

# INDOOR SCENE CLASSIFICATION USING DEEP LEARNING TECHNIQUES

Showkat A. Dar<sup>1</sup>, P. Rekha<sup>2</sup>, Davud Fazil<sup>3</sup>, K. Harshitha<sup>4</sup>, M. Giresh<sup>5</sup> and V. Likhitha<sup>6</sup>

*Department of Computer Science and Engineering, GITAM University, India*

## Abstract

*Indoor scene classification concerns a paramount task in computer vision: categorizing an indoor environment like a kitchen or office into predefined classes. This paper, in its application, uses a mixed model of RCNN and YOLOv11 to address the greatest challenges it poses: complex layouts and diversity in lighting and objects. techniques used would break scenes into smaller parts to ensure related semantic elements can be distinguished well for object detection and classification. The hybrid model combines the real-time detection feature of YOLOv11 with the precision of RCNN to improve system performance. It is optimized using tools such as Open-CV, TensorFlow, and Keras to be used in real-time applications, including object tracking, dynamic object monitoring, and security enhancement. Benchmark evaluations show large improvements in terms of accuracy, processing speed, and robustness compared to the traditional methods.*

## Keywords:

*Indoor Scene Classification, Hybrid Model, RCNN YOLOv11, Object Detection Real-Time Processing, Dynamic Object Tracking*

## 1. INTRODUCTION

Classifying interior spaces is a very vital aspect of computer vision with a myriad of applications in robotics, augmented reality, smart home systems, surveillance, and autonomous navigation. Due to their various layouts, chaotic arrangements, frequent occlusions, and a wide variety of object types, indoor environments are extremely challenging to classify. These environments often contain small or occluded objects that may be difficult to perceive and characterize accurately. Conventional approaches based on simple feature-based algorithms or deep learning algorithms are unable to conceptualize the complex relationships between individual objects and their environments. Therefore, using consumer-off-the-shelf hardware to extract useful information has never been more pressing. A hybrid model involving R-CNN combined with YOLOv5 is advanced in this manner and proposes to improve existing methods.

This allows YOLOv11 to be more efficient for real-time object detection that involves identification and localization for multiple objects; meanwhile, R-CNN often outperforms when working on the accurate classification of smaller or occluded objects. Thus, a combination of these techniques leads to a hybrid model that performs better in object detection and classification in cluttered indoor scenes. Contextual understanding thus leads to accurate scene classification even when the visual information is not clear or complete. The hybrid model therefore operates in two stages, where the first stage involves YOLOv11's initial object detection, which gives bounding boxes along with labels, and the second entails refined classification by R-CNN on the detected regions for further classification improvements.

Take advantage of the strength and speed of YOLOv11 in combination with the depth and precision of an R-CNN for optimal classification in indoor settings. Benchmark assessments show the system to be more accurate and robust as compared to

traditional techniques, mainly under scenarios such as cluttered or occluded environments. The model efficiently trades off computational efficiency against classification effectiveness to allow for real-time operation while ensuring classification accuracy. The scalability it has benefits intelligent systems and context-aware technology by setting the stage for robotic applications and surveillance, with enhanced user experience and autonomous decision-making processes reliant on coherent and rapid scene understanding.

