

RoDLA: Benchmarking the Robustness of Document Layout Analysis Models (Supplementary Materials)

<https://yufanchen96.github.io/projects/RoDLA/>

1. Implementation Details

Hardware Setup. In this work, we have trained all models (including reproduced models) on machines equipped with four A100, each having 40 GB of memory. Each node would also with 300 GB CPU memory.

Training Settings. After input batch normalization, we have applied flip with a 0.5 flip ratio and crop with a crop size (384, 600) as data augmentation method. For a fair comparison in robustness benchmark for DLA models, we have trained all models (including reproduced models) using the same training strategy as in Table 1.

Table 1. Training settings.

Configurations	Parameter
Optimizer	AdamW
Learning Rate	$2e^{-4}$
Weight Decay	$1e^{-4}$
Scheduler	step-base
Training Epochs	24
Warm-up Step	{16, 22}
Warm-up Ratio	$1e^{-3}$
Batch-size per GPU	2

To create the benchmark, we have re-trained 38 models for this robustness benchmark for DLA models: On PubLayNet [17] dataset, we have re-trained 24 models (including ablation study). On DocLayNet [11] and M⁶Doc [4] datasets, we have re-trained 7 models each, as we have only re-trained the models with representative performance, *i.e.*, high **mRD** or **mAP** for specific perturbation, on the robustness benchmark for PubLayNet [17] dataset.

2. Detail of Perturbation Taxonomy

In this section, we provide more details about our 12 document image perturbations in 3 severity levels.

(P1) Rotation. We apply a random rotation to document images, along with corresponding annotations. The rotation operation on an image of a document is an affine transformation, mathematically described by a rotation matrix. If θ is the angle of rotation, the transformation for rotating a

point (x, y) around the origin is given by:

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}. \quad (1)$$

Here, (x', y') are the coordinates of the point after rotation. For L1, the θ is selected randomly from the range $[-5, 5]$. For L2, the θ is chosen randomly from $[-10, -5]$ or $[5, 10]$, each with 50% probability. For L3, the θ is taken randomly from $[-15, -10]$ or $[10, 15]$ for simulating real-world scenarios where object orientations vary.

(P2) Warping. We apply a pixel-wise displacement defined by a displacement field D . This field is typically generated using the Gaussian smoothing of random noise to simulate elastic deformation on document paper. The warping operation is as follows:

$$D(x, y) = \alpha \cdot G_\sigma(R(x, y)), \quad (2)$$

$$\begin{cases} x' = x + D_x(x, y) \\ y' = y + D_y(x, y) \end{cases}, \quad (3)$$

where $R(x, y)$ is a random field for displacement in both the x and y directions. G_σ is a Gaussian function with standard deviation σ ; the intensity or amplitude of the displacement is controlled by a factor α . D_x and D_y are the x and y components of the displacement field D .

(P3) Keystoning. We apply a 3D transformation to a 2D plane through a 3×3 matrix H , preserving lines but not necessarily the actual angles or lengths. This operation maps the homogeneous coordinates of a point in the source image to its new coordinates in the destination image:

$$\begin{pmatrix} x' \\ y' \\ w' \end{pmatrix} = H \cdot \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}. \quad (4)$$

Then the actual position in the transformed image is given by normalizing with w' by $(\frac{x'}{w'}, \frac{y'}{w'})$. The elements of H are typically derived from corner-point correspondences between the source and destination images. The coordinates of destination images are selected randomly from a Gaussian distribution centered around the original coordinates.

The standard deviation of the Gaussian distribution is determined by the level.

(P4) Watermark. The process of adding a watermark involves several steps, primarily dealing with image composition and potential rotation. The watermark image is rotated by a random angle θ from range $[0^\circ, 360^\circ]$ before the blending process. Then the watermark W is blended onto the original image O using a technique called alpha blending. The resulting pixel value I is calculated as:

$$I = \alpha_w \cdot W + (1 - \alpha_w) \cdot O, \quad (5)$$

where α_w is the transparency level of the watermark, which allows the original image to show through to varying severity levels.

