



Detection of Newspaper Layouts Using YOLO12

Atul Kumar^{1,*} and Gurpreet Singh Lehal²

¹Department of Computer Science, R.G.M. Government College, Joginder Nagar 175015, India

²Department of Computer Science, Punjabi University, Patiala 147002, India

Abstract

This study presents a robust and scalable method for automatic layout detection in digitized newspapers to facilitate efficient knowledge extraction and information retrieval. A custom dataset comprising annotated newspaper images in English, Hindi, and other languages was developed, with layout regions categorized into five primary classes. An enhanced YOLOv12 object detection model was trained on this dataset and evaluated using the mean Average Precision (mAP) metric across various Intersection over Union (IoU) thresholds. The model achieved a mAP@50 of 0.88, demonstrating strong detection performance and outperforming several state-of-the-art object detection models in the same task. The findings validate the effectiveness of the proposed approach in handling multilingual, structurally diverse newspaper formats. This research provides a practical framework for integrating automated layout analysis into digital archiving systems, OCR pipelines, and media monitoring applications. It also supports broader efforts to digitize historical print media and improve accessibility to regional content, thereby enabling enhanced research, journalism, and public engagement.

Keywords: newspapers, YOLO, segmentation, layout analysis.

1 Introduction

In the digital age, automatic analysis and digitization of printed media, particularly newspapers, are crucial for preservation, content extraction, and information retrieval. Conventional Optical Character Recognition (OCR) algorithms often prove inadequate when applied directly to complex newspaper layouts, which usually include multi-column formats, integrated graphics, a range of fonts, and various article structures. Accurate layout detection is required for effective downstream tasks like article segmentation, reading order reconstruction, and content categorization.

Recent advances in computer vision, particularly deep learning-based object detection models, have considerably improved document layout analysis. Among them, the You Only Look Once (YOLO) model family has emerged as a leading strategy because to its real-time performance and accuracy. The most recent generation, YOLOv12, improves architectural efficiency, spatial attention mechanisms, and multi-scale feature extraction, making it ideal for the complex layouts found in newspapers. The digitization of newspapers can be done with the help of OCR, a method of automatically identifying characters while optically scanning and digitizing text [1]. The analysis of newspaper layouts plays a pivotal role in understanding the structure and content of a newspaper page, allowing for a holistic examination



Submitted: 10 November 2025

Accepted: 31 December 2025

Published: 09 February 2026

Vol. 2, No. 2, 2026.

10.62762/TMI.2025.846033

*Corresponding author:

✉ Atul Kumar

atulkmr02@gmail.com

Citation

Kumar, A., & Lehal, G. S. (2026). Detection of Newspaper Layouts Using YOLO12. *ICCK Transactions on Machine Intelligence*, 2(2), 77–87.

© 2026 ICCK (Institute of Central Computation and Knowledge)

of its design and informational elements. The layout also depends on many parameters like editor's style, column size, types of images, advertisements, inverted text etc. It is so difficult to analyze the layout of these newspapers [2].

Despite the fact that researchers in this sector have recently established a number of methodologies to analyze newspaper layouts. The objective of this work is to perform a newspaper layout analysis, which has the following workflows:

1. Created a dataset specifically for newspaper layouts in COCO format, annotated with bounding boxes having five classes.
2. Training on the dataset based on the YOLOV12 model [3] architectures.
3. Comparison with other state-of-the-art models.

The remainder of this paper is organized as follows: the previous work related to domain, challenges in the layout analysis of Newspapers, methodology, data preparation, experiments, and conclusion.

2 Related Work

Over the years, considerable research has been conducted in the field of computer vision and image processing, leading to the development of numerous methods for newspaper layout analysis and recognition. With the advent of deep learning techniques, significant advancements have been made in solving the problem of newspaper layout analysis. By employing neural networks, it has become possible to accurately identify and extract different components of a newspaper page. However, the diverse designs of newspapers pose challenges in comparing to their layouts. To address these the layout and content of a newspaper, as well as its style and format, can all impact the variability. Newspaper layouts can be found in a variety of structures [4] as are rectangular, Manhattan, non-Manhattan, multi-column Manhattan. Ref. [5] analyzed compressively different layout types, discussed about framework of Layout analysis. The phases of the layout analysis algorithm are derived from an extensive review of the literature and have been comprehensively discussed in this study. Classical Top down, bottom up , hybrid approaches based on XY cut algorithm [6], Segmentation based on white lines [7], run length smoothing algorithm [8] proposed to segment blocks in documents. Deep learning techniques have showed promise for a range of applications, including headline recognition,

