# Robustness Evaluation of OCR-based Visual Document Understanding under Multi-Modal Adversarial Attacks

**Dong Nguyen Tien**
VinUniversity
Hanoi, Vietnam
25dong.nt@vinuni.edu.vn
ntdong@cmc.com.vn

**Dung D. Le**
VinUniversity
Hanoi, Vietnam
dung.ld@vinuni.edu.vn

## Abstract

Visual Document Understanding (VDU) systems have achieved strong performance in information extraction by integrating textual, layout, and visual signals. However, their robustness under realistic adversarial perturbations remains insufficiently explored. We introduce the first unified framework for generating and evaluating multi-modal adversarial attacks on OCR-based VDU models. Our method covers six gradient-based layout attack scenarios, incorporating manipulations of OCR bounding boxes, pixels, and texts across both word and line granularities, with constraints on layout perturbation budget (e.g., IoU $\geq 0.6$) to preserve plausibility.

Experimental results across four datasets (FUNSD, CORD, SROIE, DocVQA) and six model families demonstrate that line-level attacks and compound perturbations (BBox + Pixel + Text) yield the most severe performance degradation. Projected Gradient Descent (PGD)-based BBox perturbations outperform random-shift baselines in all investigated models. Ablation studies further validate the impact of layout budget, text modification, and adversarial transferability.

## 1 Introduction

Recent advances in *Visual Document Understanding* (VDU) have enabled automated information-extraction and question-answering pipelines for banking, taxation, legal compliance, and e-government services. Multimodal Transformers that jointly encode text, layout, and visual cues-such as LayoutLMv2 (Xu et al., 2022), LayoutLMv3 (Huang et al., 2022), DocFormer (Appalaraju et al., 2021), ERNIE (Zhang et al., 2019), and GeoLayoutLM (Luo et al., 2023)-achieve high positions in the leaderboards on datasets like FUNSD (Jaume et al., 2019), CORD (Park et al., 2019), SROIE (Huang et al., 2019), and DocVQA (Mathew et al., 2021).

**OCR-based VDU families.** One can categorize the existing OCR-based VDU models into one of two following branches, based on the required input modalities:

1) **Text + Layout + Image:** LayoutLMv2 (Xu et al., 2022), LayoutLMv3 (Lu et al., 2024b), ERNIE(Zhang et al., 2019) and GeoLayoutLM (Luo et al., 2023).

2) **Text + Layout:** LayTextLLM (Lu et al., 2024b) and a prompt-engineered Llama(Grattafiori et al., 2024) baseline that serialises every bbox as a single token.

**Research Gap and Motivation.** Prior work on VDU robustness has tackled distribution shift (He et al., 2023), image corruption (Chen et al., 2024a), or unicode text attacks (Boucher et al., 2023). However, none of the mentioned work has investigated a unified attacking framework of BBox, text, and image constraint-based perturbations. To this end, we introduce a multi-modal adversarial framework targeting OCR-based VDU models, with considerations of perturbation budgets of related modalities (BBox, text, and image).

OCR-based pipelines remain dominant in practice due to their traceability. APIs like APIs-Amazon Textract (Amazon Web Services, 2024), Azure Document Intelligence (Microsoft Azure, 2024), and Google Document AI (Google Cloud, 2024) return bounding boxes crucial for auditability in domains like finance. OCR-free models (Chen et al., 2024b; Liu et al., 2024a; Bai et al., 2025) are emerging but still lag in fine-grained extraction and spatial grounding.

Our unified attacks expose severe vulnerabilities in OCR-based models (up to 29.18% $F_1$ drop), highlighting the urgent need for robustness tools in layout-aware VDU systems.

We summarize our contributions as follows:

- **Unified Multi-Modal Attack Framework.** We propose the first framework for adversarial attacks on VDU models across *layout*, *text*, and *image* modalities under a shared budget. PGD-based layout attacks leverage a differentiable mIoU loss with IoU $\in \{0.6, 0.75, 0.9\}$.

- **Scenario-Based Benchmarking.** We define six attack scenarios across word- and line-levels, evaluated on four standard datasets: FUNSD (Jaume et al., 2019), CORD (Park et al., 2019), SROIE (Huang et al., 2019), DOCVQA (Mathew et al., 2021).

- **Robustness Analysis and Insights.** Our findings:
  - Line-level attacks consistently outperform word-level.
  - PGD is more effective than random shift, even under tight layout budgets.
  - Unicode diacritic attacks cause larger degradation than random text edits.
  - PGD attacks transfer well to models without visual input (e.g., LayTextLLM (Lu et al., 2024b), Llama (Lu et al., 2024b)).

