



Detection and Classification of Oral Cancer using YOLO Object Detection Algorithm

Kavyashree C¹, H S Vimala¹ and Shreyas J²

¹Department of Computer Science and Engineering, UVCE, Bangalore University, Bangalore, India- 560001

²Department of Information Technology, Manipal Institute of Technology Bengaluru, Manipal Academy of Higher Education, Manipal, Karnataka, India - 576104

Received 4 August 2024, Revised 18 January 2025, Accepted 29 January 2025

Abstract: The early detection of oral cancer plays a pivotal role in enhancing the survival rate of the patients. Recent advancement in artificial intelligence have made diagnosis rapid and precise. The advent of deep learning has transformed medical image analysis, facilitating more precise, efficient and automated evaluations of medical images. It serves the purpose of identifying and locating particular objects within medical images. The aim of this research is to develop a deep learning-powered system for diagnosing oral cancer, capable of distinguishing between cancerous and non-cancerous areas in a provided image. The YOLO (You look only Once) is a cutting edge deep learning model employed for object detection, segmentation and classification. The system was retrained for the oral cancer dataset. The images are annotated with the help of the experts. A balanced dataset is created by data augmentation by rotating and flipping the images. The blurring is used to pre-process the images. The YOLOv8 architecture has been enhanced through the integration of EfficientNet-B0 for the generation of feature maps, along with the implementation of a Feature Pyramid Network (FPN), which facilitates the detection of objects across various scales. Following that, the model is trained with the images and then validated using YOLOv8 model. The normal and abnormal part of an images are identified with a precision of 0.901. The mean Average Precision (mAP) obtained for the model is 0.913. The YOLOv8 model is compared with other objection detection model such as YOLOv7, Mask RCNN (Region based Convolutional Neural Network) and Faster RCNN. YOLOv8 is found to be the fastest object detection and classification framework compared to the other three models. These results greatly help the medical practitioner to perform the initial investigation and help in the early detection of an oral cancer.

Keywords: Deep Learning methods, YOLO models, Object Detection, Classification, Oral Cancer

1. INTRODUCTION

Oral cancer is a form of cancer that develops in the tissues of the mouth or throat and is categorized as one of the head and neck cancers. It is essential to comprehend the causes, symptoms, diagnosis, treatment options, and preventive measures in order to effectively manage the condition and improve outcomes. Squamous Cell Carcinoma (SCC) is the predominant form of oral cancer, typically manifesting on the lining of the oral cavity. Tobacco usage is one of the largest risk factors for mouth cancer that includes chewing tobacco, smoking cigarettes [1]. Also, genes can occasionally cause cancer due to DNA mutations. According to the data, men have a higher likelihood of developing oral cancer than women [2]. Oral cancer is linked to a poor prognosis, with roughly 5 year survival rate of around 40%. Early detection not only enhances the survival rate by 80% but also aids in the formulation of an effective treatment plan. [3]. The oral cancer can be

diagnosed through biopsy which is considered as the highest standard in the cancer detection. Also imaging methods like MRI, CT and PET scans can be used for the cancer detection. A photographic images also can be used for the cancer detection. This study mainly focuses on the oral cancer detection using photographic images.

The development in artificial intelligence has assisted in automating the cancer detection by delivering precise and affordable outputs. Deep learning entails instructing artificial neural networks to identify patterns in data and make predictions. Processing input data through the network, calculating the output, and then evaluating it against the anticipated output helps in automating the process [4]. Object detection plays a vital role in computer vision by identifying and locating objects within medical images. YOLO, an acronym for You Only Look Once, is a widely utilized and effective deep learning model employed for



real-time object detection. It utilizes a classifier across different areas of an image, approaching the task as a unified regression problem [5]. In 2016, Ross Girshick, Ali Farhadi, Santosh Divvala, and Joseph Redmon [6] introduced YOLO to detect and categorize objects in images as well as in videos. YOLO receives a single image or video frame as input, it instantly creates bounding boxes around any objects that are discovered. The YOLO algorithm partitions an image into a grid, making predictions for bounding boxes and class probabilities in each grid cell. Combining these predictions yields the whole set of bounding boxes and the class with the probabilities. Deep Convolutional neural networks (CNN) are used to train the algorithm, which enables it to learn features that are characteristic of many object categories. The system is trained using massive datasets of labeled photos. There have been numerous iterations of YOLO, each with unique additions and improvements. An overview of some of the more popular variations is provided below: YOLOv1 was the original iteration of YOLO, released in 2016. Despite its speed and accuracy, there were some limitations, including challenges in identifying small objects. YOLOv2 is the second version of YOLO was released in 2017 and, in addition to a new architecture and the addition of features like batch normalization and anchor boxes, it addressed some of the shortcomings of the original version. It surpassed its predecessor regarding both accuracy and speed. YOLOv3 was released in 2018 with a few modifications to the model's performance, including the usage of residual blocks and feature pyramid networks and faster and more accurate. It also unveiled a novel feature extractor known as Darknet-53. YOLOv4 was released in 2020 with the implementation of Scaled-YOLOv4 and the usage of the CSPDarknet53 backbone, multi-input weighted residual connections, and other features. On a number of object detection benchmarks, it attained cutting-edge performance. YOLOv5 is a completely redesigned and re-implemented version of YOLO that was released in 2020. The use of anchor-free object detection, improved data augmentation methods, and more effective training are just a few of the enhancements it makes.