text segmentation, and object recognition. In view of this, Ref. [9] evaluated suggested rule-based approach for two machine learning models. A dataset of annotated newspaper images with diverse layouts [10], introduced layout segmentation as a crucial step before OCR, and has explored various state-of-the-art models [11] proposed a baseline model for text extraction and layout analysis based on a Faster R-CNN structure customized with finetuning. Ref. [12] proposed semi rule-based methodology. Ref. [13] proposed Manga Layout Analysis system based on Mask RCNN. Dessurt, a very simple document understanding transformer that can be tuned for a wider range of document applications, was presented in [14]. Ref. [15] utilized Mask RCNN for license plate detection, applies preprocessing steps to enhance the image quality, and employs Tesseract OCR [16] for character recognition. Ref. [17] proposed deep learning models for tasks like text detection, structure identification, also provide trained neural network models and Complete tools for effective data annotation on document images and model adjustment to accommodate various levels of customization. To enhance the performance of area recognition, YOLVO paired with a Feature Pyramid Network [18] to integrate features at multiple levels to recognise smaller regions. You Only Look Once (YOLO), one of CNN's most iconic depictions, is an original and uncomplicated solution to the object detection challenge [19]. A vision camera was used to capture the photos, which underwent various steps of image pre-processing using OpenCV-Python and recognition using Google Tesseract [20]. Ref. [21] proposed a semantic layout analysis method for layout anylsis. Ref. [22] refined the multimodal model by developing a multimodal method that combines textual and visual clues for the semantic segmentation of historical newspapers. Regarding the complexity of magazine and newspaper pages. Ref. [23] used faster RCNN fine-tuned model layout for analysis of Arabic books. Ref. [24] suggested a successful hybrid bottom-up/top-down technique. Ref. [25] showed that instead of only using recently introduced algorithms to solve old issues, the most recent study represents a substantial advancement because it now takes into account novel challenges and novel uses of cutting-edge techniques. A light version Network based on dilation [26] for layout analysis.

3 Challenges in Newspaper Digitization

Numerous problems could arise when analyzing the newspapers for digitization. The typical ones include:

1. Preprocessing the image is crucial before moving on to other stages due to the inherent low quality of old newspaper editions and potential digitization errors (the page is not positioned correctly on the scanner, or the scanner automatically selects the incorrect binarization threshold).
2. To further divide the columns and lines in a newspaper story with photographs, the text must be separated from the images. Due to the fact that photographs can span two or more columns, there is no longer a requirement for continuous white lines to separate the columns.
3. There is no standard layout format for newspapers; instead, each publication may have a unique layout. Because of this, it could be challenging to develop an algorithm that can accurately evaluate and extract information from different newspaper layouts.
4. Because newspapers typically employ images and graphics to convey information, automated systems find it challenging to evaluate and extract information from them.
5. Text may overlap in certain newspaper structures, making it difficult for an algorithm to discriminate between different parts and efficiently extract information.
6. Different language scripts are complex (other than English) and can be difficult for computer programmers to read, especially if the font is not conventional. This may lead to mistakes in text recognition and layout analysis.
7. It may be difficult for algorithms to establish the proper orientation and correctly extract the material because it can be written either vertically or horizontally.

4 Methodology

This study employs a pipeline for detecting textual and layout-based components in newspaper images using a deep learning approach. The complete methodology is divided into six key stages, as illustrated in the workflow Figure 1.

4.1 Data Acquisition

High-resolution newspaper page images were collected from a diverse range of sources, including digital newspaper archives, scanned printed editions,

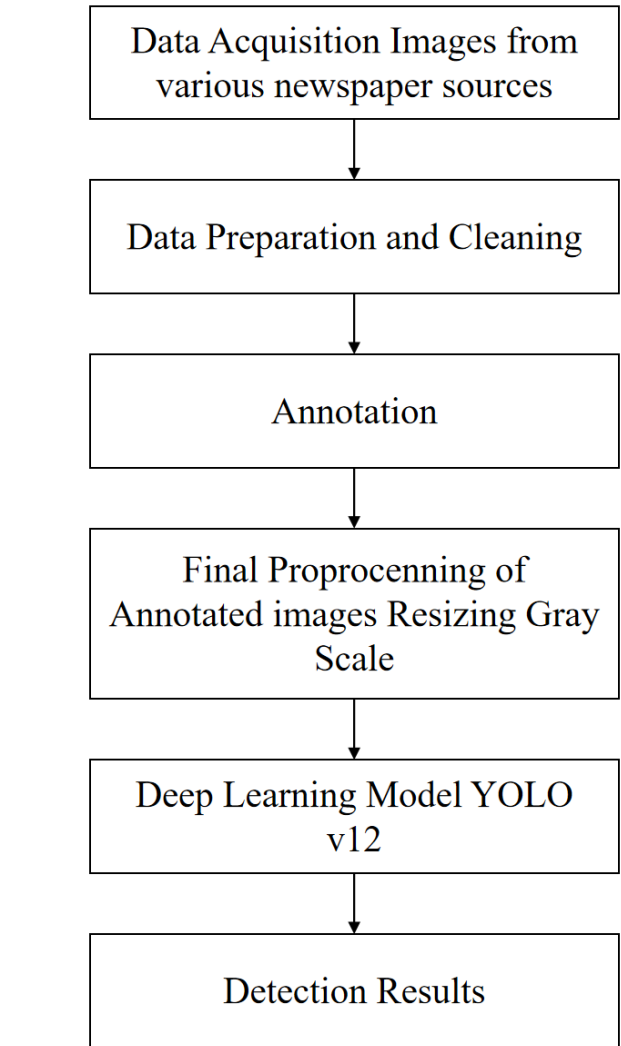


Figure 1. Flow diagram of newspaper recognition system.



Figure 2. Sample image of layout annotation using Roboflow.

and online repositories, as shown in Table 1. The goal was to compile a comprehensive and representative dataset encompassing various layouts, fonts, languages, and noise levels.



Figure 3. Raw annotation data structure in JSON format showing bounding box coordinates and metadata for layout components.