## 2 Related Work

### 2.1 OCR-based VDU Models

OCR-based VDU models can be grouped by input modality as following:

**Text + Layout + Image** models including LayoutLMv2 (Xu et al., 2022), LayoutLMv3 (Huang et al., 2022), ERNIE-Layout (Zhang et al., 2019) and GeoLayoutLM (Luo et al., 2023)-concatenate token embeddings with 2-D positional encodings and CNN/ViT image features.

**Text + Layout** models such as LayTextLLM (Lu et al., 2024b) and a prompt-engineered LLaMA-3 (Grattafiori et al., 2024) baseline serialise each bounding box as an additional token, removing the image branch while retaining spatial structure. Although these architectures achieve state-of-the-art accuracy on FUNSD (Jaume et al., 2019), CORD (Park et al., 2019), SROIE (Huang et al., 2019) and DocVQA (Mathew et al., 2021), they still rely on *discretised bounding-box embeddings* that can be shifted at inference time -a vulnerability we exploit in this work.

### 2.2 Robustness and Adversarial Attacks

**Document-specific robustness.** Do-GOOD (He et al., 2023) probes distribution shift *but keeps the original bounding boxes intact*, leaving layout robustness unexplored. RoDLA (Chen et al., 2024a) introduces a robustness suite consisting of 12 image perturbations applied to documents. The work of (Qu et al., 2023) focuses on tampering detection and localization but does not generate adversarial examples.

**Multimodal LLM Safety.** Recent safety evaluations for multimodal large language models (MLLMs) expose vulnerabilities to image or prompt-based adversarial triggers. Key works include the safety benchmark for MLLMs (Liu et al., 2024b), ImgTrojan (Tao et al., 2025), and adversarial jailbreaks through visual examples (Qi et al., 2023). We also note ongoing work in 2025 exploring new variants of ImgTrojan and visual AE-based attacks, highlighting the lack of grounded layout-aware benchmarks.

**Text-in-image attacks.** (Boucher et al., 2023) present a genetic algorithm to perturb words inside images that fool OCR systems and ViT backbones, highlighting the brittleness of current OCR-based perception. These efforts underscore the need for a layout-aware, budget-controlled benchmark, which we introduce in this work.

**OCR-free vision–language models.** InternVL (Chen et al., 2024b), LLaVA-1.5 (Liu et al., 2024a), Qwen-VL (Bai et al., 2025) and DeepSeek-VL (Lu et al., 2024a) bypass explicit OCR by decoding text directly from pixels. Because they do not expose token-level bounding boxes, their attack surface is fundamentally different from that of OCR-based models. We therefore focus our robustness study on the still-dominant OCR-based family and leave the adaptation of our layout-budget concept to patch-level perturbations for OCR-free systems to future work (see Section 1 for a detailed motivation).

### 2.3 Bounding-Box Perturbations in Vision and Document Layout

**Generic object detection.** Distortion-Aware BBox Attack (Phuc et al., 2024) and ABBG (Nokabadi et al., 2024) perturb bounding boxes to mislead detectors and trackers. However, these methods are designed for natural images and do not consider textual semantics or document layout constraints. To the best of our knowledge, *no prior work systematically evaluates adversarial robustness of OCR-based VDU models to bounding-box layout shifts*-let alone in combination with pixel and text perturbations. We close this gap by propos-

ing a budget-controlled, multi-modal attack suite and a differentiable PGD strategy tailored to discrete layout embeddings.

## 3 Method

Our goal is to devise a unified **budgeted multimodal attack framework** that enables simultaneous perturbations of *layout*, *text* and *pixel* channels of OCR-based VDU models while keeping each perturbation within an interpretable budget. Figure 1 gives an overview.

### 3.1 Threat Model and Budget Constraints

We define a unified threat model $\mathcal{T} = (\mathcal{B}_{\text{layout}}, \mathcal{B}_{\text{text}}, \mathcal{B}_{\text{pixel}})$:

- **Layout budget:**

$$\mathcal{B}_{\text{layout}} : \text{IoU}(B, \tilde{B}) \geq \tau, \ \tau \in \{0.9, 0.75, 0.6\}.$$

IoU measures the overlap between the original bounding box $B$ and perturbed box $\tilde{B}$, ensuring layout perturbations remain spatially consistent. Adversarial bounding boxes $\tilde{B}$ are generated by either (i) randomly shifting the original box $B$ in one of four directions or scaling it slightly, or (ii) optimizing $\tilde{B}$ via PGD to minimize task loss, under the constraint of a minimum IoU.

