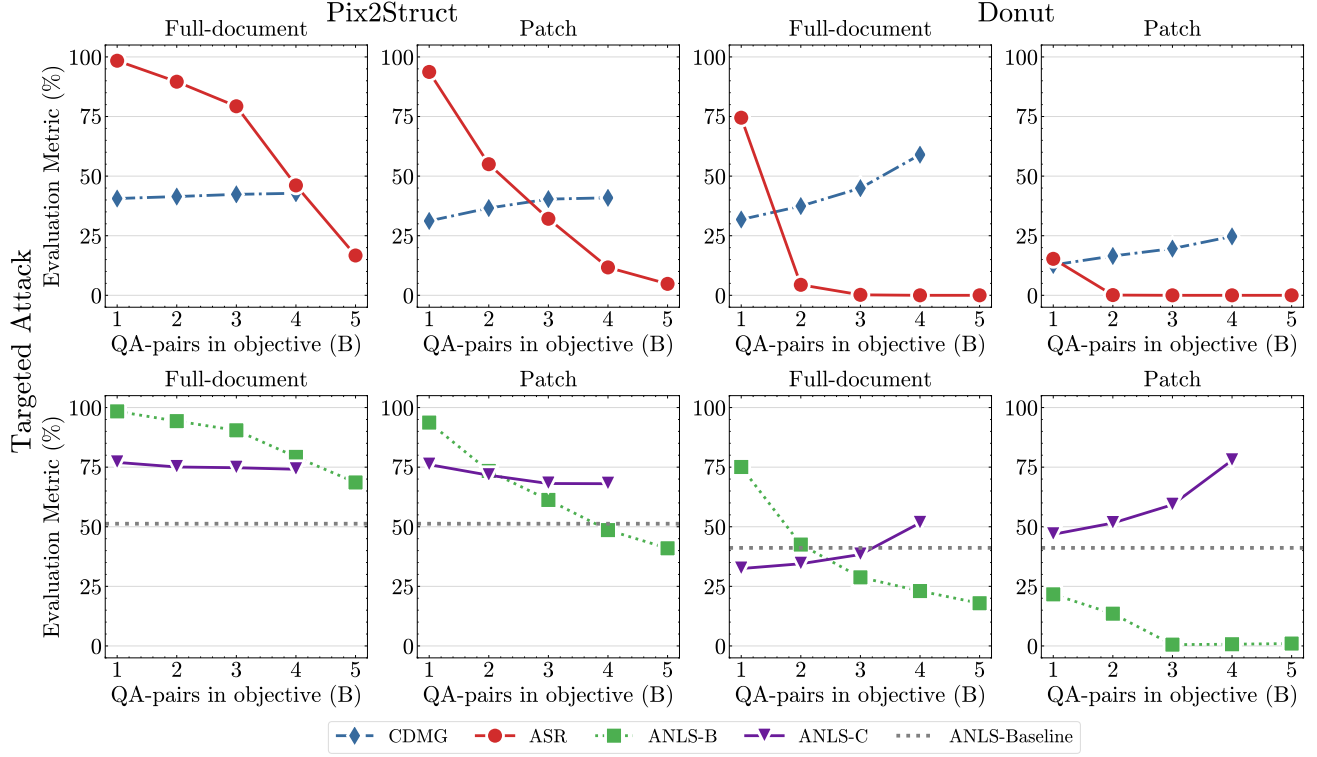
Figure 4: **Targeted attack results.** All metrics are reported for Pix2Struct and Donut across a different number