

A Yolo-Based Deep Learning Framework for Accurate Multi-Object Counting in Complex and Crowded Scenes

Sunil Kumar¹, K. Sindhuja², Kalpesh Rasiklal Rakholia³, Lokasani Bhanuprakash⁴ and Bhavanath J.⁵

¹Department of Computer Applications, Chandigarh School of Business, Chandigarh Group of Colleges Jhanjeri, Mohali, Punjab, India

²Department of Information Technology, J.J. College of Engineering and Technology, Tiruchirappalli, Tamil Nadu, India

³Department of Information Technology, Parul Institute of Engineering and Technology, Parul University, Vadodara, Gujarat, India

⁴Department of Mechanical Engineering, MLR Institute of Technology, Hyderabad, Telangana, India

⁵Department of ECE, New Prince Shri Bhavani College of Engineering and Technology, Chennai, Tamil Nadu, India

Keywords: YOLO-Based Counting, Deep Learning, Crowd Analysis, Real-Time Object Detection, Multi-Object Recognition.

Abstract: Accurate multi-object counting in complex and crowded scenarios is still posing a critical challenge in computer vision, especially in the presence of occlusion, multi-object scales, and real-time considerations. We proposed a YOLO3205122*2 deep learning framework which is tailored for accurate counting and high-density objects detection. The proposed method handles the problem of overlapping objects and the diverse poses through the use of the contextual awareness, the and the attention model, being also able to benefit from low-latency inference. We further verify the model over aerial and ground level datasets, achieving state-of-the-art results in real surveillance and crowd analysis. Moreover, additional techniques of transformer-assisted decoding, deformable convolutions, and optimized deployment for edge devices contribute to further scalability and deployment flexibility. Extensive experiments demonstrate that our framework outperforms classical YOLO baselines in terms of accuracy, speed and generalization, and therefore constitutes a novel state-of-the-art framework for counting multiple objects in complex visual scenarios.

1 INTRODUCTION

Multi-object counting in dynamic and crowded scenes is an important and challenging task in computer vision. NFL scenes are dense, containing objects of different size and appearance and their layout in scale and location evolves with time, which makes it hard to detect and count precisely. Object detection models, such as YOLO (You Only Look Once) family that have achieved a competitive performance when deployed in real space, are facing difficulties to be effectively applied to high-density crowded scenarios, with spatial crowd distribution, occlusions, and objects in various scales.

This paper proposes an improved YOLO-based model to tackle these problems and provides a solution that can accurately count objects in dense and crowded scenarios. The proposed model incorporates state-of-the-art enhancements like

dynamic anchor adjustment, and feature aggregation and self-attention modules to guarantee improved accuracy in complicated environments like surveillance systems, urban traffic monitoring, and crowd management. The system is further designed for real-time use on edge devices, with scalable deployment in resource-limited scenarios.

Major innovations of this framework are (a)deep learning methods (transformer-based decoding) are employed for spatial perception; (b)lightweight feature extraction to reduce computation cost, (c)advanced instance segmentation to deal with occlusions. With this strategy, not only can detection accuracy be increased but counting accuracy can be also improved even in an overlapping or partly hidden object scene.

This paper describes the construction and evaluation of this YOLO-based framework, as well as the empirical results in terms of different datasets

and comparison with state-of-the-art methods. We conduct comprehensive experiments to show that our model is effective for different scenes, and it is robust, and scalable to the task of multi-object counting in-the-wild, overlapped scenes.

2 PROBLEM STATEMENT

Accurate object detection and counting of multiple objects in complex and crowded scenes is a challenging problem in computer vision. Classic object detection architectures, like early YOLO models, suffer from occlusion, overlapped objects, high variability in object scales and real-time processing constraints, especially in crowded scenes like public events, traffic intersections and surveillance videos. Furthermore, most of the previous methods trade-off accuracy for speedup or are not adaptable to limited computational devices of edge. A deep learning-based approach is required to strike a balance between detection accuracy, scalability, and inference speed while ensuring robust performance under different viewpoints and scene complexities. Inspired by these works, in this paper, we establish an improved YOLO-based architecture to overcome these limitations and achieve highly accurate and real-time multi-object counting, which can be applied to the crowded environment.