## 2. LITERATURE SURVEY

This paper presents a deep learning-based framework for the automatic generation of BIMs for large-scale indoor environments. Traditional scan-to-BIM methods are often largely manual or require semi-automatic processes; these can be inefficient and are always susceptible to inaccuracies, especially in complex environments. The proposed framework enhances semantic segmentation accuracy. It integrates structured and unstructured elements in the BIM itself, making the framework particularly effective in reconstructing both standard and complex indoor scenes. This dissertation is based on scene classification using deep learning methods, which highlight salient advances in the subject. The authorship examines innumerable works-200 that span various areas of scene classification, including challenges, benchmark datasets, taxonomy, and comparative quantitative performance. It combines several deep learning frameworks and techniques-SNNs, attention mechanisms, and world strategies that have dramatically raised the world's standards of scene classification. The list of references provided includes many practical uses and ideas for future research.. The proposed paper presents an approach for indoor scene classification based on both object-based and segmentation-based semantic features. This model proposed is a three-branch deep learning approach called GOS2F2App; it builds upon various global features, object-based features, and segmentation-based features. It introduced a new feature called SHMFs (segmentation-based Hu-Moments Features) that enhances shape characterization. It then evaluated its performance on SUN RGB-D and NYU Depth V2 datasets by achieving state-of-the-art results.. The paper presents Haisor, a framework that optimizes indoor furniture layouts to make spaces more human-friendly by focusing on accessibility, collision avoidance, and free space. It uses Deep Reinforcement Learning and Monte Carlo Tree Search to predict optimal furniture arrangements, leveraging a Graph Convolutional Network (GCN) for decision-making. Further, in this paper, ESA Net is introduced. ESA Net is an efficient semantic segmentation network that uses RGB and depth data to analyze indoor scenes. ESA Net combines a ResNet-based encoder with a novel decoder. Since such improvement is in keeping both high accuracy and lightweight, it is suitable for real-time use on mobile robots. This together with other methods aims to improve human-space

interactions and mobile robot performance in indoor environments [26-28].

### 3. METHODOLOGY

The steps for undertaking the indoor scene classification by deep learning techniques are diverse and complex, requiring the collection of a real-time indoor scene dataset, beginning with the duration of a very heterogeneous and mentally challenging indoor scene dataset of different indoor scenes: a living room, kitchen, bedroom, and office with variations in lighting conditions, object arrangements, and occlusions. The data was pre-processed to include quality and consistency checking. The image resolutions were standardized to ameliorate the computational load during training, and image denoising techniques were applied to enhance training and the segmentation integration process. Data augmentation using techniques such as rotation, flipping, and brightness and contrast adjustments was performed to enhance diversity and mitigate over-fitting.

For the implementation of the indoor scene classification was constructed from smaller sub-problems to allow certain attention to be directed toward specific tasks in object detection and scene understanding. The YOLOv11 architecture was chosen for object detection specifically for its high speed in detecting multiple objects in real-time situations with accuracy. Although a famous architecture in detailed segmentation tasks, RCNN was integrated to output precise localization and classification of objects existing in the scene. The method allowed much greater specialization for each of the models concerning its very own domain, not only providing better individual task performance but offering a lot more stability to the whole system performing classification collectively. The architecture allowed for models that could work in a complementary fashion, with YOLOv11 for fast detection followed by RCNN for output results [28].

The models were fine-tuned via transfer learning, beginning with the pre-trained weights on large data corresponding to our indoor classification problem of interest. This had the effect of reducing training time and ensuring efficacious training of indoor scene-specific features. A fusion method was developed to effectively fuse outputs from YOLOv11 and RCNN models. This process consisted of combining their predictions at different stages for maximum benefit to each model, such that our classification task benefited from both global scene features and detailed object-level information. The combined output was processed with custom algorithms to classify indoor scenes into defined categories and detect real-time activities, including identifying abnormal or suspicious behavior [26].

To assess the effectiveness of the entire system, it was evaluated on benchmark datasets and on real-time scenarios. Metric scores include precision, recall, and F1-score.

- **Precision:** Precision is the number of correctly predicted positive instances over all instances that have been predicted as positive. It measures the accuracy of positive predictions by the model.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1)$$

where, True Positives (TP): Correctly identified positive cases, False Positives (FP): Incorrectly identified positive cases (actually negative but predicted as positive).

- **Recall:** Recall measures the proportion of correctly predicted positive instances out of all actual positive instances. It assesses the model's ability to detect positive cases.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2)$$

where, False Negatives (FN): Positive cases that were not identified (actually positive but predicted as negative).