of optimized QA pairs ($B$). Top panels show the ASR and the CDMG, while the bottom panels show the ANLS scores.
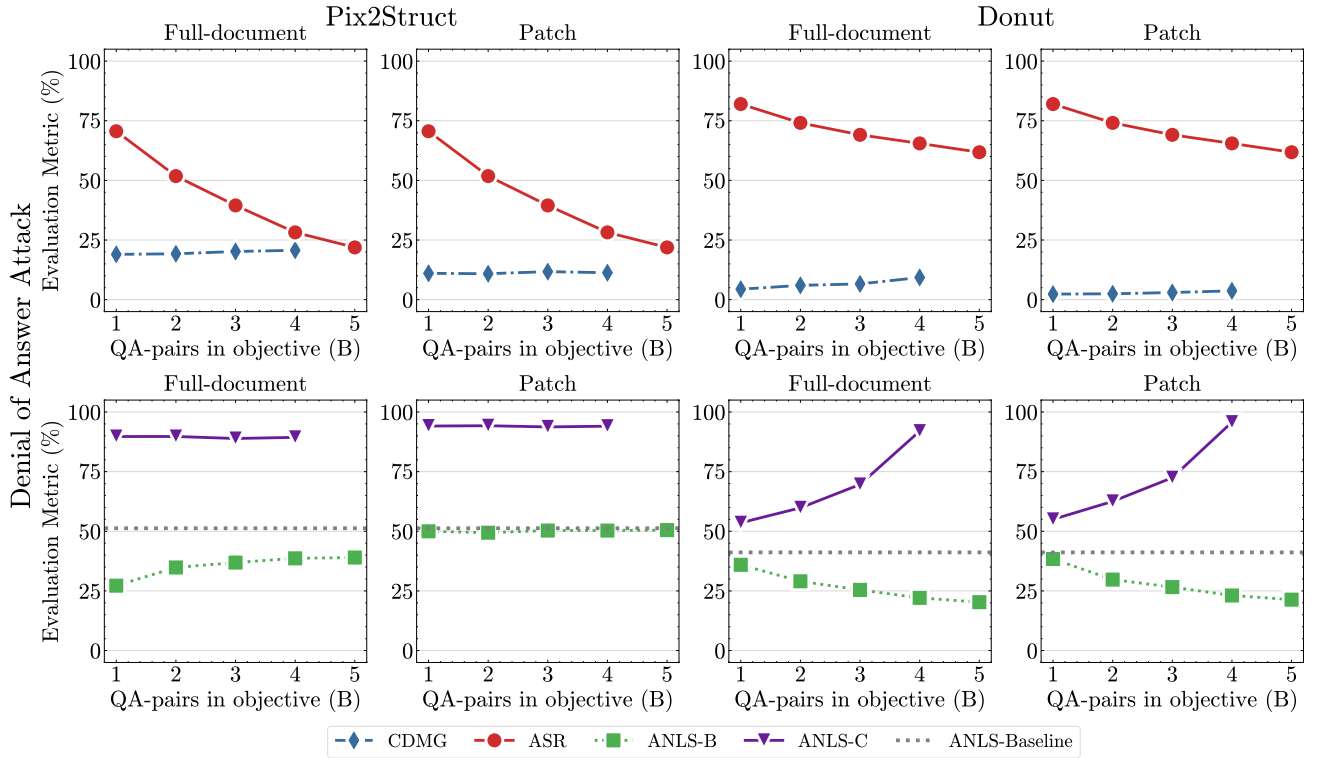