- **Text budget:**

$$\mathcal{B}_{\text{text}} : \text{edit\_rate}(x, \tilde{x}) \in \{0, 0.1\}.$$

The edit_rate quantifies character-level replacements between clean text $x$ and adversarial text $\tilde{x}$, constrained to visually plausible noise. We randomly replace characters in the original text $x$ with other characters at a fixed rate (0.1 or change Unicode). We do not allow insertions or deletions, ensuring that token positions remain aligned.

- **Pixel budget:**

$$\mathcal{B}_{\text{pixel}} : T \in \mathcal{T}_{\text{RoDLA}}.$$

$\mathcal{T}_{\text{RoDLA}}$ denotes a transformation methods from RoDLA (Chen et al., 2024a), a set of 12 document-specific augmentations (e.g., blur, noise, occlusion). The image region inside a bounding box is first shifted according to the adversarial box $\tilde{B}$, then optionally augmented using one randomly sampled transformation from the 12 document-centric visual effects defined in Chen et al. (2024a) (e.g., blur, contrast, noise, shadow).

### 3.2 Learning-Based Bounding-Box Reparameterisation

**BBox Predictor.** To enable gradient-based layout attacks, we train a compact BBox predictor $g_\theta$ that maps each token embedding $e_i$ to a tuple $\langle c_x, c_y, \log w, \log h \rangle$. The architecture consists of a 2-layer MLP for input projection, a 4-layer Transformer encoder, and a 2-layer output MLP. We optimize using a combined SmoothL1 and GIoU loss:

$$\mathcal{L}_{\text{box}} = \mathcal{L}_{\text{SmoothL1}} + \lambda_{\text{GIoU}} \cdot \mathcal{L}_{\text{GIoU}}(\hat{b}, b) \quad (1)$$

where $\lambda_{\text{GIoU}} = 2.0$.

### 3.3 PGD with an mIoU-Budget Loss

Let $\hat{\mathbf{B}} = g_\theta(\mathbf{e})$ be the boxes predicted from the clean token embeddings $\mathbf{e}$. To craft an adversarial embedding $\tilde{\mathbf{e}}$ we *maximise*

$$\mathcal{L}_{\text{adv}} = \mathcal{L}_{\text{task}} - \lambda_{\text{box}} \big[ 1 - \text{IoU}(\hat{\mathbf{B}}, \tilde{\mathbf{B}}) \big],$$

where $\tilde{\mathbf{B}} = g_\theta(\tilde{\mathbf{e}})$ and $\lambda_{\text{box}} > 0$. At each PGD step ($T{=}10$, $\alpha{=}0.05$) we update $\mathbf{e}$ along $\nabla_{\mathbf{e}} \mathcal{L}_{\text{adv}}$ and **project** the resulting boxes back into the feasible set $\{\text{IoU} \geq \tau\}$, with $\tau \in \{0.9, 0.75, 0.6\}$. Among the ten candidates we keep the one that maximises $\mathcal{L}_{\text{task}}$ while still meeting the IoU budget.

### 3.4 Text and Pixel Modules

- **Text.** Two strategies: (i) random character replacement ($\rho = 0.1$), and (ii) Unicode-combining genetic optimisation (Boucher et al., 2023).

- **Pixel.** For any shifted box we (a) translate the enclosed pixels to match the new box and optionally (b) apply one RoDLA transform (Chen et al., 2024a), yielding visually coherent perturbations.

### 3.5 Attack Scenarios and Granularities

We evaluate six scenarios using both *word-level* and *line-level* granularity, where word-level refers to assessing the accuracy of individual words recognized by the OCR model, while line-level considers the correctness of entire lines of text, including word order and spacing.