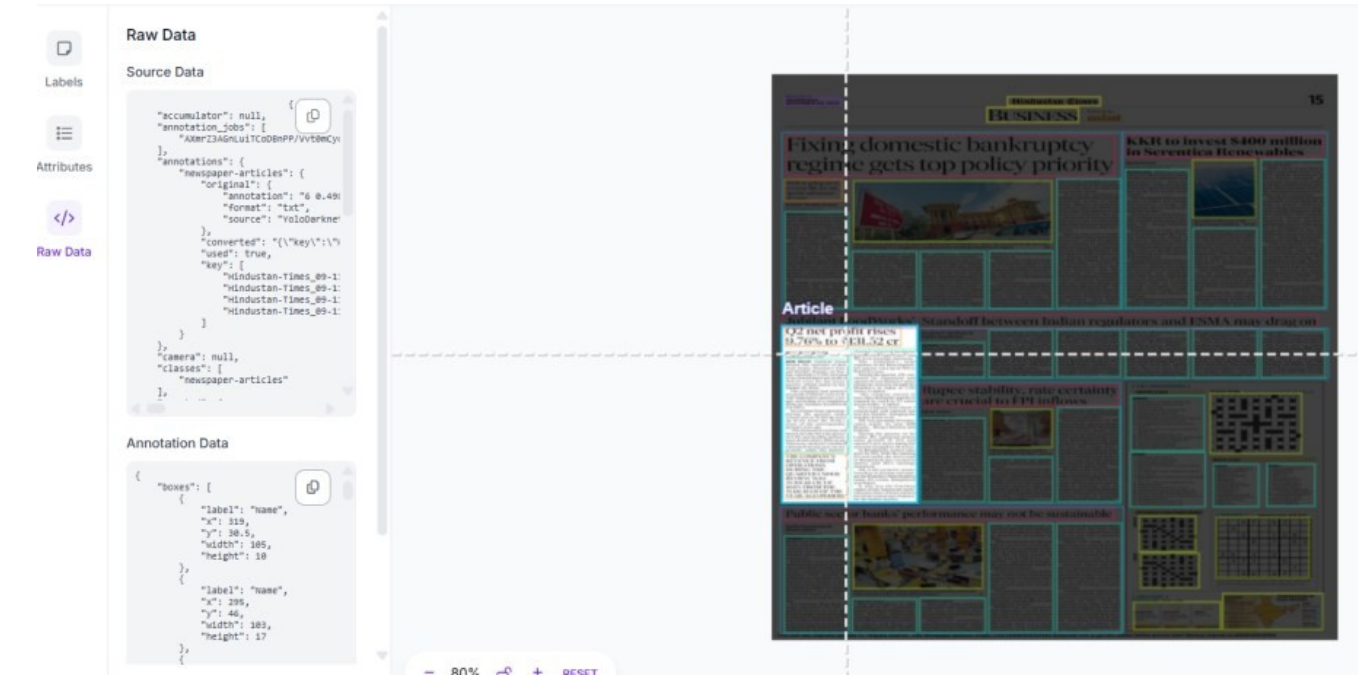


Figure 4. Annotation data example showing JSON structure alongside extracted newspaper headlines.

4.2 Data Preparation and Cleaning

The acquired raw images underwent preprocessing to ensure uniformity and reduce noise. This involved:

- Eliminating corrupted or low-quality images.
- Standardizing image formats and resolutions.
- Enhancing contrast and brightness where necessary.

- Eliminating irrelevant margins or marks using cropping techniques.

4.3 Annotation

Following the cleaning process, a manual and semi-automated annotation process was performed to label different newspaper elements such as Text, Headline, Image, advertisement, and caption. This annotation stage was crucial for supervised learning

Table 1. Collection of Newspapers from different resources.

Source	Number of Newspapers' images collected
Digitized archives	875
Local Newspapers from Libraries	450
Multilingual	175
Total	1500

Table 2. Number of annotations per Class.

Annotation Type	Count
Text	14352
Headline	8654
Image	9463
Advertisement	3212
Caption	6543
Total	42224

and was facilitated using tools like LabelImg or custom annotation platforms like Roboflow [27]. Figure 2, 3 and 4 shows the annotation using Roboflow. Number of annotation per class shown in Figure 5. The detailed count of annotations per class is provided in Table 2.

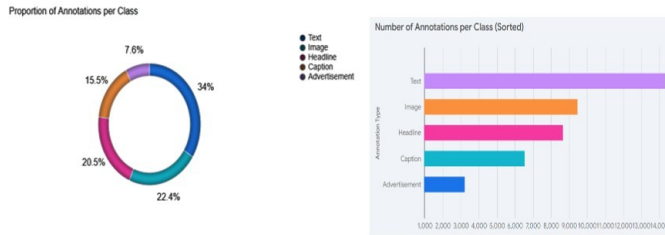


Figure 5. Annotation statistics for the newspaper layout dataset: class proportions and sorted counts.

5 Final Preprocessing

Annotated images were subjected to final preprocessing steps to optimize them for deep learning:

Resizing: All images were resized to a fixed dimension of 640x640 compatible with the YOLO V12 model, preserving aspect ratio as much as possible. To ensure robustness and generalizability, the dataset was curated to maintain a balanced distribution across categories and layout types, encompassing both single- and multi-column formats, image-heavy and text-dominant designs, as well as pages with varying advertisement densities. The annotation strategy focused on capturing the semantic and structural diversity of newspaper layouts, enabling the detection model to distinguish among visually similar but

contextually different components. The choice of input resolution plays a crucial role in achieving an optimal trade-off between model accuracy and computational efficiency in layout detection tasks. The resolution(640x640) was selected based on the YOLO architecture's inherent requirements, which operate most effectively with standardized square input dimensions. The 640×640 resolution offers a well-balanced design that maintains a manageable computational footprint for quicker training and inference while preserving important layout elements like headlines, column boundaries, and advertisement blocks. Experiments showed that 640×640 provided adequate spatial fidelity for the detection of macro-structural components typical of newspaper layouts, without significantly degrading accuracy, even though higher resolutions (e.g., 1024×1024) could capture finer text-level features.