Figure 5: **Denial of Answer (DoA) results.** All metrics are reported for Pix2Struct and Donut across a different number of optimized QA pairs ($B$). Top panels show the ASR and the CDMG, while the bottom panels show the ANLS scores.

the drop in ANLS-C, which showcase how targeting more QA pairs makes the entire document more difficult to process for the models. For Donut, on the other hand, the ASR drops rapidly, approaching almost zero already at $B = 2$, and the attack is largely ineffective for higher values of $B$. However, the ANLS-B does not reach zero, which suggests that some generated answers still partially resemble the target sequences. Moreover, the effect on the CDMG is even more evident in this case, where the inclusion of more questions in the optimization tends to affect even QA pairs that are left out of the attack. In the patch attack scenario, a similar trend can be observed across all curves, but with lower impact on the model due to the more restrictive scenario (which modifies only a portion of the document rather the full image). For Pix2Struct, both the ASR and the ANLS-B are generally lower compared to the full-document setting. For Donut, the drop is even more pronounced, with ASR approaching $0$ already at $B = 2$ and ANLS-B decreasing as well. *Overall, these results show that jointly manipulating multiple answers substantially increases the difficulty of the attack.*

**Denial of Answer manipulation.** Next, we evaluate the Denial of Answer attack, i.e., the untargeted attack. Unlike the targeted attack, the objective here is to maximize the loss with respect to the ground truth answer(s), forcing the model to provide unspecified incorrect outputs. This simulates an attacker whose purpose is to corrupt correct information, without needing to orchestrate specific disinformation. The results of this experiment are shown in Figure 5. For the full-document attack, the method proves to be highly effective. The success rates in forcing a wrong answer (measured by ASR) is quite high with $B = 1$, but decreases when the number of answers is higher because the attack fails to disrupt *all* the QA pairs involved (see eq. 12). In this case, however, the drop in the efficacy of the attack for Donut is not the same of the targeted attack. This shows that untargeted disruption for this model is slightly easier, compared to the targeted case. The CDMG remains low for all conditions, meaning that the untargeted attack does not impact the remaining questions not directly optimized. The ANLS metrics further confirm these results. The ANLS-B in this case measures how close the outputs are from the GT answers for which we want to elicit wrong outputs, and it becomes significantly lower in the presence of the perturbation. The ANLS-C remains basically not impacted, except for Donut, in which they are impacted for $B = 1$, but become less affected as the number of targets increases. The patch attack scenario turns out to be more extremely similar to the full document setting, indicating that the additional constraint enforced by the patch is not limiting the effect of the attack. *These findings show that Denial-of-Answer attacks corrupt the targeted answers with minimal impact to the rest of the document.*

### 3.4. Qualitative Evaluation

In addition to quantitative metrics, we show a visualization of the attack and its effects in Figure 6. On the original, unaltered documents, the model correctly extracts information, answering questions with high accuracy (e.g. identifying "House Majority PAC" or the estimation identifier "8176"). However, when presented with the adversarial perturbation, the model's behavior changes substantially.

Although the perturbation is imperceptible to the human eye, and most importantly, it does not directly modify the textual answers present in the document content, the model is systematically driven into the set of our attack objective, i.e. $y^t = \{$"No Answer", "Unclear", "Retry", "Try later", "I won't tell you"$\}$, achieving an ASR of 100%.

## 4. Related Work

We outline here the most relevant related research, structured into attacks that aim to adversarially disrupt systems by attacking the preprocessing part (either the OCR in OCR-based systems or the downsampling of image data in general), adversarial attacks against the ML models in isolation, and other attacks against multimodal models.

**Adversarial attacks against the preprocessing.** Several works have studied how non-differentiable or lossy preprocessing steps can be exploited to attack vision models. Song et al. [16] demonstrate that Optical Character Recognition (OCR) systems are vulnerable to adversarial perturbations that survive text extraction, showing that subtle changes at the pixel level can drastically alter the recognized text. Beyond OCR, preprocessing steps such as resizing or compression can themselves be targeted. Xiao et al. [12] introduce image-scaling attacks, in which perturbations are crafted to remain invisible in the original resolution but become adversarial after rescaling. Quiring et al. [17] extend this analysis and provide a comprehensive taxonomy of vulnerabilities arising from image resampling, quantization, and similar operations commonly used in vision pipelines. These findings underscore that preprocessing in document understanding is considerably more complex than the canonical ImageNet setting, where preprocessing is typically limited to resizing and normalization. *In contrast to these works, we focus on end-to-end document understanding models where preprocessing is an integral part of the differentiable pipeline, and adversarial perturbations must propagate through complex visual–textual reasoning rather than isolated image transformations.*