(P5) Background. For complex background simulation, we overlay multiple images onto the original image. Before background alpha composition, multiple background images are resized and placed on a copy image B of the original image A . The placement is defined by the position (x_{pos}, y_{pos}) , which is randomly generated. The alpha composition can be described as:

$$I = \alpha_A \cdot A + (1 - \alpha_A) \cdot \alpha_B \cdot B, \quad (6)$$

where α_A and α_B are the alpha values of the original image and the background image, respectively.

(P6) Illumination. We introduce non-uniform illumination into document images, simulating effects such as shadows or glare. Mask M is created with random polygons filled with black on a white canvas, which is then blurred using a Gaussian filter. The illumination adjustment can be described mathematically as a pixel-wise multiplication of the image I with mask M :

$$I'(x, y) = V \cdot I(x, y) \cdot M(x, y), \quad (7)$$

where V is the illumination scaling factor, determined by the severity levels and type of illumination adjustment, *i.e.*, shadow with V_s and glare with V_l .

(P7) Ink-Bleeding. We apply an erosion operation for ink-bleeding simulation with an elliptical structuring element. The kernel size K_e determines the extent of erosion, depending on severity levels. The basic mathematical formula for erosion \ominus of an image A by a structuring element B is:

$$(A \ominus B)(x, y) = \min_{(b_x, b_y) \in B} \{A(x + b_x, y + b_y)\}. \quad (8)$$

To improve image quality during erosion, we upscale the image tenfold in both dimensions before applying the erosion. This is followed by downscaling to the original size, ensuring enhanced detail and quality in the final image.

(P8) Ink-Holdout. To simulate Ink-Holdout, which is the opposite of Ink-Bleeding, we use the dilation operation, the inverse of erosion. The parameters for the dilation process,

including the kernel size and the number of iterations, remain the same as those used for the erosion operation to maintain consistency in simulating these opposing ink behaviors. The mathematical formula for dilation \oplus of an document image A by a elliptical structuring element B is:

$$(A \oplus B)(x, y) = \max_{(b_x, b_y) \in B} \{A(x - b_x, y - b_y)\}. \quad (9)$$

(P9) Defocus. The simulation of defocus blur is inherently complex due to the variability of point spread functions (PSFs) within diverse photographic conditions. Nevertheless, given that document images are frequently captured at close quarters, it is feasible to approximate the PSF with a Gaussian kernel function for simulating defocus blur which demonstrated as follows:

$$I_{\text{defocus}}(x, y) = (I * G)(x, y), \quad (10)$$

with

$$G(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}. \quad (11)$$

where parameters of Gaussian kernel G are correspond to the level of severity which calibrated to manipulate the scope and the depth of field of the blur.

(P10) Vibration. Document vibration is simulated by motion blur. The kernel for motion blur is a matrix with non-zero values along a line. This line simulates the path of motion. The kernel for a horizontal motion blur is:

$$K = \frac{1}{n} \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}, \quad (12)$$

where n is the number of non-zero elements in the kernel. This kernel K is then rotated through a random angle θ within a predetermined range, which emulates the directional motion effect. Similar to defocus blur, the motion blur effect is applied using a convolution operation between the image and the rotated kernel:

$$B(x, y) = (I * K_{\text{rotated}})(x, y), \quad (13)$$

where I is the original image, K is the motion blur kernel, and B is the blurred image.

(P11) Speckle. Document speckle is generated by superimposing random light (background) and dark (foreground) blobs onto a document image. We generate random blobs based on density, size, and roughness through randomly placed points and Gaussian smoothing. These foreground and background blobs are combined with the original image I as:

$$I_{\text{modified}} = \min(\max(I_{\text{original}}, N_{\text{fg}}), 1 - N_{\text{bg}}), \quad (14)$$

where N_{fg} , N_{bg} represent the foreground and background blob noise intensity. In the mathematical simulation of speckle and blotch noise on document images, Gaussian distributed blob noise are generated within the image domain, modulated by a blob density factor D_b , which is parametrically governed by designated severity levels.