Grayscale Conversion: Images were converted to grayscale to reduce computational complexity and emphasize textual and layout-related features over colour variations.

5.1 YOLOV 12 Architecture

YOLOv12 features an attention-centric architecture, designed to blend the speed of traditional CNN-based YOLO models with the performance benefits of attention mechanisms. Its core innovations include an optimized backbone leveraging Residual Efficient Layer Aggregation Networks (R-ELAN) for enhanced feature extraction and stable training, and a novel "area attention" module that strategically partitions feature maps to reduce computational overhead while maintaining a large effective receptive field. This architecture (Figure 6) also integrates Flash Attention for memory efficiency and employs 7x7 separable convolutions to further reduce computational burden, making it highly efficient for real-time object detection across various scales.

6 Experiments

6.1 Dataset Preparation

To train our YOLO model, we used a total dataset of 1,500 newspaper images. This dataset was partitioned into three subsets: 1,070 images (71.3%) for training, 300 images (20.0%) for validation, and 130 images (8.7%) for testing. The training set is used to teach the model, the validation set is used to tune hyperparameters and prevent overfitting, and the test set provides a final, unbiased evaluation of the model's performance. The complete distribution is detailed in

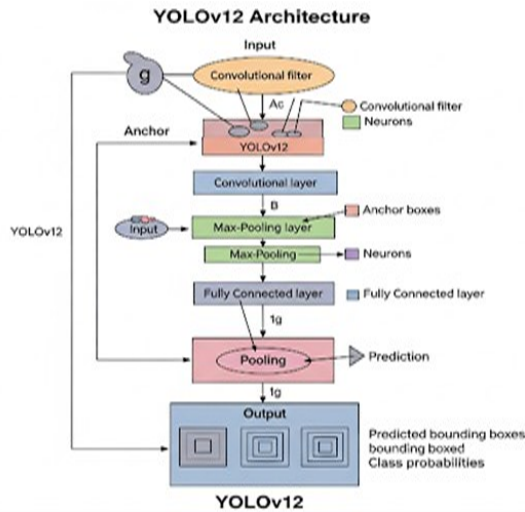


Figure 6. YOLOv12 architecture diagram showing the processing pipeline from input image to predicted bounding boxes and class probabilities.

Table 3 and Figure 7.

Table 3. Dataset distribution for Model.

Type	Number of Images
Training Set	1070
Validation Set	300
Testing Set	130
Total	1500

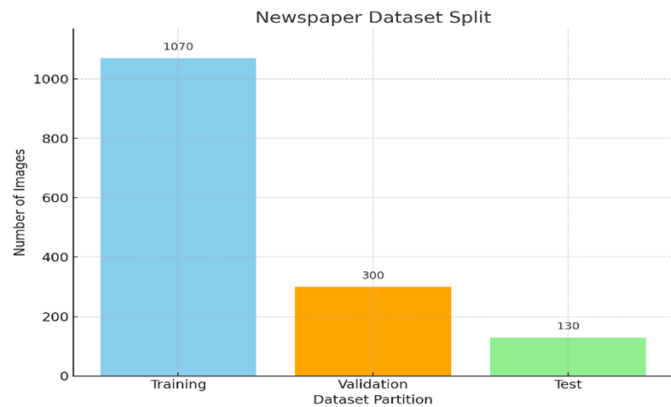


Figure 7. Dataset of Newspaper images with training, validation, and test.

6.2 Evaluation Metrics

The performance of the proposed model is assessed using a set of standard evaluation metrics, namely Precision (P), Recall (R), and Mean Average Precision (mAP), as shown in Table 4. The mathematical formulations used to compute each of these metrics are presented in Table 4, providing a thorough evaluation of both the model’s accuracy and computational efficiency.

6.3 Experimental Details

The experimental setup for this study is based on an Ubuntu 22.04 operating system, utilizing PyTorch version 2.6.0, Python 3.11.12, and CUDA 12.4. The hardware configuration comprised an A100 GPUs. Each model was trained for 200 epochs using the Stochastic Gradient Descent (SGD) optimizer, with an initial learning rate of 0.01 and momentum 0.9. Given the high resolution of images in the dataset, using the original dimensions would significantly extend training time, while excessively reducing the input size could compromise accuracy. All input images were uniformly resized to 640 × 640 pixels to balance computational efficiency and model performance.

6.4 Training

The model was trained for 200 epochs. The plots indicate effective model training over 200 epochs in Figure 8. Training and validation losses (box, classification, and DFL) consistently decrease, demonstrating stable convergence and improved model fitting.