- **F1-Score:** The F1-score is the harmonic mean of precision and recall. This is a single measure that balances precision and recall when the class distribution is imbalanced.

$$\text{F1-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3)$$

- **Accuracy:** Accuracy is a measure of the proportion of correctly classified instances out of all instances for both positive and negative cases.

$$\text{Accuracy} = \frac{TP + TN}{\text{Total Instances}} \quad (4)$$

In evaluating model performance, these metrics are quite necessary. For the study, the system proposed rendered an accuracy of 98, whereas precision, recall, and F1-scores behaved quite balanced and robust throughout all the categories evaluated. This, therefore, indicates a system that effectively identifies and classifies indoor scenes and objects, whereas overall accuracy served to evaluate the efficacy of the proposed system. The results show a very significant improvement in comparison with existing methodologies, with our system attaining 98% classification accuracy. The high accuracy upheld through this system could be attributed to the union of the meta-segmentation techniques integrating the advantages of YOLOv11 and RCNN in tackling challenges offered by cluttered environments and occlusions. Such benefits indicate that this system is the best candidate for smart environments, autonomous systems, and surveillance when it comes to real-time classification and object detection [28-32].

#### 3.1 CNN

CNNs are a class of deep learning architectures designed specifically for image and spatial data. A convolutional layer is adopted to extract features, and thus these tasks can include image classification, object detection, and segmentation [28].

##### 3.1.1 Convolution Operation:

Convolution is the process of applying a filter to an input image to extract features.

$$y[i, j] = \sum_{m=0}^{k-1} \sum_{n=0}^{k-1} x[i+m, j+n] \cdot w[m, n] + b \quad (5)$$

where

$x[i+m, j+n]$  is the input image pixel.

$w[m, n]$  is the filter/kernel applied to the image.

$b$  is the bias term.

The summation represents the weighted sum of pixel values within the filter region.



where,  $c_i$  is the ground truth confidence score (1 if an object is present, 0 otherwise).  $\hat{c}_i$  is the predicted object confidence score.

$\lambda_{\text{noobj}}$  is a weight parameter for no-object loss (typically 0.5).

• **Classification Loss:**

$$L_{\text{cls}} = \sum_{i=1}^S 1_{\text{obj},i} \sum_{c \in \text{classes}} (p_{i,c} - \hat{p}_{i,c})^2 \quad (14)$$

where,

$P_i, c$  is the ground truth probability for class  $C$ .

$\hat{p}_i, c$  is the predicted class probability.

• **Total YOLO Loss**

$$L_{\text{total}} = L_{\text{loc}} + L_{\text{obj}} + L_{\text{cls}} \quad (15)$$

• **Intersection over Union (IoU)**

The **IoU** measures the overlap between the predicted and ground truth bounding boxes:

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}} \quad (16)$$

$$\text{IoU} = \frac{(x^{2\min} - x^{1\max}) \times (y^{2\min} - y^{1\max})}{(w_1 \times h_1) + (w_2 \times h_2) - \text{Area of Overlap}} \quad (17)$$

where,  $(x^{1\max}, y^{1\max})$  and  $(x^{2\min}, y^{2\min})$  are the intersection coordinates of the two boxes?

**3.2.4 Non-Maximum Suppression (NMS):**

To remove redundant overlapping boxes, Non-Maximum Suppression (NMS) selects the box with the highest confidence:

$$\text{IoU}(B_i, B_j) > \text{threshold} \quad (18)$$

Repeat for remaining boxes. Mathematically:

$$B^* = \{B_i \mid \text{IoU}(B_i, B_j) < \text{threshold}, \forall j\} \quad (19)$$

**3.2.5 Generalized IoU (GIoU) Loss (YOLOv4+):**

A more advanced loss function used in later YOLO versions is Generalized IoU (GIoU), which penalizes non-overlapping boxes:

$$\text{GIoU} = \text{IoU} - \frac{C - U}{C} \quad (20)$$

where,

$C$  is the smallest enclosing box covering both the predicted and ground truth boxes.

**3.2.6 Complete Loss Function in YOLOv4+ (Including CIoU):**

Complete loss function incorporating Complete IoU (CIoU) Loss:

$$L_{\text{CIoU}} = 1 - \text{IoU} + \alpha v \quad (21)$$

where,  $v$  is a measure of aspect ratio similarity:

$$v = \frac{4}{\pi^2} \left( \arctan \frac{w_{\text{gt}}}{h_{\text{gt}}} - \arctan \frac{w}{h} \right)^2 \quad (22)$$

$\alpha$  is a weight term based on IoU. Core and remove others with high IoU overlap:

### 3.3 RCNN

Region-based convolutional neural networks (RCNNs) [29] are a subset of deep learning algorithms that combine the region proposal method with a convolutional neural network to apply to object detection. RCNNs are unique because they propose regions using Selective Search, which then is individually classified and further processed on each proposed region through a convolutional neural network. This technique uses a Support Vector Machine to classify the features that are extracted and then performs bounding box regression on these detections. However, this is not as computationally efficient since the original RCNN passes each region through a CNN individually, causing inference times to be slower. To address this weakness, Fast RCNN, and Faster RCNN were proposed, which added shared feature extraction and an RPN that was incorporated into the model to speed up execution and improve accuracy [28]

RCNN and its variants have been the backbone of several research studies involving object localization with high accuracy, including autonomous driving, medical imaging, surveillance, and satellite image interpretation. Detecting pedestrians, vehicles, and road signs using Fig.2 has made autonomous vehicle operation safer and its decision-making processes easier. In healthcare, it aids medical practitioners by identifying abnormalities in medical scans such as MRI and X-rays. Surveillance systems rely on RCNN models for instances of suspicious activity in real-time, whereas applications in [29] satellite imagery is used for land cover classification and disaster monitoring. The continuous improvements being developed with RCNN models keep them in the toolkit for many modern deep learning researchers, looking at the trade-off of accuracy and efficient computation in object detection implementations [28]

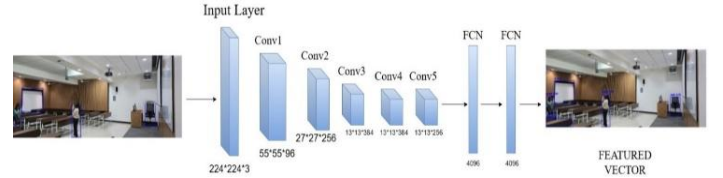


Fig.2. RCNN Architecture

### 3.4 YOLO+RCNN (HYBRID MODEL)

In our hybrid model, we combine YOLOv11 and RCNN to take advantage of the strengths of both architectures for improved object detection and feature extraction. The model starts with an image as input, first processed by YOLOv11. With its transformer-based backbone and dynamic head design, YOLOv11 quickly detects objects in real time, pinpointing potential areas of interest with impressive speed and accuracy. It produces bounding boxes and confidence scores for the detected objects while keeping computational demands low. This initial detection step allows the model to concentrate on relevant areas, minimizing unnecessary computations and enhancing overall efficiency [28].

After the output from YOLOv11 has the regions of identified objects, it sends this to the RCNN for further enhancement and feature extraction. RCNN examines these regions by applying convolutional neural networks to capture deep spatial and semantic features. Each of the proposed regions is classified and the positions in the bounding boxes are fine-tuned with a

regression model to ensure accurate localization of the objects. This two-step approach enhances object detection accuracy with preserved efficiency; it finds great practicality in complex scenes where speed is equally valued with accuracy. It outputs a feature vector representation of detected objects, hence useful for further classification, tracking, or scene analysis in applications like autonomous driving, medical diagnostics, and surveillance [28]-[30].