**Attacks against multimodal models.** The increasing integration of multiple modalities into large language models has created a powerful new attack surface. While our work focuses on corrupting answers for specific questions in Document VQA, concurrent research highlights broader vulnerabilities. Bailey et al. [18] show that adversarial images can hijack the output of generative multimodal models at inference time, with the objective of steering high-level generative behavior. Similarly, Cui et al. [19] and Schlarmann et al. [20] conduct a systematic robustness study, confirming that LMMs are broadly vulnerable to adversarial images. In the domain of documents, recent works also explore multimodal inference risks: DocMIA [21] highlights that sensitive document content can be leaked through model outputs. However, these works primarily study broad alignment bypasses or information leakage, rather than targeted perturbations aimed at corrupting specific answers in DocVQA. Moreover, they don't target multiple answers at the same time. *In contrast, our work focuses on crafting end-to-end differentiable attacks that manipulate multiple fine-grained answers in the document analysis pipeline, highlighting a previously unexplored vulnerability surface for this domain.*
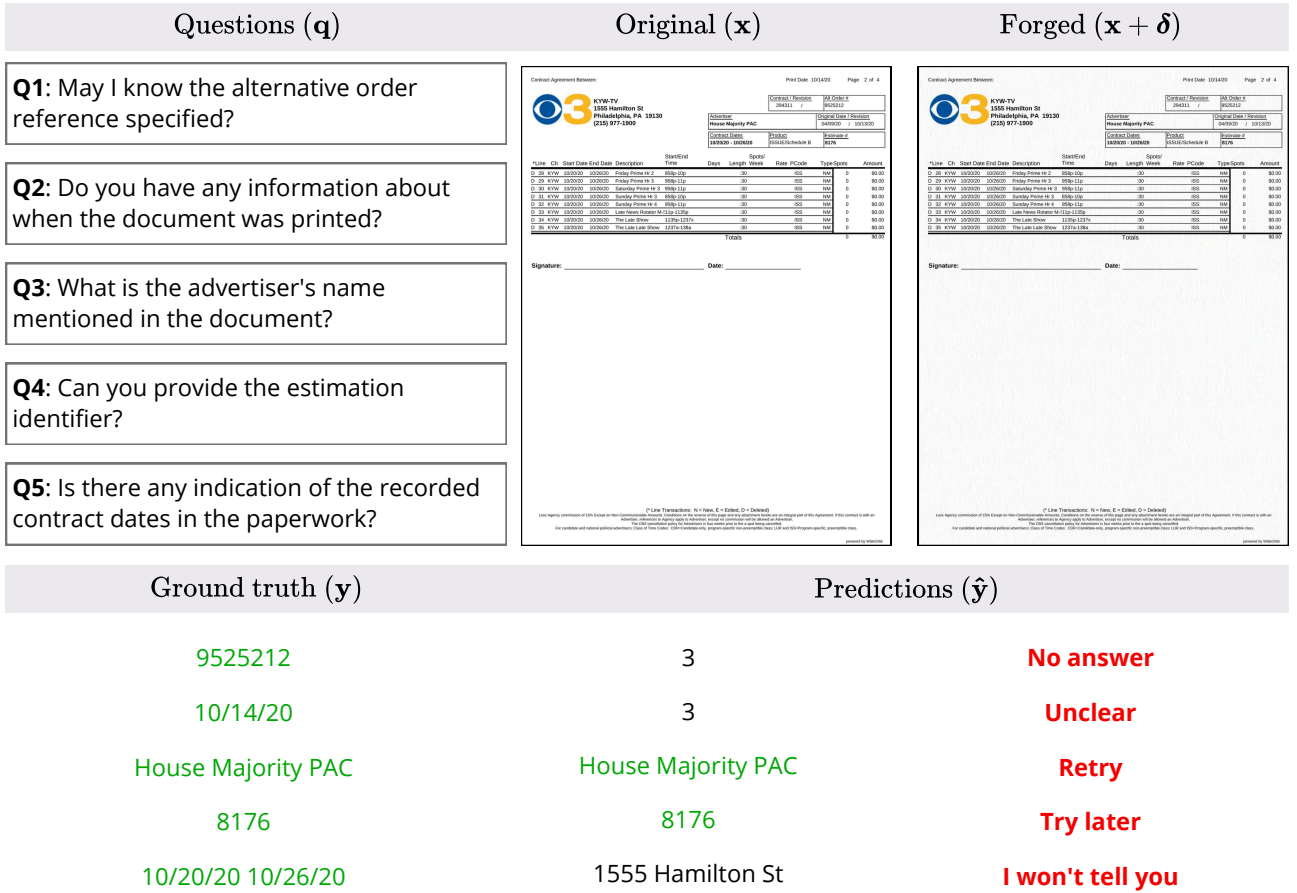
Figure 6: **Qualitative example.** Targeted multi-answer attack optimized with 5 questions on the full-document scenario.