**S1 BBox only**: shift bounding boxes; pixels and text frozen.

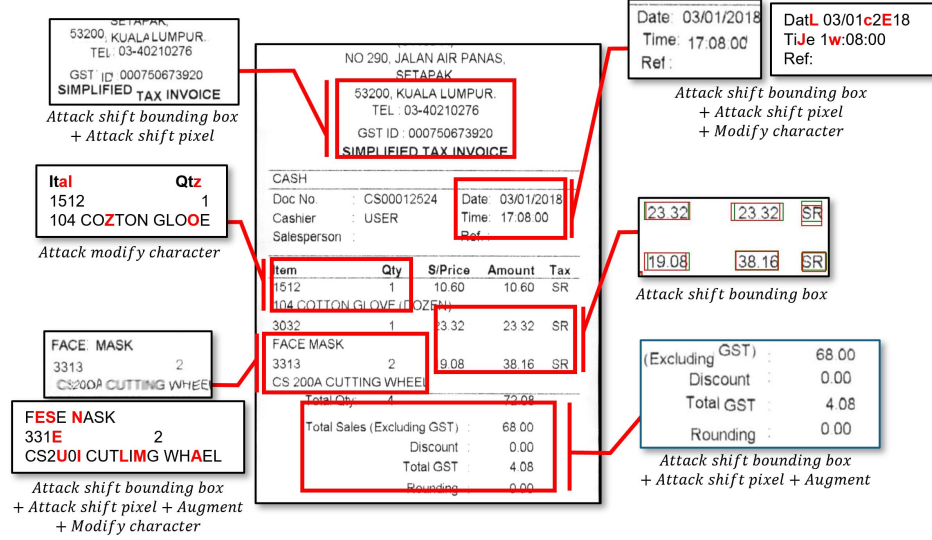**S2 BBox + Pixel**: shift boxes and translate the same pixel region.

Figure 1: Overview of the six proposed adversarial attack scenarios apply in each sample of the datasets

**S3 BBox + Pixel + Augment**: S2 plus one RoDLA transform (Chen et al., 2024a).

**S4 Text only**: mutate text under $\rho$; layout and image untouched.

**S5 BBox + Text**: change layout and text jointly.

**S6 BBox + Pixel + Text**: full multi-modal attack across all three channels.

### 3.6 Pipeline Overview

Figure 1 illustrates our modular attack pipeline. The attacker operates on OCR-based documents by modifying layout (via bounding box embeddings), text (via character-level or Unicode perturbations), and image content (via RoDLA-based pixel transformations). Perturbations are bounded by a unified budget, and gradient flow is enabled through a differentiable BBox-Predictor.

Given a document image and OCR output, the pipeline produces perturbed inputs that are fed into a frozen VDU model. The resulting adversarial document is then evaluated on key information extraction (KIE) or document visual QA metrics.

## 4 Experimental Results

### 4.1 Experimental Setup

We evaluate our adversarial framework on four widely used benchmarks: **FUNSD** (Jaume et al., 2019), **CORD** (Park et al., 2019), **SROIE** (Huang et al., 2019), and **DocVQA** (Mathew et al., 2021).

**Granularity.** All datasets provide ground-truth annotations at the **word level**. To simulate OCR

| Dataset | Level | Train BBoxes | Test BBoxes |
|---------|-------|--------------|-------------|
| FUNSD | Word | 21,888 | 8,707 |
| | Line | 7,259 | 2,270 |
| CORD | Word | 19,370 | 2,356 |
| | Line | 11,106 | 1,336 |
| SROIE | Word | 73,747 | 40,411 |
| | Line | 34,465 | 19,085 |
| DocVQA | Word | 6,202,284 | 937,786 |
| | Line | 1,806,853 | 259,582 |

Table 1: Dataset statistics at word and derived line level

post-processing under real-world conditions, we derive **line-level** segments by merging vertically aligned word boxes. This creates two granularity tiers-word and line-that reflect typical OCR system outputs. Table 1 shows the number of document images and bounding boxes at both granularities.

**Model and attack configuration.** All VDU models are finetuned for 100 epochs with AdamW (learning rate $2{\times}10^{-5}$, batch 32, weight decay $1{\times}10^{-2}$) on a single NVIDIA L40S (48 GB). Unless stated otherwise, we use text-edit budget $\rho = 0.10$ and layout budgets $\tau \in \{0.9, 0.75, 0.6\}$.

### 4.2 Bounding Box Predictor Evaluation

To enable PGD-based layout attacks, we train a lightweight 4-layer Transformer (§3.2) that maps encoder embeddings to bounding box parameters $(c_x, c_y, \log w, \log h)$. Training targets are derived from clean ground-truth annotations at word or line level.