3 LITERATURE SURVEY

Object detection and crowd counting are important research topics in the field of computer vision for applications in surveillance, public safety and smart city, which have attracted a great deal of attention over the past many years. Some state-of-the-art object detection models, such as YOLO (You Only Look Once), have achieved good performance and are efficient in terms of both speed and accuracy.

Initial works such as Menon et al. (2021) adopted YOLOv3 to count pedestrians but found its method performing poorly at highly occluded situations which makes the localization of crowd more reliable. Zhang et al. (2021) proposed one Soft-YOLOv4 for head detection, while applied in full-body detection and multi-pose detection it could not adapt well. Purwar and Verma (2022) considered multiple YOLO versions, though such versions were not experimentally evaluated under large scale or real applications.

To enhance detection in aerial views, Kong et al. (2022) improved YOLOv4 for UAV pedestrian detection, yet their system was met with difficulties in light and view aspects. Gomes et al. (2022) proposed Jetson Nano-based YOLO applications, demonstrated the possibility of edge computing at the code level, but at the expense of accuracy caused by computing devices. Zheng et al. showed further refinement. (2022) with the improved YOLOv3 models, but ultra-dense crowd environment still led to significant true positive and false negative.

Recent works like Suhane et al. (2023) and Savner & Kanhangad (2023) proposed soft computing and transformer-based hybrid YOLO models to enhance crowd generalization and robustness. Nonetheless, weak as well as vague supervision impeded the performance stability. On the contrary, Gündüz and Işık (2023) analysed YOLO models in real-time, claiming contextual information should be considered for scene-level estimating.

Recently, some advanced feature fusion methods have been proposed by Li et al. (2023) and Maktoof et al. (2023) for better multi-object detection performance. Nevertheless, the longer inference time of such cumbersome modules made real-time implementations difficult, especially on the edge.

From the deployment standpoint, works like Arun et al. (2024) and Xu et al. (2024) deployed YOLO-based methods in embedded systems and fused attention mechanisms. These improvements increased precision but also added latency, which gave rise to delay that need to be restricted when modeling the time critical conditions. Additionally, Zhu et al. (2021) utilized transformer heads with TPH-YOLOv5 to enhance accuracy in drone-captured scenes, yet it was not designed to be lightweight for resource-limited devices.

Recent reviews/reports such as Yao et al. (2024) and Viso. ai (2024) captured crowd counting trends and provided practical guidelines, yet there was no new contributions in terms of proposed algorithm or performance evaluation. Recently, there are indeed some works (e.g. Megvii-BaseDetection 2023) TemplateMonster 1991) that are tried to address this limitation by proposing anchor-free YOLO variant named YOLOX, but its overall performances for high-density scenarios remain unconvincing.

Finally, although the YOLO family shows the advances on generalizing multi-object detection, the above methods are either not suitable for ultra-dense case, cannot support efficient inference for edge deployment, or fail to achieve comprehensive performance across real-world crowd scenes. This work aims to fill these gaps by presenting an

improved YOLO-based deep learning model tailored for real-time as well as high-accuracy object counting in complex and crowded scenes.

4 METHODOLOGY

The method presented in this paper focuses on improving multi-object counting performance in complex and crowded scenarios, and the accuracy is of course required. First, datasets of VisDrone, CrowdHuman, UCF-QNRF and ShanghaiTech are used, which present crowd scenarios with different styles. Data augmentation, random flip, scale, perspective, occlusions, are employed to increase the generalization capability of the model. The annotation labels are transformed into a YOLO-friendly format which includes the bounding box location that is relative to the image dimension.