(P12) Texture. We have endeavored to replicate the texture interference patterns characteristic of document imagery. This approach aims to emulate texture interference by simulating the complex plant fiber structures historically present in archival documents. We have modeled the random walk of fiber paths as follows:

$$\text{FiberPath} = \left[\sum_{k=1}^n \cos(\theta_k) \cdot \delta, \sum_{k=1}^n \sin(\theta_k) \cdot \delta \right], \quad (15)$$

where θ_k are angles drawn from a Cauchy distribution, δ is the step length, and k is the step number. The final fibrous image is obtained by blending fibrous textures:

$$I' = (M \cdot I_{\text{ink}}) + ((1 - M) \cdot (1 - I_{\text{paper}})) \times 255, \quad (16)$$

where mask M determines the application of ink and paper textures to the original image. Within this simulation, the spatial distribution of the fibers predominantly conforms to a Gaussian distribution, thereby reflecting the randomness inherent in the physical composition of paper. To impart authenticity to the fiber noise and facilitate a more accurate representation of document wear and quality, we have modulated the fiber density across varying noise levels.

3. Evaluation Metrics

In this work, there are two types of evaluation metrics, including: (1) the metrics for quantifying the impacts of perturbations will be presented in Sec. 3.1, such as MS-SSIM, CW-SSIM, and our proposed mPE. (2) the metrics for assessing the robustness of models will be detailed in Sec. 3.2, such as mAP and our proposed mRD.

3.1. Details of Perturbation Evaluation Metrics

To elucidate the effects of different perturbations and compare perturbation evaluation metrics, we present detailed analyses in Fig. 1, showcasing the impact of various perturbation categories and levels on document images.

MS-SSIM & CW-SSIM. In our robustness benchmark, we utilize MS-SSIM (Multi-Scale Structural Similarity Index) and CW-SSIM (Complex Wavelet Structural Similarity Index) metrics, both widely recognized for assessing the similarity between two images and pertinent for evaluating the extent of information loss caused by such perturbations. These indices exhibit varying sensitivity to image perturbations, as in Fig. 1. However, in this study, we deviate from the conventional usage of MS-SSIM and CW-SSIM as mere

similarity measures. Given that these metrics yield a value of 100 for identical images, we propose using their complements relative to 100 to represent the loss of information, *i.e.*, $100 - f^{\text{MS-SSIM}}$ and $100 - f^{\text{CW-SSIM}}$. This approach enables a nuanced assessment of the impact of perturbations on document images, thereby enhancing the evaluation of model robustness in handling document perturbations.

mPE. The Mean Perturbation Effect (mPE) metric integrates the effects of image quality degradation and model performance reduction under various perturbations in DLA. Our mPE metric reveals a consistent trend, with an escalation in values corresponding to increased severity, particularly evident in Keystoning and Texture perturbations, as shown in Fig. 1. It highlights the compounded effects of perturbations, underscoring the importance of robustness in document analysis models. While all metrics show heightened impact with more severe perturbations, mPE uniquely captures the overall impact, serving as a dependable measure of model robustness against document perturbations and offering a comprehensive view of model robustness.

3.2. Details of Robustness Evaluation Metrics

mAP. The mean Average Precision (mAP) is a crucial metric in object detection, assessing a model’s performance across various classes. It is calculated as the mean of the Average Precision (AP) for each category, where AP is the area under the Precision-Recall curve.

$$\text{mAP} = \frac{1}{N} \sum_{i=1}^N AP_i. \quad (17)$$

Here, AP_i is the Average Precision for the i^{th} class. mAP is especially important in multi-class detection tasks with varying Intersection over Union (IoU) thresholds.

P-Avg. We introduce P-Avg (Perturbation Average), a novel metric based on the mAP framework, designed to evaluate a model’s robustness in document layout recognition across various levels and types of perturbations. P-Avg extends mAP to quantify a model’s ability to maintain recognition accuracy under diverse perturbation scenarios. Based on Eq. (17), the P-Avg can be mathematically expressed as:

$$\text{P-Avg} = \frac{1}{MN} \sum_{s=1}^M \sum_{p=1}^N \text{mAP}_{s,p}. \quad (18)$$

In this formula, s represents perturbation level, p represents perturbation categories, and $\text{mAP}_{s,p}$ is the mAP calculated for the s^{th} level of perturbation in the p^{th} category. This metric provides insights into the model’s adaptability and consistency in recognizing document layouts despite the presence of diverse and challenging distortions.