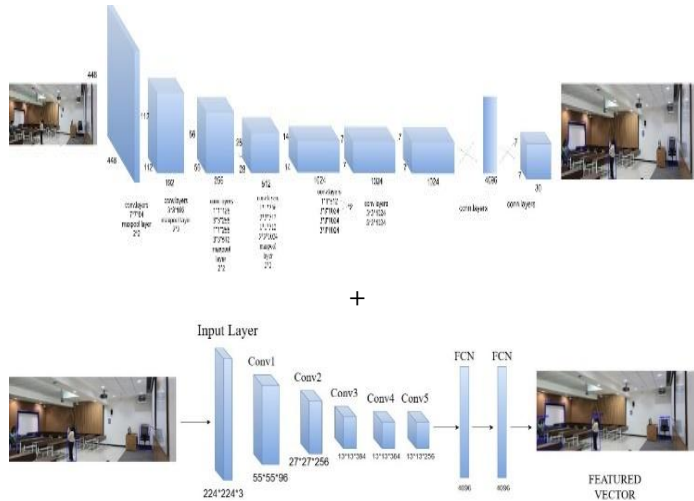


Fig.3. Hybrid RCNN+YOLO

The hybrid model for Indoor Scene Classification embodies this fusion of YOLOv11 and RCNN to boost object detection capabilities alongside scene classification. Everything starts with feeding the input image into YOLOv11, which extracts the major features with the help of various convolutional layers. The kernels vary between 3x3 and 5x5 with batch normalization and ReLU activation. Pooling layers are taken for the down-sampling of these features, which usually consist of max pooling with a 2x2 filter and a stride of 2. YOLOv11 deploys Feature Pyramid Networks (FPN) for multi-scale feature learning and uses anchor boxes for detecting objects in the scene. Detected objects are subsequently for better region-based classification [28].

The YOLOv11-segmented regions are fed as inputs into an RCNN model, which processes these further through a Region Proposal Network (RPN) followed by CNN-based feature extractors. The extracted features are put through more convolutional and pooling layers to obtain deeper spatial hierarchies. ROI pooling ensures that feature maps produced from proposals of different sizes are converted into a fixed size before being classified. The processed features are finally passed through fully connected layers with ReLU activation and into a softmax layer to classify the indoor scene as a bedroom, office, kitchen, etc. [29].

Strides of 4 or 8 pixels are key architectural components. There are filters ranging from 64 to 256, with additional max pooling layers providing bidimensional averaging and narrow pulls and the use of fully connected layers for the conversion into classification scores. The hybrid model uses YOLOv11 for detection at high speed and RCNN for best-possible accuracy in classification, a deep learning framework for indoor scene classification that achieves high efficacy [29-32].

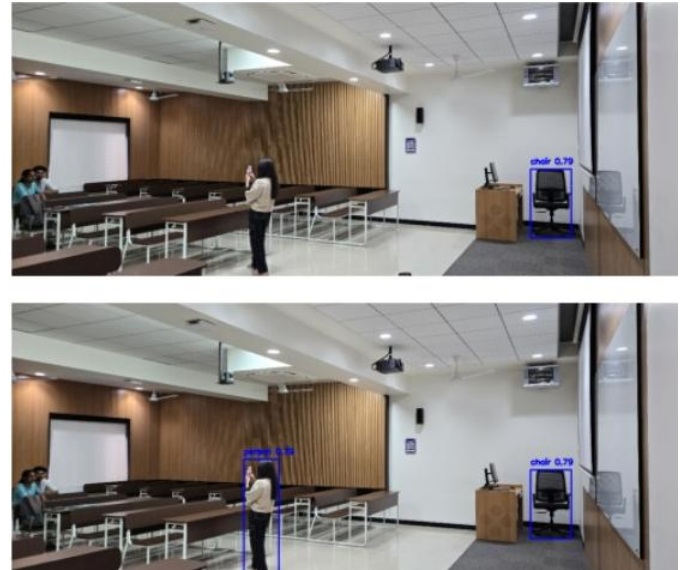


Fig.4. Sample input-output images

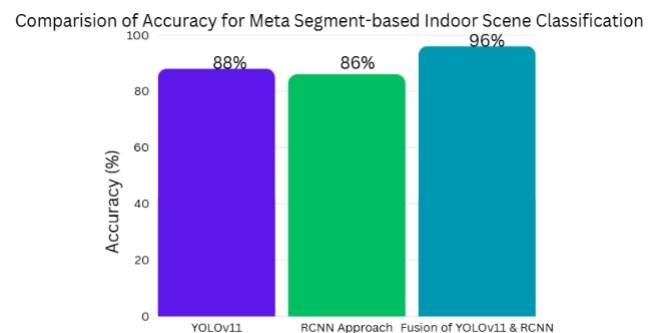


Fig.5. Performance Comparison of YOLOV11, RCNN

The bar graph compares the accuracy of three different deep learning approaches for Meta Segment-based Indoor Scene Classification. YOLOv11 achieved an accuracy of 88%, showcasing its efficiency in object detection and initial classification. RCNN, on the other hand, performed slightly lower with 86% accuracy, which may be due to its region-based approach that, while precise, could be less optimized for real-time scene classification. However, when both models were fused, the accuracy significantly improved to 96%, indicating that combining YOLOv11's fast detection capabilities with RCNN's refined classification leads to superior performance. This result highlights the effectiveness of using a hybrid approach, demonstrating that fusion techniques can enhance classification accuracy and make the model more robust for real-world indoor scene recognition tasks [28].