YOLOv6 was funded by Meituan in 2022 and used in robots. YOLOv7 implements pose estimation on the COCO keypoints dataset [7]. The most recent version of YOLO is called YOLOv8 by Ultralytics. YOLOv8 expands on the success of earlier editions state-of-the-art model by enhancing the features and performance. YOLOv8 supports a wide range of visual AI tasks such as detection, estimation, segmentation, tracking, and classification. Because of its adaptability, YOLOv8 can be employed in numerous applications[8] and we are using it for the oral cancer detection.

This study is focused on YOLO based technique to detect and classify the cancer affected areas in an oral cancer images. The YOLO framework is selected due to its ability to identify objects in real time with remarkable accuracy and reliability. It is capable of automating the screening

processes and can be effectively integrated with imaging techniques. This encourages the research to integrate YOLO for the early detection of oral cancer in the images.

The contributions made by this study are as follows:

- 1) Annotation of Oral cancer images using experts.
- 2) Augmentation techniques to improve the quality of the dataset.
- 3) This study employs YOLOv8 architecture in integration with EfficientNet B0 and FPN for object detection in images.
- 4) The performance analysis is performed comparing with other object detection models.

2. LITERATURE REVIEW

The objective of processing the medical images is to create computer aided algorithms to forecast the future developments in the health care systems. Thus, there are many algorithms developed using several deep learning techniques for detecting different cancers, including oral, prostate, skin, and breast cancer. This section provides the overview of various deep learning techniques used in the cancer detection such as YOLO models, RCNN (Region-based CNN) and transformer based methods. YOLO is an object detection model that operates in a single stage, segmenting an image into a grid while concurrently predicting bounding boxes and class probabilities. This model is particularly advantageous for applications that necessitate rapid and real-time analysis. CNN is employed to extract features and is utilized for tasks related to image classification. RCNN is capable of generating region proposals followed by classification. Although it offers greater accuracy compared to YOLO, it operates at a slower pace and is less appropriate for real-time analysis. Transformer-based approaches effectively identify relationships within images and demonstrate strong performance when applied to extensive datasets.

Gao et al. [16] integrated the YOLOv7 version with a coordinate attention mechanism to enhance the feature extraction process, ensuring that significant features are not overlooked. The streamlined Feature Pyramid Network (FPN) and the anchor-free model decrease the overall complexity. The sophisticated loss function contributes to enhancing both the accuracy and the resilience of the model. This comprehensive model is capable of identifying irregularities within dental images. Hsu et al. [14] used YOLOv7 model for oral mucous lesion detection and classification. The model is trained using 50,000 macroscopic images, each representing various grading. It categorizes images not solely as benign or malignant; it is also capable of identifying potentially malignant lesions. The YOLOv7-E6 variant demonstrated strong performance in both precision and recall, while the YOLOv7-D6 variant excelled in achieving a high F1-score. In addition to these metrics, the proposed model demonstrated enhanced accuracy in the detection of lesions. Mammeri et al. [17] used Yolov7 for

TABLE I. Summary of deep learning techniques used in Cancer detection

Author	Methodology	Modality	Results
[9]	Yolov5	Mamograms	mAP: 0.621
[10]	Mask RCNN	X- ray images	mAP : 0.900, F1-scores : 0.630 and precision : 0.960
[11]	Yolov8	Histopathological images	Accuracy: 0.970
[12]	Faster RCNN	CT Scan images	Accuracy : 0.986, sensitivity : 0.975, specificity : 0.96.8
[13]	Yolov3	Mamograms	mAP: 0.942, Accuracy :0.846
[14]	Yolov7	Macroscopic images	mAP : 0.638, Precision: 0.689, Recall : 0.658
[15]	Faster RCNN	Histopathological images	Accuracy :0 83.3 , recall : 0.718
[6]	Yolov8	Ultrasound images	mAP: 0.741, Accuracy : 0.880