mRD. The mathematical underpinning of mRD pivots on the interplay between degradation D and the Mean Perturbation Effect (mPE). The metric is designed to normalize

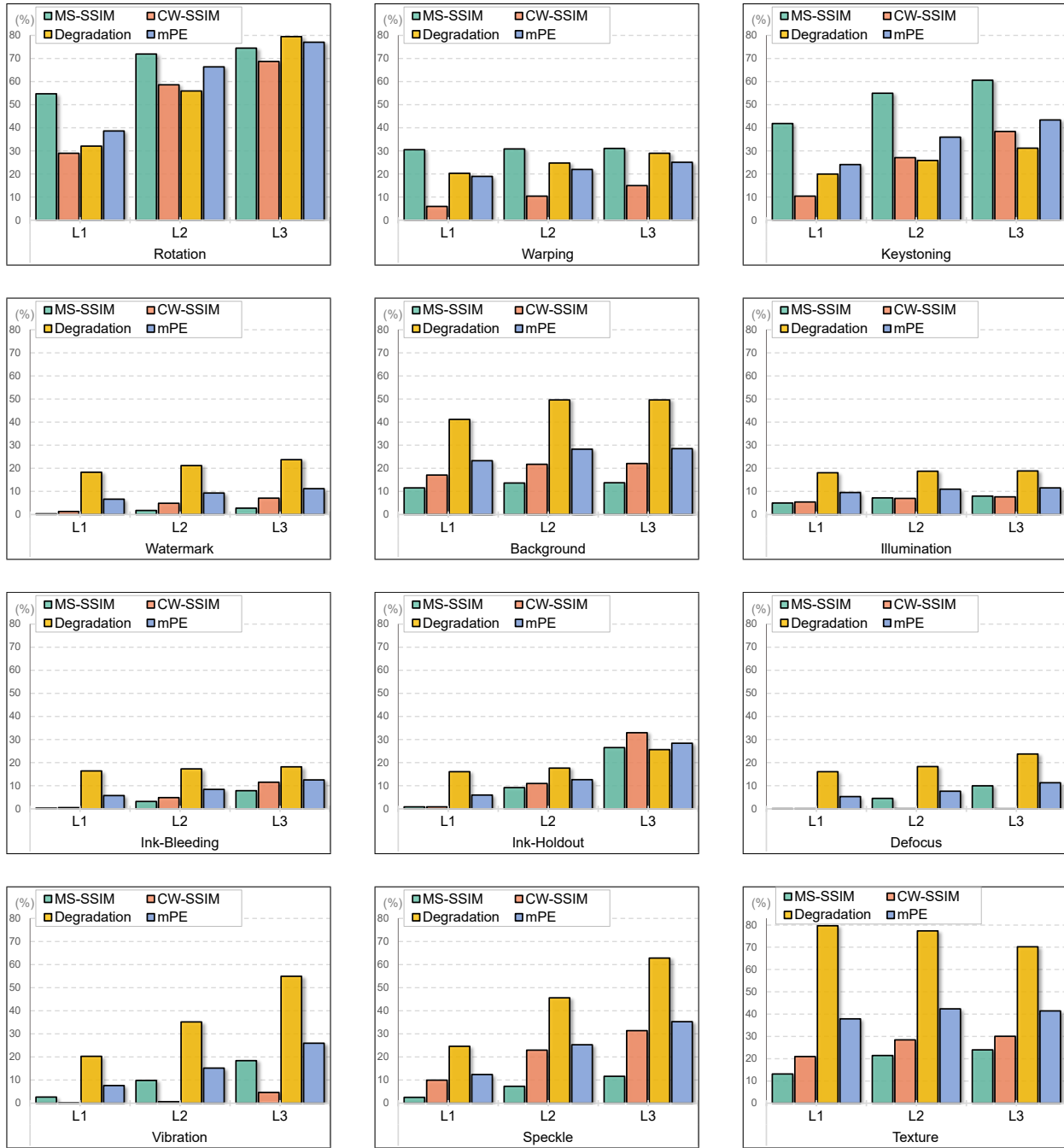


Figure 1. **Comparison between perturbation evaluation metrics** on 12 perturbation categories and 3 severity levels, including Image Quality Assessment methods (MS-SSIM and CW-SSIM), Degradation *w.r.t* a baseline, and the proposed mean Perturbation Effect (mPE). Other metrics cannot assess specific perturbations, for example MS-SSIM is insensitive to *warping* perturbation, and Degradation inversely measures *texture* perturbation across levels. In contrast, mPE is more balanced and inclusive to all perturbations and severity levels.