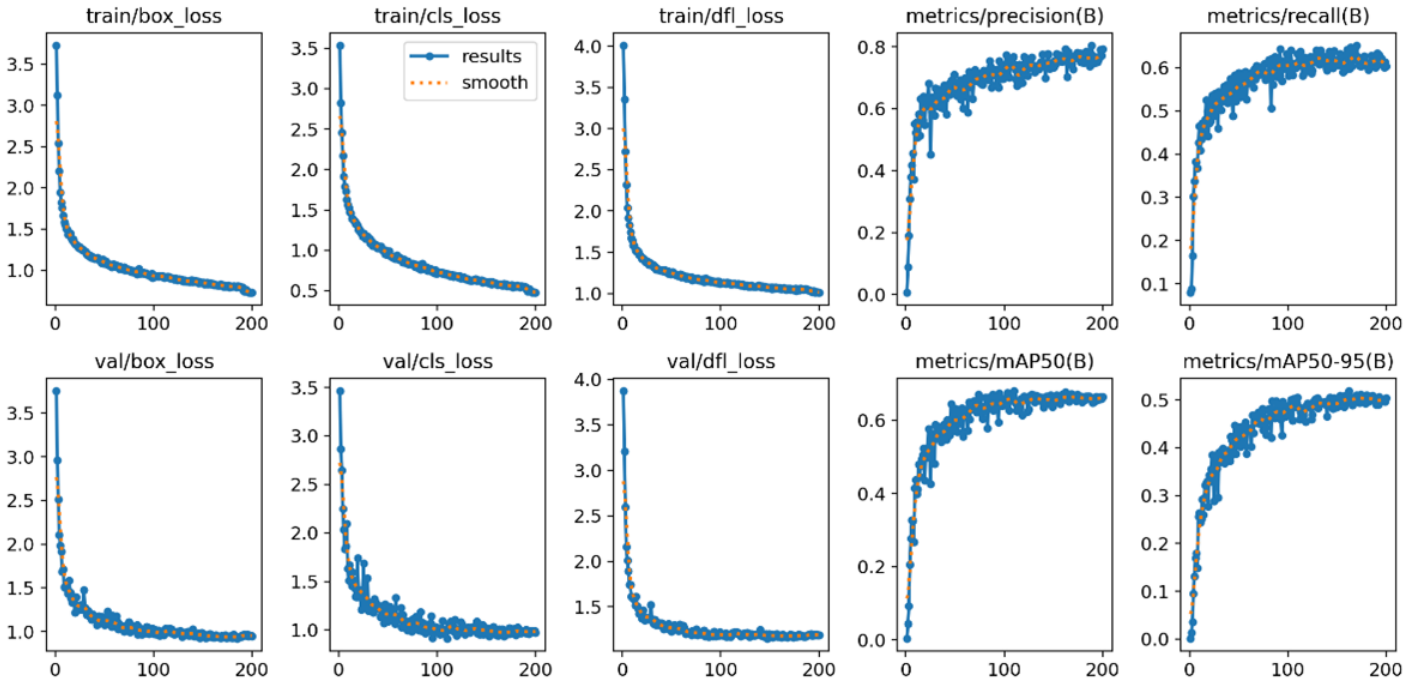
We have tested the model performance for around 130 newspaper images of test set as shown in Table 5.

To assess the effectiveness of the proposed method, the enhancement strategy was applied to various YOLOv12 model variants and compared against their respective baseline counterparts. The experimental outcomes are presented in Table 6, with the highest-performing results for each variant emphasized in bold. A detailed analysis of the results demonstrates that YOLOv12 consistently surpasses the baseline models in terms of both detection accuracy and model efficiency.

Table 6 compares three YOLOv12 model variant: YOLOv12n, YOLOv12s, and YOLOv12m based on their detection performance and model complexity. YOLOv12n, the most lightweight model with 2.6 million parameters, achieves a mAP@50 of 71% and mAP@50–90 of 76%, making it suitable for resource-constrained environments. YOLOv12s, with 9.3 million parameters, offers improved accuracy with a mAP@50 of 74% and mAP@50–90 of 81%, balancing performance and efficiency. YOLOv12m, the most complex model at 20.2 million parameters, attains the highest mAP@50 of 78% but slightly lower mAP@50–90 of 80%, indicating strong performance at lower IoU thresholds.

Table 4. Summary of key evaluation metrics.

Metric	Formula	Purpose
Precision	$\text{Precision} = \frac{TP}{TP+FP}$	Measures the correctness of positive predictions
Recall	$\text{Recall} = \frac{TP}{TP+FN}$	Measures coverage of actual positives
Average Precision	$AP = \frac{1}{n} \sum_{i=1}^n AP_i$	Area under the precision-recall curve for a class
Mean Average Precision	$mAP = \frac{mAP}{n}$ $= \frac{1}{n} \sum_{i=1}^n AP_i$	Mean of AP over all classes

**Figure 8.** Training and validation performance metrics of the object detection model.**Table 5.** Detection of the Model on Different Categories

Category	Original category(Count)	Detected category(Count)	Accuracy (%)
Text	1563	1430	91.49
Headline	476	438	92.01
Image	520	504	96.92
Advertisement	76	69	90.78
Caption	274	252	91.97

Table 6. Performance comparison between YOLO12n, YOLO12s and YOLO12m.