We train a bounding box predictor for each model separately, using the spatial embeddings extracted from that model when fine-tuned on each corresponding dataset.

| Model | Gran. | FUNSD | CORD | SROIE | DocVQA |
|---|---|---|---|---|---|
| LayoutLMv2 | Line | 71.78 | 73.29 | 70.23 | – |
| LayoutLMv2 | Word | 66.45 | 69.45 | 67.95 | – |
| LayoutLMv3 | Line | **89.34** | **94.73** | **94.05** | – |
| LayoutLMv3 | Word | 84.55 | 93.17 | 91.41 | – |
| ERNIE | Line | 70.84 | 86.09 | 88.71 | – |
| ERNIE | Word | 72.84 | 80.74 | 88.95 | – |
| GeoLayoutLM | Line | – | – | – | – |
| GeoLayoutLM | Word | – | – | – | – |

Table 2: BBox prediction accuracy (mIoU, %) across datasets and granularities. "–" indicates cases where no model was trained, or where PGD adversarial samples could not be generated due to predicted boxes falling below the budget threshold (IoU $< 0.6$).

Table 2 reports mIoU of the bounding box predictors used to generate PGD adversarial examples. Since LayoutLMv3 (Huang et al., 2022) yields the most accurate bounding box predictions, we exclusively use its line-level bbox predictor to generate PGD adversarial samples across all evaluations. For other models, missing values ("–") indicate either no training was performed or the predicted boxes do not meet the minimum IoU budget constraint for PGD attack (i.e., no adversarial box can be constructed with IoU $\geq 0.6$).

## 4.3 Effectiveness of Budgeted Attacks

We comprehensively analyze how budget-controlled perturbations affect model robustness. We focus on six adversarial scenarios (S1–S6), and evaluate across attack methods (Random vs. PGD), granularities (word vs. line), and tasks (KIE vs. VQA). All attacks are conducted under a fixed budget (IoU $\geq 0.6$), and results are averaged over 5 random seeds.

### 4.3.1 Random Shift vs. PGD

We evaluate Random Shift and PGD attacks across six scenarios at the **line level**, using LayoutLMv3 (Huang et al., 2022) on four datasets. According to Table 3, PGD-based attacks consistently yield greater performance degradation than Random Shift in compound scenarios such as S5 and S6. This is attributable to the gradient-driven nature of PGD as well as the higher precision of the underlying bbox predictor. In our case, the BBox predictor from LayoutLMv3 embeddings achieves the highest mIoU, enabling more targeted and effective adversarial shifts. We observe similar trends on LayoutLMv2, ERNIE, and GeoLayoutLM, confirming the robustness and generality of these findings across architectures.

### 4.3.2 Line vs. Word Granularity (LayoutLMv3)

We compare line-level and word-level attacks on LayoutLMv3 (Huang et al., 2022) across all scenarios using FUNSD (Jaume et al., 2019). Table 4 shows the gap between line and word $F_1$ drops by scenario and granularity: Table 4 reports the performance gap between line-level and word-level attacks on LayoutLMv3 across all scenarios. Positive values indicate that line-level attacks are more damaging. We observe that line-based perturbations consistently lead to larger $F_1$ drops than their word-level counterparts across all compound scenarios (S2–S6), for both Random and PGD. The largest gap appears in scenario S6 (layout + pixel + text), with a 21.4 pp difference under Random shift and 13.4 pp under PGD. This reflects how line boxes—being longer and semantically denser—induce broader misalignment, impacting multiple tokens and layout cues simultaneously.

### 4.3.3 Cross-Task Robustness (KIE vs. VQA)

Table 3 further compares the performance of LayoutLMv3 under line-level perturbation across KIE and VQA tasks. KIE datasets suffer more under layout-pixel shifts (S1–S3), while VQA exhibits greater degradation under text-related attacks (S4–S6). This divergence reflects task-specific dependencies: KIE relies on precise layout structure, while VQA depends more on accurate textual content.

## 4.4 Transferability of PGD Attacks

We test the cross-model transferability of PGD attacks generated on LayoutLMv3 (Huang et al., 2022).

### 4.4.1 Transfer to Text + Layout + Image models.

These include LayoutLMv2 (Xu et al., 2022), ERNIE-Layout (Zhang et al., 2019), and GeoLayoutLM (Luo et al., 2023). Despite incorporating visual encoders, these models still suffer significant $F_1$ drops under PGD transfer. FUNSD (Jaume et al., 2019) remains the most vulnerable dataset, particularly for LayoutLMv2 and GeoLayoutLM (55.5% and 53.6% drop respectively). ERNIE-Layout is the most robust overall, with all drops under 7%, even under PGD.