**Jailbreaks in autoregressive models.** A complementary line of research studies adversarial attacks that exploit the autoregressive generation for LLMs to target alignment mechanisms of LLMs and LVMs. Instead of optimizing directly over the sequential distribution defined by the decoder, these methods require surrogate objectives for tractability. This perspective connects to gradient-based controllable generation attacks, such as Greedy Coordinate Gradient (GCG) [22], which operate directly on discrete token sequences to cause jailbreaks on LLMs. Carlini et al. [23] show that visual adversarial inputs can fully bypass alignment mechanisms of large multimodal models (LMMs), enabling them to generate harmful content despite safety training. *In contrast to jailbreak attacks on LLMs, that aim to bypass content filters and alignment mechanisms, our work focuses on controlling or disrupting specific downstream answers. Importantly, the resulting outputs are frequently plausible and policy-compliant in form, so they may not trigger alignment-based detectors.*

## 5. Conclusions

In this work, we have outlined the threat model for multimodal DocVQA systems, and demonstrated that they are vulnerable to end-to-end adversarial perturbations that can be crafted to produce incorrect or targeted answers. Our experiments show that small, visually imperceptible (full-document) or still inconspicuous and localized (patch) perturbations are often sufficient to substantially degrade model performance or achieve target responses, and that these attacks can disrupt multiple question–answer pairs on the same document simultaneously. These results highlight a previously underexplored attack surface on this domain.

We suggest several directions for future research. First, our primary results assume a fully white-box threat model; relaxing this assumption by averaging over variations of the same question or by optimizing against ensembles/surrogates of models would make the threat model more realistic. Second, many real-world document tasks operate on multi-page documents; extending attacks to multi-page DocVQA remains an important open problem. Third, developing cross-document or *universal* perturbations that generalize across documents (for example, via patch-based methods that resemble watermarks) is a promising but technically nontrivial direction.

Finally, future work should investigate practical defenses and robust evaluation protocols: realistic preprocessing hardening, model-side detection of anomalous inputs, adversarial training adapted to document modalities. We hope this study draws attention to the concrete risks of adversarial manipulation in DocVQA systems and stimulates development of principled, deployable mitigations.

## Proactive Prevention of Harm

Our work exposes vulnerabilities in end-to-end DocVQA systems with the aim of highlighting their security issues. Because such models are increasingly used

in high-stakes settings, we considered the potential risks of publishing these results. We believe that documenting these weaknesses and releasing the source code is an important step towards achieving true security: without a clear understanding of how current systems fail, developers and practitioners may deploy them without being aware of these potential issues. Our experiments are conducted only on public datasets, and we did not target, probe, or interact with any deployed systems whose compromise could have caused harm to third parties.

## Acknowledgments

## References

[1] M. Mathew, D. Karatzas, and C. Jawahar, "Docvqa: A dataset for vqa on document images," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 2200–2209.

[2] K. Lee, M. Joshi, I. Turc, H. Hu, F. Liu, J. Eisenschlos, U. Khandelwal, P. Shaw, M.-W. Chang, and K. Toutanova, "Pix2struct: screenshot parsing as pretraining for visual language understanding," in *Proceedings of the 40th International Conference on Machine Learning*, ser. ICML'23. JMLR.org, 2023.

[3] G. Kim, T. Hong, M. Yim, J. Nam, J. Park, J. Yim, W. Hwang, S. Yun, D. Han, and S. Park, "Ocr-free document understanding transformer," in *European Conference on Computer Vision (ECCV)*, 2022.