the degradation observed for a given perturbation by the perturbation’s inherent difficulty as captured by mPE. This normalization is crucial as it accounts for the perturbation’s baseline impact on the images, thus offering a relativized robustness measure. The degradation D represents how much a model’s performance deviates from its unperturbed state

when subjected to a specific perturbation and severity level.

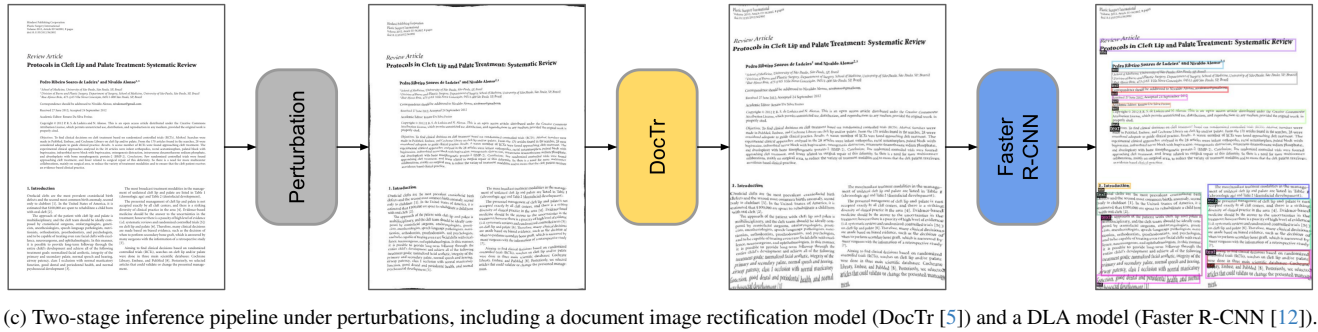
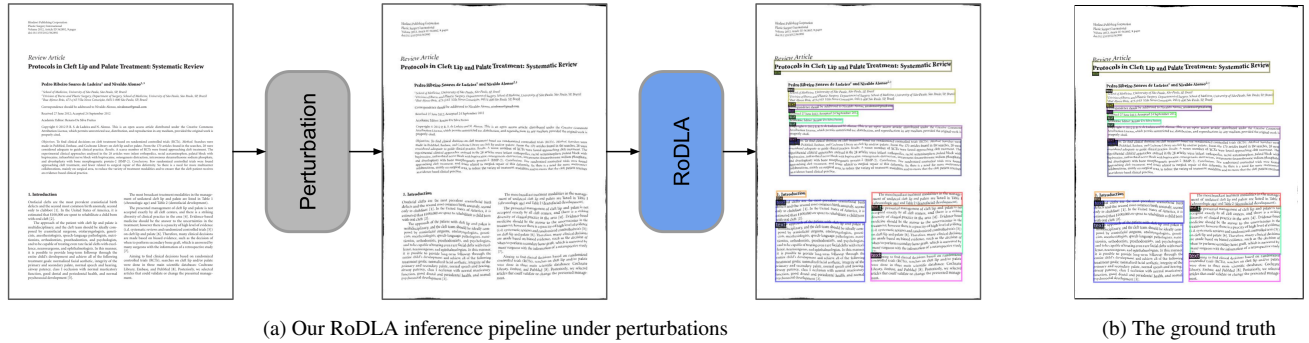


Figure 2. **Pipeline comparison** between (a) our one-stage inference (RoDLA) and (c) two-stage inference (first rectification, then DLA). Compared to (b) the ground truth, our RoDLA can obtain better DLA results.

Table 14. Result comparison between two-stage and one-stage pipelines. The **P-Avg** are evaluated on **PubLayNet-P**, **DocLayNet-P**, and **M⁰Doc-P** datasets. ‘-’ means no document image rectification model has been implemented. Here, **P-Avg** only refers to the result for four types of perturbations: Rotation, Warping, Keystoning, and Illumination.