Table 5 shows PGD consistently causes higher performance degradation than Random Shift, particularly in scenario S6. This reinforces the ef-

| | | FUNSD (F1 Drop %) | | CORD (F1 Drop %) | | SROIE (F1 Drop %) | | DocVQA (ANLS Drop %) |
|---|---|---|---|---|---|---|---|---|
| Model | Scenario | Rand. | PGD | Rand. | PGD | Rand. | PGD | Rand. |
| *LayoutLMv3* | | | | | | | | |
| | S1: BBox only | 7.94 | **13.32** | 1.28 | **4.77** | 0.46 | **5.59** | 0.10 |
| | S2: BBox + Pixel | 16.24 | 13.24 | 3.14 | **4.99** | 0.29 | **6.34** | 0.01 |
| | S3: S2 + Augment | 17.80 | 14.38 | 3.14 | **5.24** | 0.39 | **6.40** | 0.08 |
| | S4: Text only | 7.31 | – | 7.13 | – | 20.78 | – | 34.50 |
| | S5: BBox + Text | 16.55 | **22.78** | 11.67 | **11.80** | 23.02 | **26.12** | **35.75** |
| | S6: BBox + Pixel + Text | 28.91 | **29.18** | 13.23 | **18.37** | 23.42 | **28.18** | 35.42 |

Table 3: Drop in $F_1$ or ANLS (%) under line-level attack for all six scenarios. For FUNSD, CORD, and SROIE, each pair of columns represents Random and PGD variants. Bold values indicate strongest degradation across settings. PGD values will be filled in separately if available.

| Scenario | Gap (Rand) | Gap (PGD) |
|---|---|---|
| S1 | 7.60 | 11.98 |
| S2 | 15.61 | 11.77 |
| S3 | 15.93 | 13.29 |
| S4 | 0.09 | – |
| S5 | 8.56 | 13.04 |
| S6 | 21.37 | 13.44 |

Table 4: Gap in $F_1$ performance between line-level and word-level attacks on LayoutLMv3 (line drop minus word drop) under each scenario on FUNSD (Jaume et al., 2019). Positive values indicate stronger degradation from line-level attacks.

fectiveness of gradient-based optimization when paired with accurate bbox predictors. In this setting, LayoutLMv3 provides the highest mIoU predictor, enabling highly effective transfer of perturbations. Meanwhile, ERNIE-Layout appears more robust, potentially due to architectural differences or reduced reliance on spatial features.

### 4.4.2 Transfer to Text + Layout models

Despite lacking visual encoders, models like Lay-TextLLM (Lu et al., 2024b), LLaMA3 (Grattafiori et al., 2024), and ChatGPT 4.1 mini (OpenAI, 2024) still suffer non-trivial $F_1$ degradation under layout-based PGD attacks. Table 6 shows that PGD transfer leads to drops up to 3.4 pp on FUNSD and 7.9 pp on SROIE, while Random Shift causes negligible or even slightly negative effects. This confirms the strong cross-modality and cross-architecture transferability of PGD, even to models with no visual input or explicit layout supervision.

### 4.5 Ablation Study

We conduct additional ablation studies to analyze the sensitivity of VDU models to three key factors: bounding box budget, adversarial transferability, and text modification strategies.

| Model | Sce. | Rand. F1 Drop% | PGD Transfer F1 Drop% |
|---|---|---|---|
| *LayoutLMv2* | | | |
| FUNSD | S1 | 18.72 | **20.21** |
| FUNSD | S6 | 40.82 | **55.54** |
| CORD | S1 | 6.42 | **6.70** |
| CORD | S6 | 34.57 | **40.26** |
| *ERNIE-Layout* | | | |
| FUNSD | S1 | 5.03 | **7.50** |
| FUNSD | S6 | 4.51 | **6.72** |
| CORD | S1 | 2.16 | **3.00** |
| CORD | S6 | 3.10 | **3.33** |
| *GeoLayoutLM \*Word level* | | | |
| FUNSD | S1 | **1.14** | 0.87 |
| FUNSD | S6 | 49.08 | **53.56** |
| CORD | S1 | 0.04 | **0.58** |
| CORD | S6 | 10.44 | **12.75** |

Table 5: Transferability of LayoutLMv3-generated PGD examples to other *text+layout+image* models under scenarios S1 and S6 (line-level). Each block corresponds to a model with dataset and scenario breakdown. Values are absolute $F_1$ drops (%).