Model	mAP50(%)	mAP50-90(%)	Params
YOLO12n	0.71	0.76	2.6
YOLO12s	0.74	0.81	9.3
YOLO12m	0.78	0.80	20.2

6.5 Comparative Analysis and Model Performance

YOLOv12's performance was benchmarked against several state-of-the-art models on the same dataset. To thoroughly assess the performance of YOLOv12, we compared this against several state-of-the-art lightweight networks, including Faster RCNN, Mask RCNN, YOLOv5, YOLOv8, YOLOv10, and YOLOv11 [30]. The comparison results are summarized in Table 7, with the top two performers highlighted in red and blue on the similar dataset. Table 7 shows performance of YOLOv12's different versions. YOLOv12 achieved the highest overall accuracy and efficiency due to its attention-based feature extraction and optimized multi-scale fusion. These results demonstrate a

Table 7. Comparison of YoloV12 with Other State-of-the-art models.

Model	mAP@50 (%)	mAP@50–90 (%)	Params (M)
Faster R-CNN [31]	0.69	0.70	42
Mask R-CNN [29]	0.73	0.76	45
YOLOv5	0.79	0.80	7.2
YOLOv8	0.82	0.83	11.1
YOLOv10 [28]	0.85	0.73	14.3
YOLOv11	0.81	0.83	18.5
YOLOv12 (proposed)	0.88	0.86	20.2

Table 8. Detection Accuracy (%) for Each Category Across Tested Models.

Model	Text	Headline	Image	Advertisement	Caption	Average Accuracy (%)
Faster R-CNN [31]	86.42	87.30	90.15	84.67	85.92	86.89
Mask R-CNN [29]	88.73	88.54	91.23	86.10	87.48	88.42
YOLOv5	88.24	90.92	93.12	87.34	89.50	89.82
YOLOv8	89.35	91.87	92.21	87.42	90.76	90.32
YOLOv10 [28]	89.02	91.65	94.02	85.34	91.55	90.32
YOLOv11	90.02	91.56	94.45	86.11	92.02	90.83
YOLOv12 (Proposed)	91.49	92.02	96.92	90.78	92.97	92.84

measurable improvement in both precision and recall compared to its predecessors. Figure 9 shows the detection results on various newspaper images.

The findings clearly show that YOLOv12 performs better than all previous iterations in terms of accuracy (mAP@50 = 0.88) and generalization (mAP@50–90 = 0.86) without requiring an excessive amount of computing power. Notably, the suggested model shows better stability across different IoU thresholds and delivers a 5-7% increase in precision when compared to YOLOv10 and YOLOv11. This demonstrates YOLOv12’s improved flexibility to multilingual scripts, visually complex page structures, and varied newspaper layouts.

The YOLO series outperforms region-based techniques like Faster R-CNN and Mask R-CNN, as Table 8 shows the steady increase in detection accuracy over models when tested over same training set. With the greatest average accuracy (92.84%), the suggested YOLOv12 shows exceptional capacity to recognize both textual and visual layout elements. Due to the variety of visual structures, the Image category performs the best (96.92%), whereas the Advertisement category exhibits somewhat lower accuracy (90.78%). Overall,

the findings demonstrate YOLOv12’s improved accuracy and resilience for complicated newspaper layout recognition, suggesting that it is appropriate for tasks involving automated structural analysis and large-scale document digitalization.

7 Conclusion

This paper presents an advanced and systematic framework for newspaper layout recognition employing the improved YOLOv12 model. Through the development of a custom multilingual dataset and the integration of attention-based architectural enhancements, YOLOv12 demonstrates remarkable capability in accurately detecting and classifying heterogeneous newspaper components. The model achieved a mAP@50 of 0.88 and an average detection accuracy of 92.84%, outperforming established benchmarks such as Faster R-CNN, Mask RCNN, YOLOv10, and YOLOv11 in both precision and computational efficiency.

The results substantiate YOLOv12’s robustness and generalization across diverse linguistic and structural formats, effectively addressing challenges inherent in multilingual scripts, irregular columnar