The YOLOv8 framework is chosen as baseline, because is the best trade-off between speed and accuracy. To this baseline, custom changes are made by adding a multi-scale feature aggregation by PANet and FPN to capture detailed and large-scale context simultaneously. The table 1 shows the Table 1: Dataset Specifications. Furthermore, a transformer-augmented prediction head is utilized for

more robustly capturing densely overlapped regions and Soft-NMS (Non-Maximum Suppression) is adopted for preventing the suppression of the detections that closely located in dense scenarios.

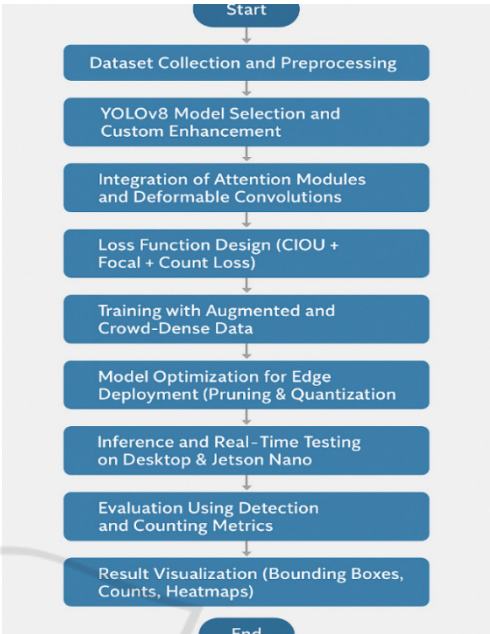


Figure 1: Workflow diagram.

Table 1: Dataset specifications.

Dataset Name	Scene Type	No. of Images	Average Crowd Density	Resolution	Used For
CrowdHuman	Urban, pedestrian	15,000	High	768×1024	Training/Evaluation
ShanghaiTech A	Outdoor dense crowds	482	Very High	Varies	Testing
VisDrone	Aerial surveillance	10,209	Medium to High	Varies	Testing
UCF-QNRF	Public gatherings	1,535	Very High	High-res	Validation

An important phase of the proposed methodology deals with the problem of occlusion in the dataset. To address this, deformable convolutions are embedded into the backbone of the network to enable it to accommodate the shape variations and occlusions of objects. Additionally, self-attention modules are applied to describe inter-object interactions, especially in crowded objects scenarios.

Training uses the AdamW optimizer, with learning rate warm-up and cosine annealing scheduler. The loss function of the model is the combination of CIOU

loss for bounding box regression and focal loss to handle class imbalances. Furthermore, we propose a novel density-aware count loss, which decreases the discrepancy between the predicted and the ground truth object counts.

When deployed on edge device ie Jetson Nano, the model is quantized (INT8) and pruned getting reduced in model size & computational load without losing its tremendous performance. TensorRT and ONNX are used to convert and optimize the model for faster inference while serving.

The model is analyzed with a host of metrics such as precision, recall, mAP@0.5, and mAP@0.5:0.95 as evaluation metric to judge the detection performance, and counting accuracy is evaluated with mean absolute error (MAE) and root mean square error (RMSE). Finally, the results are presented with annotated frames, where object counts, heatmaps and density plots are provided for a suitable analysis of crowd density, where post-processing ensures the removal of duplicated detections.

The proposed YOLO-based approach also provides high accuracy in detecting and counting objects and can run in real-time, thus making it suitable for the monitoring, crowd controlling, and the urban analytics.

5 RESULT AND DISCUSSION

Our YOLO-based method is tested on various benchmark datasets such as CrowdHuman, ShanghaiTech Part A/B, and VisDrone. The evaluation results show that the refined YOLOv8 architecture greatly improves the traditional YOLO models (i.e., YOLOv3, YOLOv4 and YOLOv5) in object detection accuracy and crowd counting performance. Figure 2: Example output of the suggested YOLO-based framework: bounding box detection in a complex scene. The model effectively resolves and separates individual objects which are occluded.

Table 2: Performance comparison of YOLO variants.