**Effect of Bounding Box Budget.** Table 7 indicates that lowering the IoU threshold increases attack strength for both PGD and Random methods. Notably, PGD maintains higher effectiveness even under stricter constraints (e.g., 6.5% drop at IoU 0.9 on FUNSD), while Random attacks quickly lose impact as the budget tightens (e.g., only 0.54% drop). This highlights PGD's ability to find optimal perturbations within a tight search space.

**Text Modification Strategies.** We compare two text-only attack methods: (1) random character replacement with 10% edit rate, and (2) a Unicode diacritic attack following (Boucher et al., 2023), which uses visually confusable glyphs via combining diacritical marks. Results in Table 8 show that Unicode attacks consistently cause greater performance degradation-up to 22.4% on CORD-highlighting their stronger disruption of semantic and visual consistency.

| Target Model | Rand. F1 Drop% | PGD Transfer F1 Drop% |
|---|---|---|
| LayTextLLM / FUNSD | 0.73 | **2.87** |
| LLaMA3-3B / FUNSD | -1.12 | **2.50** |
| ChatGPT 4.1 mini / FUNSD | 1.92 | **3.37** |
| LLaMA3-1B / FUNSD | -0.17 | **0.77** |
| LayTextLLM / SROIE | -0.05 | **0.93** |
| LLaMA3-3B / SROIE | 0.89 | **3.14** |
| ChatGPT 4.1 mini / SROIE | 4.34 | **6.83** |
| LLaMA3-1B / SROIE | 2.41 | **7.86** |

Table 6: Transferability of LayoutLMv3 (Huang et al., 2022) PGD examples to *text+layout* models under scenario S6 (line-level). PGD leads to greater degradation than Random shift, despite the lack of visual modality.

| Level | IoU Budget | Attack | F1 Drop% (FUNSD / CORD) |
|---|---|---|---|
| Line | 0.6 | Random | 7.94 / 1.28 |
| Line | 0.75 | Random | 2.94 / 0.22 |
| Line | 0.9 | Random | 0.54 / 0.15 |
| Line | 0.6 | PGD | 13.32 / 4.77 |
| Line | 0.75 | PGD | 6.60 / 0.92 |
| Line | 0.9 | PGD | 6.50 / 0.51 |

Table 7: $F_1$ drops under BBox attacks (scenario S1), across different IoU budgets and attack methods. PGD remains more effective under tight constraints.

# 5 Conclusion

We present a unified adversarial framework for evaluating the robustness of OCR-based Visual Document Understanding (VDU) systems across layout, pixel, and text modalities. Our method incorporates attack granularity (word vs. line) and budget constraints to simulate realistic perturbations.

Our contributions include: (1) a budgeted attack strategy combining layout, pixel, and text; (2) a differentiable BBox predictor enabling gradient-based layout attacks despite discrete spatial encodings; (3) a six-scenario benchmark over four datasets; and (4) in-depth analysis of modality combinations, granularity effects, and transferability.

Experiments reveal major vulnerabilities in state-of-the-art VDU models. Layout perturbations (S1) degrade performance significantly, with compound attacks (S6) amplifying this effect. PGD-based layout attacks consistently outperform random shift, particularly at tighter IoU budgets. Line-level attacks are more destructive than word-level due to broader context impact, and Unicode diacritic text attacks show stronger disruption than random edits. PGD samples also transfer well across models, including to those without image inputs.

## Limitations.

Our work focuses on OCR-based systems; models like LayTextLLM (Lu et al., 2024b) exhibit

| Model | Attack | FUNSD Drop% | CORD Drop% |
|---|---|---|---|
| LayoutLMv3 | S3 text only 0.1 | 7.31 | 7.13 |
| LayoutLMv3 | Unicode diacritic | 16.75 | 22.35 |

Table 8: Text-only attacks: Unicode-based modifications cause larger degradation compared to random character replacement.

resilience due to shuffled box encoding. We also do not evaluate OCR-free models (e.g., Qwen-VL (Bai et al., 2025), DeepSeek-VL (Lu et al., 2024a)), which rely on pure visual grounding and may require different attack strategies. Furthermore, our analysis is limited to white-box access; extending to black-box settings where only model outputs are observable remains future work.

## Ethical Considerations

Our study aims to evaluate the robustness of Visual Document Understanding systems by exposing vulnerabilities through controlled adversarial perturbations. All attacks are simulated under strict budget constraints to reflect plausible real-world issues. The adversarial examples are generated for research purposes only and are not intended for malicious use. No sensitive personal data are used or exposed in our experiments.