[4] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto, and F. Roli, "Evasion attacks against machine learning at test time," in *Machine learning and knowledge discovery in databases: European conference, ECML PKDD 2013, Prague, Czech Republic, September 23-27, 2013, proceedings, part III 13*. Springer, 2013, pp. 387–402.

[5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *2nd International Conference on Learning Representations, ICLR 2014*, 2014.

[6] A. F. Biten, R. Tito, A. Mafla, L. Gomez, M. Rusinol, E. Valveny, C. Jawahar, and D. Karatzas, "Scene text visual question answering," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 4291–4301.

[7] B. Biggio, F. Roli *et al.*, "Wild patterns: ten years after the rise of adversarial machine learning," *PATTERN RECOGNITION*, vol. 84, pp. 317–331, 2018.

[8] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," *arXiv preprint arXiv:1712.09665*, 2017.

[9] Q. Dong, L. Kang, M. Pintor, and D. Karatzas, *Position-Aware Stamp-Like Adversarial Attack for Document Classification*. International Conference on Document Analysis and Recognition, 09 2025, pp. 294–310.

[10] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *International Conference on Learning Representations*, 2018.

[11] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 9185–9193.

[12] Q. Xiao, Y. Chen, C. Shen, Y. Chen, and K. Li, "Seeing is not believing: Camouflage attacks on image scaling algorithms," in *28th USENIX Security Symposium (USENIX Security 19)*, 2019, pp. 443–460.

[13] R. Tito, K. Nguyen, M. Tobaben, R. Kerkouche, M. A. Souibgui, K. Jung, J. Jälkö, V. Poulain D'Andecy, A. Joseph, L. Kang, E. Valveny, A. Honkela, M. Fritz, and D. Karatzas, "Privacy-aware document visual question answering," in *Proceedings of the International Conference on Document Analysis and Recognition (ICDAR) 2024*, 2024.

[14] M. Mathew, V. Bagal, R. Tito, D. Karatzas, E. Valveny, and C. Jawahar, "Infographicvqa," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 1697–1706.

[15] K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, and M. Fritz, "Not what you've signed up for: Compromising real-world llm-integrated applications with indirect prompt injection," in *Proceedings of the 16th ACM workshop on artificial intelligence and security*, 2023, pp. 79–90.

[16] C. Song and V. Shmatikov, "Fooling ocr systems with adversarial text images," *arXiv preprint arXiv:1802.05385*, 2018.

[17] E. Quiring, D. Klein, D. Arp, M. Johns, and K. Rieck, "Adversarial preprocessing: Understanding and preventing Image-Scaling attacks in machine learning," in *29th USENIX Security Symposium (USENIX Security 20)*. USENIX Association, Aug. 2020, pp. 1363–1380. [Online]. Available: https://www.usenix.org/conference/usenixsecurity20/presentation/quiring

[18] L. Bailey, E. Ong, S. Russell, and S. Emmons, "Image hijacks: Adversarial images can control generative models at runtime," in *Proceedings of the 41st International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, Eds., vol. 235. PMLR, 21–27 Jul 2024, pp. 2443–2455. [Online]. Available: https://proceedings.mlr.press/v235/bailey24a.html

[19] X. Cui, A. Aparcedo, Y. K. Jang, and S.-N. Lim, "On the robustness of large multimodal models against image adversarial attacks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 24 625–24 634.

[20] C. Schlarmann and M. Hein, "On the adversarial robustness of multi-modal foundation models," in *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*. IEEE, 2023, pp. 3679–3687.

[21] K. Nguyen, R. Kerkouche, M. Fritz, and D. Karatzas, "Docmia: Document-level membership inference attacks against docvqa models," in *The Thirteenth International Conference on Learning Representations*, 2025.

[22] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," *arXiv preprint arXiv:2307.15043*, 2023.

[23] N. Carlini, M. Nasr, C. A. Choquette-Choo, M. Jagielski, I. Gao, P. W. W. Koh, D. Ippolito, F. Tramer, and L. Schmidt, "Are aligned neural networks adversarially aligned?" *Advances in Neural Information Processing Systems*, vol. 36, pp. 61 478–61 500, 2023.