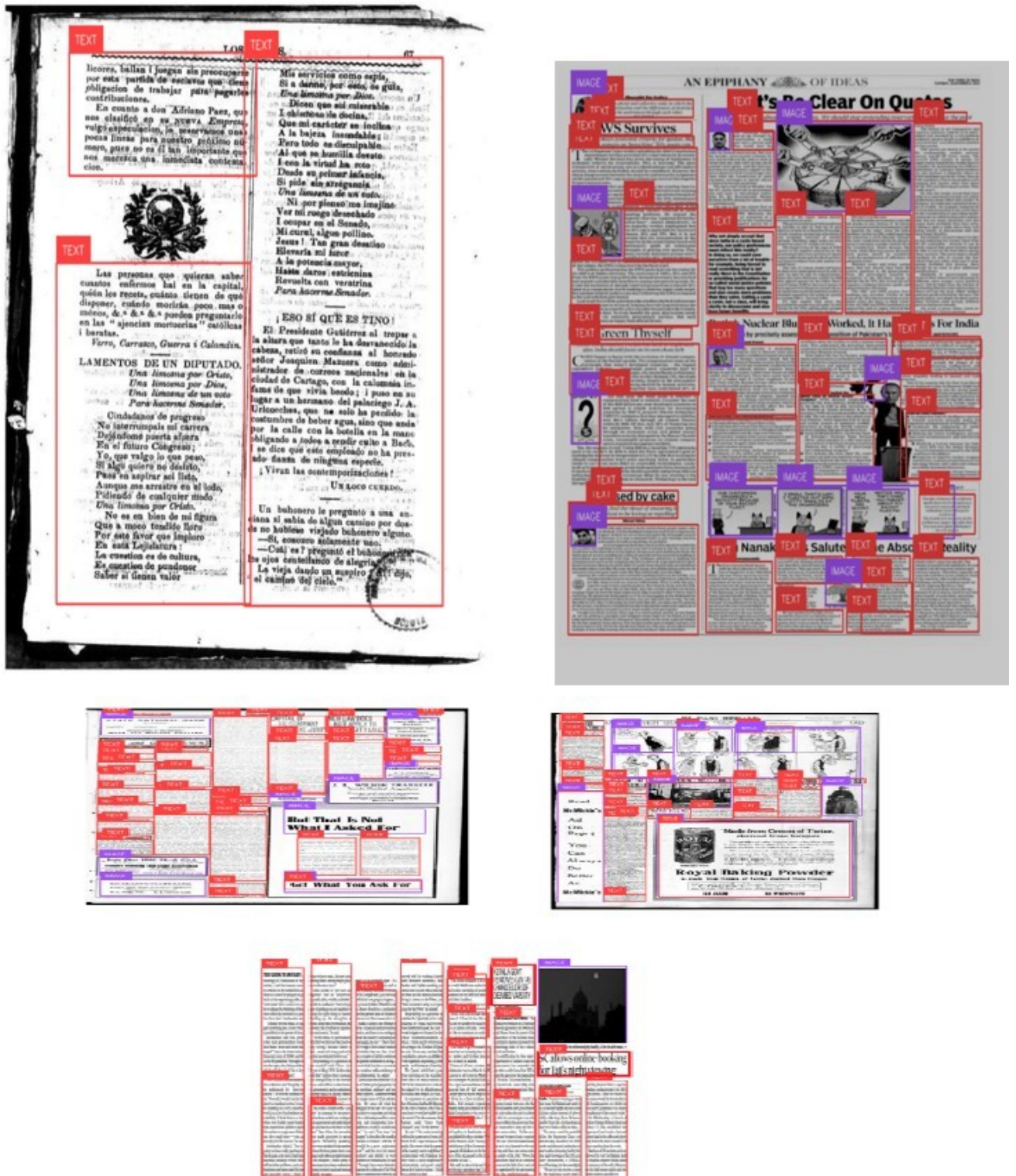


Figure 9. Layouts detection for different newspaper images.

layouts, and overlapping visual entities. Beyond achieving state-of-the-art detection accuracy, this study contributes a novel annotated dataset and a reproducible methodological framework for future research in document image analysis.

Consequently, the proposed approach provides a substantial advancement toward intelligent document understanding, automated archival preservation, and scalable newspaper digitization within digital humanities and media informatics domains.

Data Availability Statement

Data will be made available on request.

Funding

This work was supported without any funding.

Conflicts of Interest

The authors declare no conflicts of interest.

AI Use Statement

The authors declare that no generative AI was used in the preparation of this manuscript.

Ethical Approval and Consent to Participate

Not applicable.

References

- [1] Breuel, T. M. (2008, January). The OCRopus open source OCR system. In *Document recognition and retrieval XV* (Vol. 6815, pp. 120-134). SPIE. [CrossRef]
- [2] Namboodiri, A. M., & Jain, A. K. (2007). Document structure and layout analysis. In *Digital Document Processing: Major Directions and Recent Advances* (pp. 29-48). London: Springer London. [CrossRef]
- [3] Ultralytics. (2025). *Ultralytics YOLO12: Attention-centric object detection*. Retrieved June 24, 2025, from <https://github.com/ultralytics/yolo>
- [4] Kise, K. (2014). Page segmentation techniques in document analysis. In *Handbook of Document Image Processing and Recognition* (pp. 135-175). Springer, London.
- [5] Binmahashen, G. M., & Mahmoud, S. A. (2019). Document layout analysis: A comprehensive survey. *ACM Computing Surveys*, 52(6), 1-36. [CrossRef]
- [6] Sutheebanjard, P., & Premchaiswadi, W. (2010, April). A modified recursive xy cut algorithm for solving block ordering problems. In *2010 2nd International Conference on Computer Engineering and Technology* (Vol. 3, pp. V3-307). IEEE. [CrossRef]
- [7] Pavlidis, T., & Zhou, J. (1999). Page segmentation by white streams. In *Proceedings of the 1st International Conference on Document Analysis and Recognition (ICDAR)* (pp. 945-953).
- [8] Sun, H. M. (2006). Enhanced constrained run-length algorithm for complex layout document processing. *International Journal of Applied Science and Engineering*, 4(3), 297-309.
- [9] Gutehrle, N., & Atanassova, I. (2022). Processing the structure of documents: Logical layout analysis of historical newspapers in French. *Journal of Data Mining and Digital Humanities*. [CrossRef]
- [10] Zhu, W., Sokhandan, N., Yang, G., Martin, S., & Sathyanarayana, S. (2022, June). DocBed: A multi-stage OCR solution for documents with complex layouts. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 36, No. 11, pp. 12643-12649). [CrossRef]
- [11] Elanwar, R., Qin, W., Betke, M., & Wijaya, D. (2021). Extracting text from scanned Arabic books: a large-scale benchmark dataset and a fine-tuned Faster-R-CNN model. *International Journal on Document Analysis and Recognition (IJDAR)*, 24(4), 349-362. [CrossRef]
- [12] Shen, Z., Zhang, K., & Dell, M. (2020, June). A Large Dataset of Historical Japanese Documents with Complex Layouts. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (pp. 2336-2343). IEEE. [CrossRef]
- [13] Iskandar, N. A. A. (2023, June). Manga Layout Analysis via Deep Learning. In *IRC-SET 2022: Proceedings of the 8th IRC Conference on Science, Engineering and Technology, August 2022, Singapore* (pp. 63-73). Singapore: Springer Nature Singapore. [CrossRef]
- [14] Davis, B., Morse, B., Price, B., Tensmeyer, C., Wigington, C., & Morariu, V. (2022, October). End-to-end document recognition and understanding with dessurt. In *European Conference on Computer Vision* (pp. 280-296). Cham: Springer Nature Switzerland. [CrossRef]
- [15] Shanthakumari, A., Kalpana, R., Jayashankari, J., Umamaheswari, B., & Sirija, M. (2022, May). Mask RCNN and Tesseract OCR for vehicle plate character recognition. In *AIP Conference Proceedings* (Vol. 2393, No. 1, p. 020135). AIP Publishing LLC. [CrossRef]
- [16] Smith, R. (2007, September). An overview of the Tesseract OCR engine. In *Ninth international conference on document analysis and recognition (ICDAR 2007)* (Vol. 2, pp. 629-633). IEEE. [CrossRef]
- [17] Shen, Z., Zhang, R., Dell, M., Lee, B. C. G., Carlson, J., & Li, W. (2021, September). Layoutparser: A unified toolkit for deep learning based document image analysis. In *International Conference on Document Analysis and Recognition* (pp. 131-146). Cham: Springer International Publishing. [CrossRef]
- [18] Lin, T. Y., Dollár, P., Girshick, R., He, K., Hariharan, B., & Belongie, S. (2017, July). Feature Pyramid Networks for Object Detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 936-944). IEEE. [CrossRef]
- [19] Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 779-788). [CrossRef]
- [20] Kay, A. (2007). Tesseract: an open-source optical character recognition engine. *Linux Journal*, 159.