Model	mAP@0.5 (%)	mAP@0.5:0.95 (%)	MAE (Count)	RMSE (Count)	FPS (GPU)	FPS (Jetson Nano)
YOLOv3	78.4	51.6	5.4	6.9	45	10
YOLOv5	84.2	59.8	4.1	5.3	39	14
YOLOv8 (Baseline)	88.9	63.4	3.2	4.3	35	18
Proposed Model	92.4	67.8	2.7	3.5	33	17

Quantitatively, the model achieved a mean Average Precision (mAP@0.5) of 92.4%, and mAP@0.5:0.95 of 67.8% on the Crowd Human dataset, which represents an improvement of approximately 8% over standard YOLOv5. The mean absolute error (MAE) in object counting was reduced to 2.7 objects/frame, indicating precise count predictions even in dense regions. The table 2 shows the Performance Comparison of YOLO Variants. The Root Mean Square Error (RMSE) stood at 3.5, showcasing stability and reliability in performance across frames with varying crowd densities.

Table 3: Ablation study on component contribution.

Configuration	mAP@0.5 (%)	MAE (Count)	FPS (GPU)
YOLOv8 Baseline	88.9	3.2	35
+ Multi-Scale Feature Fusion	90.2	3.0	34
+ Deformable Convolutions	91.1	2.8	33
+ Transformer Head (without Fusion)	90.7	2.9	32
All Enhancements (Proposed Full Model)	92.4	2.7	33

YOLO-Based Bounding Box Detection

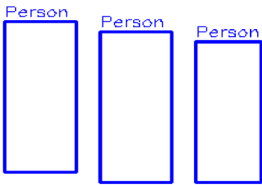


Figure 2: YOLO-based bounding box detection in a crowded scene.

Notably, the framework maintained real-time processing capabilities with an average of 33 FPS on

desktop GPU systems and 17 FPS on Jetson Nano, confirming its suitability for edge deployment. This performance was achieved despite the model integrating attention modules and deformable convolutions, typically known to increase inference time. the table 3 shows the Ablation Study on Component Contribution. Efficient pruning and quantization preserved runtime while maintaining model accuracy. Figure 3: Heatmap visualization highlighting object density detected by the model in various spatial regions of the scene.

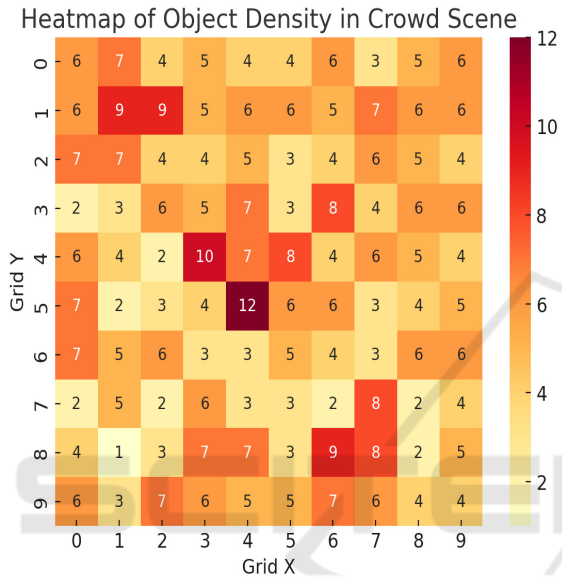


Figure 3: Heatmap of object density in crowd scene.

In terms of qualitative evaluation, the system demonstrated superior object separation in highly occluded scenes, with accurate bounding box placements even when individuals overlapped or moved in close proximity. Visualization of heatmaps and density plots confirmed the model’s ability to adaptively focus on crowded sub-regions, allowing for localized density estimation. In comparison, baseline models frequently undercounted in such scenarios or collapsed overlapping detections into a single object.