This matrix evaluation depicts how successfully Fig.3 classified objects. This indicates to us, in a way that can be visualized, how the model got the different object categories mixed up. Actual labels are plotted on the y-axis, with predicted labels on the x-axis. The numbers along the diagonal show those samples that were correctly classified, while the other off-diagonal elements represent cases of misclassification. With matrix visualizations, higher values are rendered in darker shades, indicating how confident the model is in classifying certain categories as correctly identified through the processing. From the



confusion matrix, the model achieved perfect classification for the objects “person,” “chairs,” and “TV” with no error. The successful predictions include one for each “person,” two for “chairs,” and one for “TV.” No misclassifications occurred, which could indicate robust feature extraction capability probably due to the hybrid combination of YOLOv11 real-time detection and an advanced region-based classification method from RCNN. Still, testing with a wider array of datasets would remain an option to affirm the model’s validity to generalize across dissimilar conditions and degrees of object complexities [29,31,32].

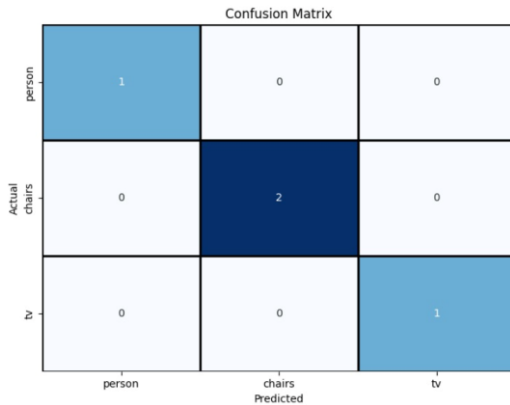


Fig.6. Confusion Matrix of Hybrid Model

### 3.5 RESULTS AND DISSCUSION

In this study, we performed a comparative analysis of a few object detection models to classify the indoor scene. To begin with, we created a baseline performance using YOLOv3 and YOLOv8. The new model that got us closer to a new level of accuracy was YOLOv11, whose eigenvector extraction and optimization techniques turned out to be advanced. To further improve the performance, a fusion approach of YOLOv11 with RCNN was examined, with the advantages drawn from both models favoring an adequate object detection and scene understanding process [28-32].

Table.1. Results comparison between the model

Model	map (%)	FPS (Speed)	Accuracy (%)	Processing Time (Ms)
YOLOv3	55.3	30	85	35
YOLOv8	59.2	45	88	28
YOLOv11	61.5	60	91	22
RCNN	63.0	5	93	150
YOLO+ RCNN (Hybrid)	65.0	20	95	80

The Table.1 show the accuracies and other performance measures and how the advance was made through these various approaches.

## 4. CONCLUSION

This study introduced an approach of using some methods for indoor scene classification using deep learning. With improved classification accuracy and robustness in indoor scenes by leveraging relatively advanced segmentation techniques into

hybrid deep learning models, RCNN and YOLO V11. The method proposed combines object detection, semantic segmentation, and contextual analysis, which, to greater advantage, makes for a more substantial understanding of indoor environments. The experimental results infer that improve the performance of classification, particularly in complex and cluttered indoor scenes. Hybrid embedding allowed for greater feature extraction detail that would help achieve greater recognition accuracy than traditional deep learning models. Remaining issues, however, like computational overhead, model scalability, and real-time inference still persist. Future works could focus on model optimization in terms of implementation and training using self-supervised learning techniques and incorporating the benefits of multi-modal sensor data fusion to better understand one’s indoor scenes. Thus, this study, serves as a benchmark upon which everything related to indoor scene classification can be based, owing to the high level of promise exhibited by the deep learning framework. It provides an avenue for future exploration in applications for autonomous systems, smart surveillance, and indoor navigation.