## References

Amazon Web Services. 2024. BoundingBox — Amazon Textract. https://docs.aws.amazon.com/textract/latest/dg/API_BoundingBox.html. Accessed: 2025-05-16.

Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R. Manmatha. 2021. Docformer: End-to-end transformer for document understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 993–1003.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. Qwen2.5-vl technical report. *Preprint*, arXiv:2502.13923.

Nicholas Boucher, Jenny Blessing, Ilia Shumailov, Ross Anderson, and Nicolas Papernot. 2023. When Vision Fails: Text Attacks Against ViT and OCR.

Yufan Chen, Jiaming Zhang, Kunyu Peng, Junwei Zheng, Ruiping Liu, Philip Torr, and Rainer Stiefelhagen. 2024a. Rodla: Benchmarking the robustness of document layout analysis models. In *CVPR*.

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,

Xizhou Zhu, Lewei Lu, and 1 others. 2024b. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.

Google Cloud. 2024. Enterprise Document OCR — Google Cloud Document AI. https://cloud.google.com/document-ai/docs/form-parser. Accessed: 2025-05-16.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Jiabang He, Yi Hu, Lei Wang, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. 2023. Do-good: Towards distribution shift evaluation for pre-trained visual document understanding models. *Preprint*, arXiv:2306.02623.

Yupan Huang, Yiheng Xu, Lei Cui, and et al. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Findings of the Association for Computational Linguistics: EMNLP 2022*.

Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and C. V. Jawahar. 2019. Icdar2019 competition on scanned receipt ocr and information extraction. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE.

Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. 2019. Funsd: A dataset for form understanding in noisy scanned documents. In *ICDAR*.

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024a. Llava-next: Improved reasoning, ocr, and world knowledge.

Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. 2024b. Safety of multimodal large language models on images and texts. *arXiv preprint arXiv:2402.00357*.

Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024a. Deepseek-vl: Towards real-world vision-language understanding. *Preprint*, arXiv:2403.05525.

Jinghui Lu, Haiyang Yu, Yanjie Wang, Yongjie Ye, Jingqun Tang, Ziwei Yang, Binghong Wu, Qi Liu, Hao Feng, Han Wang, and 1 others. 2024b. A bounding box is worth one token: Interleaving layout and text in a large language model for document understanding. *arXiv preprint arXiv:2407.01976*.

Chuwei Luo, Changxu Cheng, Qi Zheng, and Cong Yao. 2023. Geolayoutlm: Geometric pre-training for visual information extraction. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Minesh Mathew, Dimosthenis Karatzas, and C.V. Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*.

Microsoft Azure. 2024. Read API — Azure AI Document Intelligence. https://learn.microsoft.com/en-us/azure/ai-services/document-intelligence/prebuilt/layout?view=doc-intel-4.0.0&tabs=rest%2Csample-code. Accessed: 2025-05-16.

Fatemeh Nourilenjan Nokabadi, Jean-Francois Lalonde, and Christian Gagné. 2024. Adversarial bounding boxes generation (abbg) attack against visual object trackers. *Preprint*, arXiv:2411.17468.

OpenAI. 2024. Gpt-4.1-mini model overview. https://platform.openai.com/docs/models/gpt-4.1-mini. Accessed: 2025-05-20.

Seungjae Park, Seunghyun Shin, Bohyung Lee, and et al. 2019. Cord: A consolidated receipt dataset for post-ocr parsing. In *Document Intelligence Workshop at NeurIPS*.

Pham Phuc, Son Vuong, Khang Nguyen, and Tuan Dang. 2024. Distortion-aware adversarial attacks on bounding boxes of object detectors. *Preprint*, arXiv:2412.18815.

Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. 2023. Visual adversarial examples jailbreak aligned large language models. *Preprint*, arXiv:2306.13213.

Chenfan Qu, Chongyu Liu, Yuliang Liu, Xinhong Chen, Dezhi Peng, Fengjun Guo, and Lianwen Jin. 2023. Towards robust tampered text detection in document image: New dataset and new solution. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5937–5946.

Xijia Tao, Shuai Zhong, Lei Li, Qi Liu, and Lingpeng Kong. 2025. Imgtrojan: Jailbreaking vision-language models with one image. *Preprint*, arXiv:2403.02910.

Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, Min Zhang, and Lidong Zhou. 2022. Layoutlmv2: Multi-modal pre-training for visually-rich document understanding. *Preprint*, arXiv:2012.14740.

Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. Ernie: Enhanced language representation with informative entities. *Preprint*, arXiv:1905.07129.