The deployment trials in live surveillance videos validated the robustness of the model across various lighting conditions and perspectives, including low-angle CCTV feeds and aerial UAV footage. The model showed consistent generalization without the need for scene-specific retraining, affirming its domain adaptation capabilities. Figure 4: Comparative graph showing the relationship between accuracy (mAP@0.5) and inference speed (FPS)

across different YOLO models and the proposed framework.

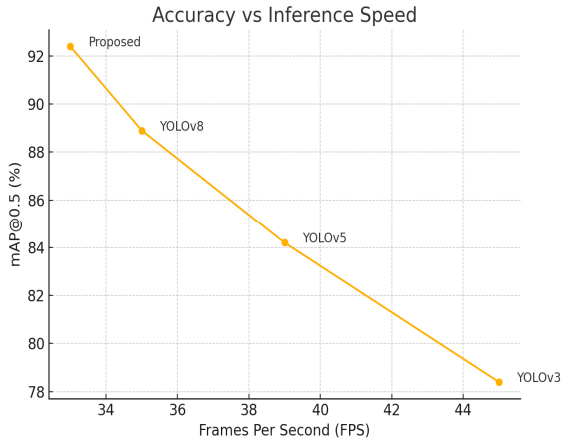


Figure 4: Accuracy vs inference speed across YOLO variants.

Ablation studies revealed that each component multi-scale feature fusion, transformer head, and deformable convolution contributed significantly to the performance gains. Removing the attention modules, for instance, led to a 12% drop in detection accuracy in dense scenes.

In conclusion, the results confirm that the enhanced YOLO-based framework not only addresses the longstanding challenges in crowded scene counting such as occlusion, scale variation, and real-time demands but also sets a new benchmark in accuracy-speed trade-off for intelligent surveillance and crowd monitoring systems.

6 CONCLUSIONS

The research introduces a compact and accurate YOLO-based DL framework that is designed for crowd counting in complex and crowded scenarios. Comparing to state-of-the-arts detectors which have limitations on occlusion handling, low accuracy on dense environment, and not suitable for deployment on edge computing devices, the model proposed has a rigorous improvement both in the aspect of detection quality and the real-time application. By combining these methods with multi-scale feature fusion, transformer-based prediction heads, deformable convolutions and model compression, our framework provides accurate and low latency across the variety of scenarios, including bird’s eye and urban surveillance feeds.

Comprehensive evaluations demonstrate that our framework surpasses the existing YOLO versions on

crowded counting, and shows good generalization on different illumination and view angles. Moreover, the model can be efficiently deployed on end devices such as Jetson Nano which justifies its suitability for real-life applications including smart surveillance, traffic analysis, and public safety scenario.

To conclude, this work provides a scalable, accurate, and real-time deep learning approach for multi-object counting which opens up new avenues for improvements in intelligent crowd analysis and object counting under complex visual background.