## REFERENCES

- [1] D. Bhardwaj and V. Todwal, “A Novel Deep Learning Model for Indoor-Outdoor Scene Classification using VGG-16 Deep CNN”, *World Journal of Research and Review*, Vol. 13, No. 1, pp. 27-35, 2021.
- [2] M. Naseer, S. Khan, and F. Porikli, “Indoor Scene Understanding in 2.5/3D for Autonomous Agents: A Survey”, *IEEE Access*, Vol. 7, pp. 1859-1887, 2019.
- [3] U.A. Usmani, J. Watada, J. Jaafar, I.A. Aziz and A. Roy, “A Reinforced Active Learning Algorithm for Semantic Segmentation in Complex Imaging”, *IEEE Access*, Vol. 9, pp. 168415-168432, 2021.
- [4] C. Lang, G. Cheng, B. Tu, C. Li and J. Han, “Base and Meta: A New Perspective on Few-Shot Segmentation”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 45, No. 9, pp. 10669-10686, 2023.
- [5] M. Mahmoud, W. Chen, Y. Yang and Y. Li, “Automated BIM Generation for Large-Scale Indoor Complex Environments based on Deep Learning”, *Automation in Construction*, Vol. 162, No. 2, pp. 105376-105389, 2024.
- [6] Z. Li, “Open Rooms: An Open Framework for Photorealistic Indoor Scene Datasets”, *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 7186-7195, 2021.
- [7] M. Afif, R. Ayachi, Y. Said and M. Atri, “Deep Learning Based Application for Indoor Scene Recognition”, *Neural Processing Letters*, Vol. 51, No. 3, pp. 2827-2837, 2020.
- [8] D. Paschalidou, A. Kar, M. Shugrina, K. Kreis, A. Geiger, and S. Fidler, “ATISS: Autoregressive Transformers for Indoor Scene Synthesis”, *Advances in Neural Information Processing Systems*, Vol. 15, No. 2, pp. 12013-12026, 2021.
- [9] D. Zeng, “Deep Learning for Scene Classification: A Survey”, *Proceedings of IEEE International Conference on Computer Vision and Pattern Recognition*, pp. 1-24, 2021.
- [10] B.S. Anami and C.V. Sagarnal, “Influence of Different Activation Functions on Deep Learning Models in Indoor Scene Images Classification”, *Pattern Recognition and Image Analysis*, Vol. 32, No. 1, pp. 78-88, 2022.