Rectification Model	DLA Model	Clean	Rotation			Warping			Keystoning			Illumination			P-Avg \uparrow
			L1	L2	L3	L1	L2	L3	L1	L2	L3	L1	L2	L3	
PubLayNet-P															
DocTr [5]	Faster R-CNN [12]	90.2	37.5	50.6	49.6	18.9	19.8	20.0	20.6	27.1	31.4	74.1	72.7	72.1	41.2
-	Faster R-CNN [12]	90.2	67.9	44.1	20.6	79.7	75.2	71.0	80.0	74.1	68.8	81.9	81.3	81.1	68.8
-	RoDLA (Ours)	96.0	71.9	19.9	02.9	89.0	80.4	68.5	88.4	81.2	72.1	92.0	91.4	91.5	70.8
DocLayNet-P															
DocTr [5]	Faster R-CNN [12]	73.4	10.9	21.0	20.9	04.5	04.2	04.0	04.8	06.5	07.8	62.1	61.4	60.8	22.4
-	Faster R-CNN [12]	73.4	37.6	11.0	01.4	65.2	60.4	58.8	62.1	52.7	46.3	70.6	70.5	70.1	50.6
-	RoDLA (Ours)	80.5	49.6	17.8	04.3	72.6	64.2	59.4	73.2	65.8	59.1	80.3	80.2	80.0	58.9
M ⁰ Doc-P															
DocTr [5]	Faster R-CNN [12]	62.0	12.8	25.1	24.2	06.2	07.0	07.0	07.0	08.7	09.1	53.7	52.6	52.0	22.1
-	Faster R-CNN [12]	62.0	44.6	24.4	06.2	60.1	58.4	57.8	56.6	50.5	48.8	59.3	58.1	56.6	48.5
-	RoDLA (Ours)	70.0	58.4	40.9	23.7	68.2	66.0	64.0	66.3	63.4	60.9	67.5	68.2	67.8	59.6

tion from Inconsistency, Defocus from Blur, and Speckle from Noise. We further compare the visualization of predictions from Faster R-CNN [12], Mask R-CNN [7], and SwinDocSegmenter [1]. Observing the contrasts in Fig. 3, it becomes evident that our RoDLA model exhibits strong robustness against various perturbations, consistently yielding better predictions that align closely with the ground truth. In contrast, the other three models, *i.e.*, Faster R-

CNN [12], Mask R-CNN [7], and SwinDocSegmenter [1], demonstrate varying degrees of erroneous predictions when confronted with different perturbations. Notably, SwinDocSegmenter displays particularly pronounced inaccuracies under the influence of these perturbations. The qualitative analysis proves the effectiveness of our proposed RoDLA in enhancing the robustness of DLA.

7. Discussion

7.1. Limitations and Future Works

Our benchmark, designed for testing the robustness of Document Layout Analysis (DLA) models, currently simulates only a subset of perturbations commonly encountered in document images. It does not account for content tampering or content replacement, which could significantly impact model robustness. Additionally, our robustness benchmark is limited to only three severity levels. This granularity might be too coarse, and further subdivision into more nuanced levels could reveal subtler variations in DLA model robustness across different intensities of perturbations. In addition to our current robustness benchmark, we have only separately identified potential perturbations and assessed their impacts individually. We have not yet combined multiple perturbations on a single image with a certain probability, which would more accurately reflect real-world scenarios. A comprehensive evaluation of these combined perturbations is necessary for a more realistic assessment of the robustness of DLA. Moreover, in the robustness evaluation metric, we have incorporated only two image quality assessment methods and have used the performance of the Faster R-CNN [12] as the reference for calculating the Mean Perturbation Effect (mPE). To enhance the objectivity of mPE and more accurately reflect the impact of perturbations on document images, it would be beneficial to include a broader range of image quality assessment methods and performances from various models.