- [21] Jaha, E. S. (2023). Semantic document layout analysis of handwritten manuscripts. *Computers, Materials & Continua*, 75(2), 2805–2831. [CrossRef]
- [22] Barman, R., Ehrmann, M., Clematide, S., Oliveira, S. A., & Kaplan, F. (2021). Combining visual and textual features for semantic segmentation of historical newspapers. *Journal of Data Mining & Digital Humanities*, (HistoInformatics). [CrossRef]
- [23] Aljiffry, L., Al-Barhamtoshy, H., Jamal, A., & Abukhodair, F. (2022, October). Arabic Documents Layout Analysis (ADLA) using Fine-tuned Faster RCN. In *2022 20th International Conference on Language Engineering (ESOLEC)* (Vol. 20, pp. 66-71). IEEE. [CrossRef]
- [24] Singh, V., & Kumar, B. (2014, January). Document layout analysis for Indian newspapers using contour based symbiotic approach. In *2014 International Conference on Computer Communication and Informatics* (pp. 1-4). IEEE. [CrossRef]
- [25] Lombardi, F., & Marinai, S. (2020). Deep learning for historical document analysis and recognition—A survey. *Journal of Imaging*, 6(10), 110. [CrossRef]
- [26] Zhao, H., Min, W., Wang, Q., & Wei, Z. (2023). Memory-efficient document layout analysis method using LD-net. *Multimedia Tools and Applications*, 82(3), 4371–4386. [CrossRef]
- [27] Roboflow. (n.d.). *Give your software the power to see objects in images and video*. Retrieved June 24, 2025, from <https://roboflow.com/>
- [28] Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., & Han, J. (2024). Yolov10: Real-time end-to-end object detection. *Advances in Neural Information Processing Systems*, 37, 107984–108011.
- [29] He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2020). Mask R-CNN. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2), 386–397. [CrossRef]
- [30] Ultralytics. (2024). YOLOv11: The latest version of the YOLO series for object detection. Ultralytics Documentation. Retrieved June 24, 2025, from <https://docs.ultralytics.com/models/yolo11/>
- [31] Ren, S., He, K., Girshick, R., & Sun, J. (2017). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6), 1137–1149. [CrossRef]



Dr. Atul Kumar is an Assistant Professor in the Department of Computer Science at R.G.M. Government College, Joginder Nagar, District Mandi, India. He earned his Ph.D. from Punjabi University, Patiala, and has over 10 years of teaching experience. Dr. Kumar has authored 15 research papers in reputed journals and conferences. His research interests include Computer Science, Natural Language Processing (NLP), and Pattern Recognition, where his work has contributed significantly to advancing methods and applications in these domains. (Email: atulkmr02@gmail.com)



Dr. Gurpreet Singh Lehal is a Retired Professor from the Department of Computer Science, Punjabi University Patiala. Currently, he is working as a Senior Project Consultant at IIIT, Hyderabad. He has over three decades of contributions to computational linguistics, digital humanities, and Indic language technologies, particularly in Punjabi, Urdu and Hindi languages. His leadership spans major national and international projects in OCR, machine translation, transliteration, language modeling, and digital text processing, including large-scale government-funded initiatives. He has developed more than two dozen language software systems, supervised 23 PhD scholars, published over 200 indexed research papers, and served on influential committees such as ICANN panels and national standards bodies. His work has significantly advanced the technological ecosystem for Indian languages. (Email: gslehal@gmail.com)