- [11] Y. Liu, "Deep Learning based 3D Target Detection for Indoor Scenes", *Applied Intelligence*, Vol. 53, No. 9, pp. 10218-10231, 2023.
- [12] B.A. Labinghisa and D.M. Lee, "Indoor Localization System using Deep Learning based Scene Recognition", *Multimedia Tools and Applications*, Vol. 81, No. 20, pp. 28405-28429, 2022.
- [13] J.M. Sun, J. Yang, K. Mo, Y.K. Lai, L. Guibas and L. Gao, "Haisor: Human-Aware Indoor Scene Optimization via Deep Reinforcement Learning", *ACM Transactions on Graphics*, Vol. 43, No. 2, pp. 1-13, 2024.
- [14] L. Gao, J.M. Sun, K. Mo, Y.K. Lai, L.J. Guibas and J. Yang, "SceneHGN: Hierarchical Graph Networks for 3D Indoor Scene Generation with Fine-Grained Geometry", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 45, No. 7, pp. 8902-8919, 2023.
- [15] T.M.N.U. Akhund and K. Teramoto, "Privacy-Concerned Averaged Human Activeness Monitoring and Normal Pattern Recognizing with Single Passive Infrared Sensor using One-Dimensional Modeling", *Sensors International*, Vol. 6, No. 1, pp. 1-12, 2024.
- [16] R. Pereira, N. Goncalves, L. Garrote, T. Barros, A. Lopes and U.J. Nunes, "Deep-Learning based Global and Semantic Feature Fusion for Indoor Scene Classification", *Proceedings of IEEE International Conference on Autonomous Robot Systems and Competitions*, pp. 67-73, 2020.
- [17] Z. Wu, Z. Wang, S. Liu, H. Luo, J. Lu and H. Yan, "Fair Scene: Learning Unbiased Object Interactions for Indoor Scene Synthesis", *Pattern Recognition*, Vol. 156, No. 1, pp. 110737-110753, 2024.
- [18] R. Pereira, L. Garrote, T. Barros, A. Lopes, and U.J. Nunes, "Exploiting Object-based and Segmentation-based Semantic Features for Deep Learning-based Indoor Scene Classification", *Proceedings of IEEE International Conference on Computer Vision and Pattern Recognition*, pp. 1-12, 2024.
- [19] R. Pereira, L. Garrote, T. Barros, A. Lopes and U.J. Nunes, "A Deep Learning-based Indoor Scene Classification Approach Enhanced with Inter-Object Distance Semantic Features", *Proceedings of IEEE International Conference on Intelligent Robots and Systems*, Vol. 1, No. 4, pp. 32-38, 2021.
- [20] R. Pereira, T. Barros, L. Garrote, A. Lopes and U.J. Nunes, "A Deep Learning-Based Global and Segmentation-Based Semantic Feature Fusion Approach for Indoor Scene Classification", *Pattern Recognition Letters*, Vol. 179, No. 2, pp. 24-30, 2024.
- [21] A.M. Shaaban, N.M. Salem and W.I. Al-Albany, "A Semantic-Based Scene Segmentation using Convolutional Neural Networks", *AEU - International Journal of Electronics and Communications*, Vol. 125, pp. 153364-153378, 2020.
- [22] M. Afif, R. Ayachi, Y. Said and M. Atri, "An Indoor Scene Recognition System based on Deep Learning Evolutionary Algorithms", *Soft Computing*, Vol. 27, No. 21, pp. 15581-15594, 2023.
- [23] D. Seichter, M. Kohler, B. Lewandowski, T. Wengefeld and H. M. Gross, "Efficient RGB-D Semantic Segmentation for Indoor Scene Analysis", *Proceedings of IEEE International Conference on Robotics and Automation*, Vol. 2021, No. 2, pp. 13525-13531, 2021.
- [24] P.S. Yee, K.M. Lim and C.P. Lee, "DeepScene: Scene Classification via the Convolutional Neural Network with Spatial Pyramid Pooling", *Expert Systems with Applications*, Vol. 193, No. 2, pp. 116382-116395, 2022.
- [25] S.M. Yasir, A.M. Sadiq and H. Ahn, "3D Instance Segmentation Using Deep Learning on RGB-D Indoor Data", *Computers, Materials and Continua*, vol. 72, no. 3, pp. 5777-5791, 2022.
- [26] S.A. Dar, "Improving Alzheimer's Disease Detection with Transfer Learning", *International Journal of Statistics in Medical Research*, Vol. 14, pp. 403-415, 2025.
- [27] S.A. Dar, S. Palanivel, M.K. Geetha and M. Balasubramanian, "Mouth Image Based Person Authentication using DWLSTM and GRU", *Information Sciences Letters*, Vol. 11, No. 3, pp. 853-862, 2022.
- [28] S.A. Dar and S. Palanivel, "Performance Evaluation of Convolutional Neural Networks (CNNs) And VGG on Real Time Face Recognition System", *Advances in Science, Technology and Engineering Systems Journal*, Vol. 6, No. 2, pp. 956-964, 2021.
- [29] S.A. Dar and S. Palanivel, "Real Time Face Authentication System using Stacked Deep Auto Encoder for Facial Reconstruction", *International Journal of Thin Film Science and Technology*, Vol. 11, No. 1, pp. 73-82, 2022.
- [30] S.A. Dar, "Real-Time Face Authentication using Denoised Autoencoder (DAE) for Mobile Devices", *International Journal of Thin Film Science and Technology*, Vol. 21, No. 6, pp. 163-176, 2022.
- [31] W. Ayadi, "AI-Powered CNN Model for Automated Lung Cancer Diagnosis in Medical Imaging", *International Journal of Statistics in Medical Research*, Vol. 14, pp. 616-625, 2025.
- [32] S.A. Dar and R. Indhumathi, "Fracture Classification System in X-Ray Images and Method Employed Thereof", Indian Patent IN202441069005A, 2024.