7.2. Societal Impacts

Our robustness benchmark and RoDLA model for document layout analysis bear significant implications. While demonstrating a strong ability to withstand various perturbations and achieving impressive robustness metrics, the current evaluation of our model is confined to three document layout datasets. This step is critical to ensuring its applicability in diverse practical scenarios, such as digitizing historical documents, streamlining administrative processes, or enhancing accessibility for visually impaired individuals. As we transition from controlled datasets to varied real-life environments, continuous model refinement is necessary to address any unforeseen challenges, ensuring the model’s relevance and effectiveness. Ethical considerations, particularly data integrity and impartiality, are paramount in avoiding biases and incorrect predictions that could have significant societal consequences. While our model demonstrates technical robustness, its deployment in real-world settings requires careful consideration of its broader societal implications to ensure it contributes positively and responsibly.

8. Acknowledgments

This work was supported in part by Helmholtz Association of German Research Centers, in part by the Ministry of Science, Research and the Arts of Baden-Württemberg (MWK) through the Cooperative Graduate School Accessibility through AI-based Assistive Technology (KATE) under Grant BW6-03, and in part by BMBF through a fellowship within the IFI programme of DAAD. This work was partially performed on the HoreKa supercomputer funded by the MWK and by the Federal Ministry of Education and Research, partially on the HAICORE@KIT partition supported by the Helmholtz Association Initiative and Networking Fund, and partially on bwForCluster Helix supported by the state of Baden-Württemberg through bwHPC and the German Research Foundation (DFG) through grant INST 35/1597-1 FUGG.