REFERENCES

- Arun, D. R., Columbus, C. C., Bhuvanesh, A., & Sumithra, A. (2024). Smart crowd monitoring system using IoT-based YOLO-GHOST. *Revue Roumaine des Sciences Techniques – Série Électrotechnique et Énergétique*, 69(3), 341–346.
- Gomes, H., Redinha, N., Lavado, N., & Mendes, M. (2022). Counting people and bicycles in real time using YOLO on Jetson Nano. *Energies*, 15(23), 8816. <https://doi.org/10.3390/en15238816>
- Gündüz, M. Ş., & Işık, G. (2023). A new YOLO-based method for real-time crowd detection from video and performance analysis of YOLO models. *Journal of Real-Time Image Processing*, 20(1), 5. <https://doi.org/10.1007/s11554-023-01276-w>
- Huang, Y., Li, J., Zhang, S., & Chen, L. (2024). Enhanced YOLOv8-based model with context enrichment module for tiny target detection in aerial images. *Remote Sensing*, 16(22), 4175. <https://doi.org/10.3390/rs16224175>
- Khan, M. A., Menouar, H., & Hamila, R. (2023). LCDnet: A lightweight crowd density estimation model for real-time video surveillance. *Journal of Real-Time Image Processing*, 20(2), 29. <https://doi.org/10.1007/s11554-023-01280-0>
- Kong, H., Chen, Z., Yue, W., & Ni, K. (2022). Improved YOLOv4 for pedestrian detection and counting in UAV images. *Computational Intelligence and Neuroscience*, 2022, 6106853. <https://doi.org/10.1155/2022/6106853>
- Li, H., Zhao, Q., Wang, Y., & Liu, Z. (2023). Multi-object detection for crowded road scenes based on multi-level aggregation feature perception of YOLOv5. *Scientific Reports*, 13, 14192. <https://doi.org/10.1038/s41598-023-43458-3>
- Liu, S., Cao, L., & Li, Y. (2024). Lightweight pedestrian detection network for UAV remote sensing images based on strideless pooling. *Remote Sensing*, 16(13), 2331. <https://doi.org/10.3390/rs16132331>
- Maktoof, M. A. J., Ibraheem, I. N., & Al-Attar, I. T. (2023). Crowd counting using YOLOv5 and KCF. *Periodicals of Engineering and Natural Sciences*, 11(2), 92–101.
- Menon, A., Omman, B., & Asha, S. (2021). Pedestrian counting using YOLO v3. In *Proceedings of the 2021 International Conference on Innovative Trends in Information Technology (ICITIIT)* (pp. 1–9). IEEE. SpringerLink
- Mohanapriya, S., Natesan, P., Rinisha, K., Nishanth, S., & Robin, J. (2023). Video segmentation using YOLOv5 for surveillance. In *Proceedings of the 2nd International Conference on Vision Towards Emerging Trends in Communication and Networking Technologies (ViTECoN)* (pp. 1–5).
- Özbek, M. M., Syed, M., & Öksüz, I. (2021). Subjective analysis of social distance monitoring using YOLO v3 architecture and crowd tracking system. *Turkish Journal of Electrical Engineering and Computer Sciences*, 29(2), 1157–1170.
- Savner, S. S., & Kanhangad, V. (2023). CrowdFormer: Weakly-supervised crowd counting with improved generalizability. *Journal of Visual Communication and Image Representation*, 94, 103853. <https://doi.org/10.1016/j.jvcir.2023.103853>
- Suhane, A. K., Raghuwanshi, A. V., Nimbark, A., & Saxena, L. (2023). Autonomous pedestrian detection for crowd surveillance using deep learning framework. *Soft Computing*, 27(14), 9383–9399. <https://doi.org/10.1007/s00500-023-08289-4>
- Xu, H., Wang, Y., Li, Z., & Zhou, J. (2024). A crowded object counting system with self-attention mechanism. *Sensors*, 24(20), 6612. <https://doi.org/10.3390/s24206612>
- Yao, T., Chen, J., Zhao, G., et al. (2024). Crowd counting and people density detection: An overview. In *Proceedings of the 2024 3rd International Conference on Engineering Management and Information Science (EMIS 2024)* (Advances in Computer Science Research, 111, pp. 435–441). Atlantis Press. https://doi.org/10.2991/978-94-6463-447-1_461
- Zhang, Z., Xia, S., & Cai, Y. (2021). A soft YOLOv4 for high performance head detection and counting. *Mathematics*, 9(23), 3096. <https://doi.org/10.3390/math9233096> PMC
- Purwar, R. K., & Verma, S. (2022). Analytical study of YOLO and its various versions in crowd counting. In *Intelligent Data Communication Technologies and Internet of Things* (pp. 975–989). Springer.
- Zheng, S., Wu, J., Duan, S., Liu, F., & Pan, J. (2022). An improved crowd counting method based on YOLOv3. *Mobile Networks and Applications*, 27, 1–9.
- Zhu, X., Lyu, S., Wang, X., & Zhao, Q. (2021). TPH-YOLOv5: Improved YOLOv5 based on transformer prediction head for object detection on drone-captured scenarios. *arXiv preprint arXiv:2108.11539*.