lung nodule detection. This technology has the potential to create bounding boxes around the nodules, thereby assisting radiologists in tracking and identifying them within the entirety of the slide images. The model attained a mAP (mean average precision) of 81.28%, despite the absence of image preprocessing. It additionally conducted classification across multiple classes utilizing a pretrained VGG16 model. It also emphasized the importance of determining the degree of malignancy, which has the potential to enhance the overall diagnostic process. Prinzi et al. [9] employed the YOLOv5 model to identify suspicious objects within mammograms, thereby assisting in the detection of breast cancer. The implementation of transformers in place of convolution layers also facilitated the detection of smaller objects. Transformers assist in reducing the dimension of a vector to achieve a more compact output. Eigen – CAM notably decreases the occurrence of false negatives, although it results in a rise in false positives. The model could able to obtain 0.621 mAP, which can subsequently assist in the evaluation process. Chou et al. [18] proposed a YOLO based framework for esophageal cancer detection using versions v5 and v8. The application of enhancers for white light images has significantly enhanced the detection of carcinoma in comparison to RGB images. YOLOv5 has demonstrated superior performance with the dataset in comparison to YOLOv8, particularly excelling in feature learning capabilities, while YOLOv8 is noted for its high precision. The model underwent rigorous training for 500 iterations, achieving a precision of 0.85 with YOLOv5 and 0.81 with YOLOv8. Salman et al. [19] created a model for the prostate cancer diagnosis and grading of biopsy images using YOLOv3 by fine tuning its activation function for ReLu to Sigmoid and Tanh activation function because of the features being non linearly distributed. Pacal et al. [20] suggested a YOLO based identification of polyps by improving the performance using Cross Stage Partial Network (CSPNet) that enables instantaneous detection. The model's performance and clinical applicability are improved by the use of a large dataset during training.

Baccouche et al. [21] used YOLO model to identify and categorize worrisome breast lesions from the whole set of mammograms. The model is validated with public as well as the privately collected database. Also suggested the fusion model approach to enhance the performance. [22] also

proposed Image-to-image translation methods for mammography pairs data reconstruction. The model performs single class as well as multi class prediction for detecting the mass lesion. Karaman et al. [23] designed a system based on real-time identification of polyps with YOLOv5 that uses artificial bee colony to optimize the activation functions. The study employs the different versions of YOLO model saves the best model in order to enhance the activation processes and hyper parameters. Nersisson et al. [24] used YOLOv2 model to extract the features from lesion for skin cancer detection. Conventional characteristics, such as texture and color features are combined with the feature obtained from CNN. This fusion model has improved the overall performance in the skin cancer detection. Hamed et al. [25] compared the machine learning models with the YOLO and RetinaNet model and found YOLO provides the highest accuracy compared to other models. Aly et al. [13] used YOLO to detect mammograms and it is found to be the better object detection model compared to CNN models. The reason for the enhanced object detection is the use of anchor boxes in yolov3 that employs the k-means clustering algorithm. It provides the average precision of 94.2% and the classification accuracy of 84.6%. The study used a small annotated dataset and the major limitation is it could not detect the small masses. Nie et al. [26] employed the YOLOv1, v2 and v3 models for skin cancer identification and achieved a mean precision average of 0.82 with the limited training images of 200. This is designed for light weight applications like mobile applications. Victor et al. [27] used YOLOv3 to identify and categorize the skin cancer. YOLOv3 is used to generate the feature map and is combined with the color features extracted from Quad histogram. The integration of these features is input into the deep convolution neural network for cancer detection, resulting in a commendable level of accuracy. Ji et al. [28] used YOLO to find lung cancer in computed tomography (CT) images. A one stage model has been developed that improve the feature layer's overall multi-scale representation capability. They compared the model to other cutting-edge models and found that it performed better in terms of particular criteria like recall. Patel et al. [11] used YOLOv8 for breast cancer detection along with ResNet50 trained on BreakHis data. Data augmentation is performed to have a balanced dataset. This model performed better with accuracy of 97.8% and false positive rate of 1.2% suitable of real