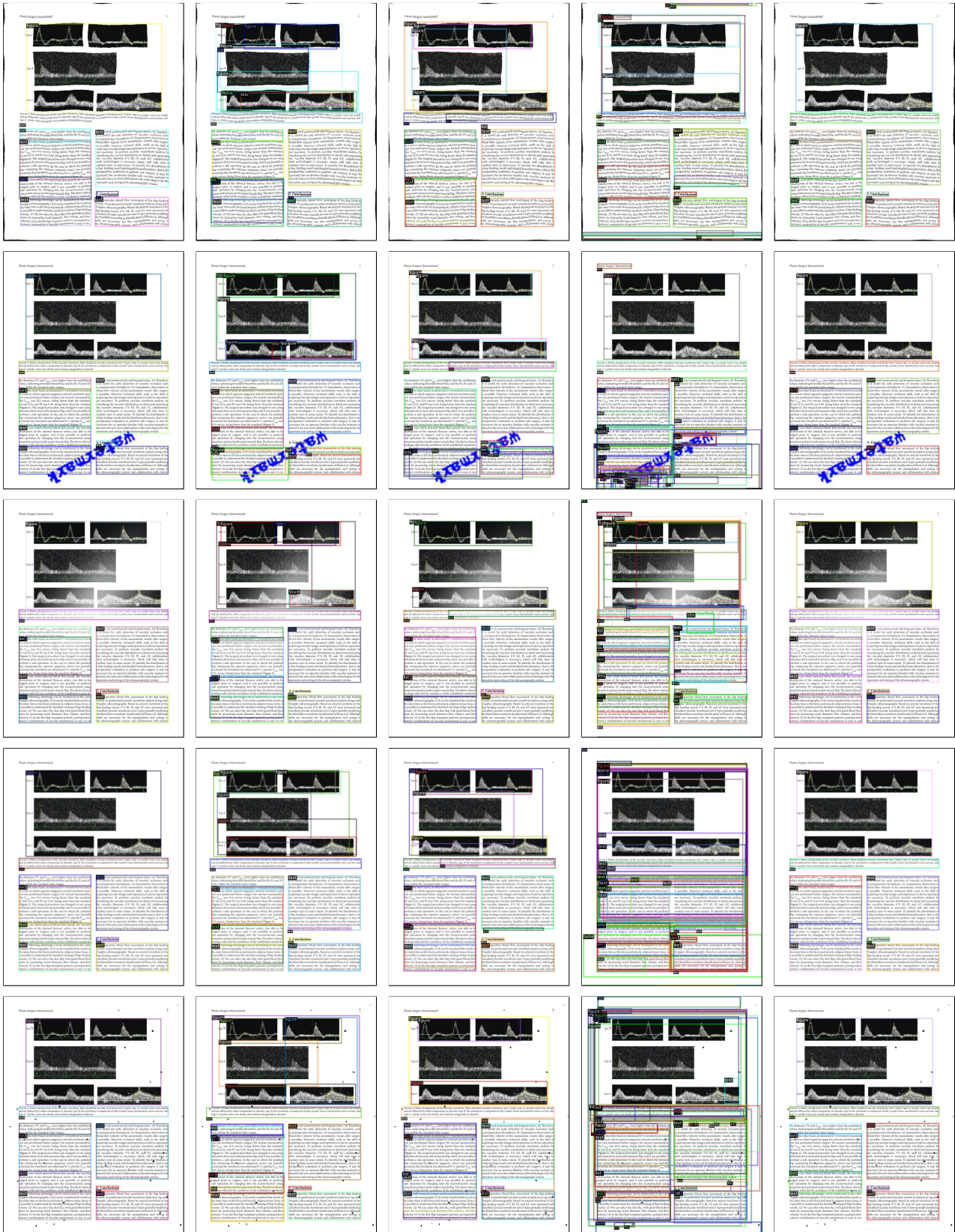


Figure 3. Visualizations on PubLayNet-P. From top to bottom: **Warping**, **Watermark**, **Illumination**, **Defocus**, and **Speckle**. From left to right: **the ground truth**, predictions from, *i.e.*, **Faster R-CNN** [12], **Mask R-CNN** [7], **SwinDocSegmenter** [1], and our **RoDLA**.

References

- [1] Ayan Banerjee, Sanket Biswas, Josep Lladós, and Umapada Pal. SwinDocSegmenter: An End-to-End Unified Domain Adaptive Transformer for Document Instance Segmentation. *ICDAR*, 2023. 5, 6, 7, 8, 10
- [2] Sanket Biswas, Ayan Banerjee, Josep Lladós, and Umapada Pal. DocSegTr: An Instance-Level End-to-End Document Image Segmentation Transformer. *arXiv preprint arXiv:2201.11438*, 2022. 5, 6, 7
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *CVPR*, 2018. 5, 6, 7
- [4] Hiuyi Cheng, Peirong Zhang, Sihang Wu, Jiaxin Zhang, Qiyuan Zhu, Zecheng Xie, Jing Li, Kai Ding, and Lianwen Jin. M6Doc: A Large-Scale Multi-Format, Multi-Type, Multi-Layout, Multi-Language, Multi-Annotation Category Dataset for Modern Document Layout Analysis. In *CVPR*, 2023. 1
- [5] Hao Feng, Shaokai Liu, Jiajun Deng, Wengang Zhou, and Houqiang Li. Deep Unrestricted Document Image Rectification. *arXiv preprint arXiv:2304.08796*, 2023. 6, 7, 8
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. 5, 6, 7
- [7] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017. 5, 6, 7, 8, 10
- [8] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *ACMMM*, 2022. 5, 6, 7
- [9] Junlong Li, Yiheng Xu, Tengchao Lv, Lei Cui, Cha Zhang, and Furu Wei. DiT: Self-supervised Pre-training for Document Image Transformer. *arXiv preprint arXiv:2203.02378*, 2022. 5, 6, 7
- [10] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9992–10002, Los Alamitos, CA, USA, 2021. IEEE Computer Society. 5, 6, 7
- [11] Birgit Pfizmann, Christoph Auer, Michele Dolfi, Ahmed S Nassar, and Peter Staar. Doclaynet: A large human-annotated dataset for document-layout segmentation. In *SIGKDD*, 2022. 1
- [12] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *TPAMI*, 2017. 5, 6, 7, 8, 9, 10
- [13] Zejiang Shen, Ruochen Zhang, Melissa Dell, Benjamin Charles Germain Lee, Jacob Carlson, and Weining Li. Layoutparser: A unified toolkit for deep learning based document image analysis. In *ICDAR*, 2021. 5
- [14] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *CVPR*, 2023. 5, 6, 7
- [15] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. *CoRR*, abs/1611.05431, 2016. 5
- [16] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. DINO: DETR with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 5
- [17] Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. PubLayNet: largest dataset ever for document layout analysis. In *ICDAR*, 2019. 1
- [18] Zhuofan Zong, Guanglu Song, and Yu Liu. DETRs with Collaborative Hybrid Assignments Training. In *ICCV*, 2023. 5