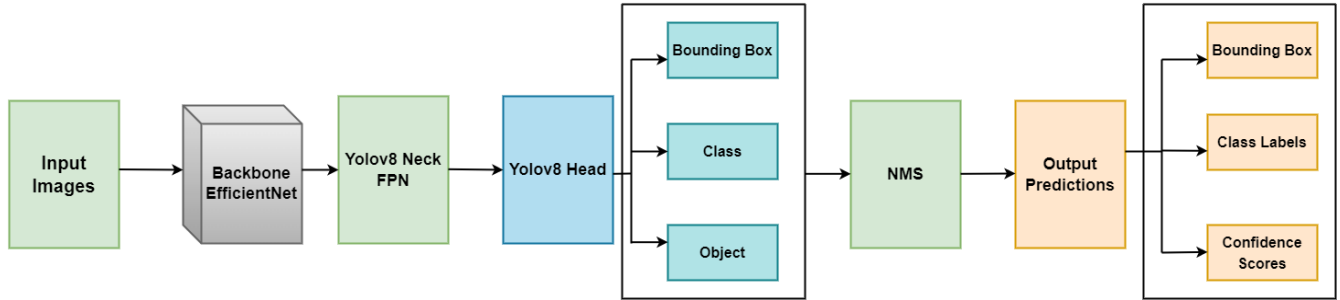


Figure 1. Proposed YOLOv8 architecture

time analysis. Pham et al. [6] used YOLOv8 to detect and classify tumors in ovary in the ultrasound images. The research indicates that YOLOv8 demonstrated a 19% increase in precision compared to YOLOv7, although it also exhibited a decrease in speed. The study showcases a comparative analysis of multiple iterations of YOLO models, highlighting the respective values observed. The research further determined that additional studies are necessary for the real-time analysis and detection of ovarian tumors and can not afford to miss even a single object. A method for object detection in panoramic X-ray images utilizing Mask RCNN has been proposed by [10]. This is predominantly used for instance segmentation. Winata et al. [15] used faster RCNN conjunction with the Adam optimizer to achieve improved outcomes. Jenipher et al. [12] used MobileNetV2 as backbone network, to localize the lung tumor for faster RCNN method. Table I summarizes the deep learning techniques used in literature study. The research provides a comprehensive understanding of the utilization of diverse input methods and the implementation of various object detection techniques. The effectiveness of cancer detection remains commendable; however, recent studies indicate a decline in average precision, suggesting a decrease in the accuracy of identifying abnormalities. Another concern is the requirement for substantial computational resources. Our research aims to explore the integration of EfficientNet, complemented by FPN and YOLOv8, as a means to enhance performance. This approach is anticipated to improve the detection of objects of varying sizes while minimizing the required computational resources.

3. METHODOLOGY

A. YOLOv8

YOLOv8 was designed by Glenn Jocher at Ultralytics and was released on January 10th, 2023. This most recent iteration of the well-known model is used for real-time object identification and image segmentation. Due to its rapid, accurate, and user-friendly design, YOLOv8 offers an excellent solution for various tasks such as object identification, classification of images and instance segmentation [29]. YOLOv8 estimates a target's center directly instead of predicting its offset from a predefined anchor box. A complex post-processing technique called Non-Maximum Suppression (NMS), which sorts through candidate detection after inference, is sped up with the absence of anchors

since it reduces the number of boxes being obtained. 3*3 convolution layers are added instead of 6*6 convolution layers.

A deep neural network called YOLOv8 employs a succession of convolution layers for the extraction of data from the obtained image in turn producing the bounding boxes and classes in the output. The architectural structure is composed of multiple essential elements, including the SPPF (Spatial Pyramid Pooling - Fast) layer, C2f module, Detection module, and the backbone. The backbone is responsible for extracting high-level features from the input image. The SPPF layer, along with subsequent convolutional layers, handles features at various scales. The C2f module integrates advanced functionalities with contextual data to enhance detection precision, while the detection module employs convolutional and linear layers to produce the bounding boxes and class probabilities. Along with these main components, there are several other layers used in the YOLOv8 architecture, include Up sample and Concat layers, which enhance the resolution of feature maps and merge feature maps from various layers, respectively. The COCO dataset is used for pretraining the Detect, Segment, and Pose models, while the ImageNet dataset is used for pretraining the Classify models. YOLOv8 offers different versions such as YOLOv8n, YOLOv8s, YOLOv8m, YOLOv8l, YOLOv8x that runs on different speed and used different parameters. YOLOv8 provides the high accuracy rate according to COCO and Roboflow 100 measurements and also easy for the developers to use. In other words, instead of forecasting the deviation of an object from a recognized anchor box, it directly forecasts the object's center. This provides a solution to the major difficult component of previous YOLO models.

B. Proposed YOLOv8 based object detection for oral cancer detection

Figure 1 shows the proposed YOLOv8 architecture used in the study. It can be integrated with various backbone like CSPDarknet53, EfficientNet, ResNet, Inception and other pretrained models for feature extraction. CSPDarknet53 proves to be highly effective and perfectly compatible with yolov8, providing the ideal combination of balance and top-notch performance. Due to its computational complexity and superior accuracy, EfficientNet has been chosen as

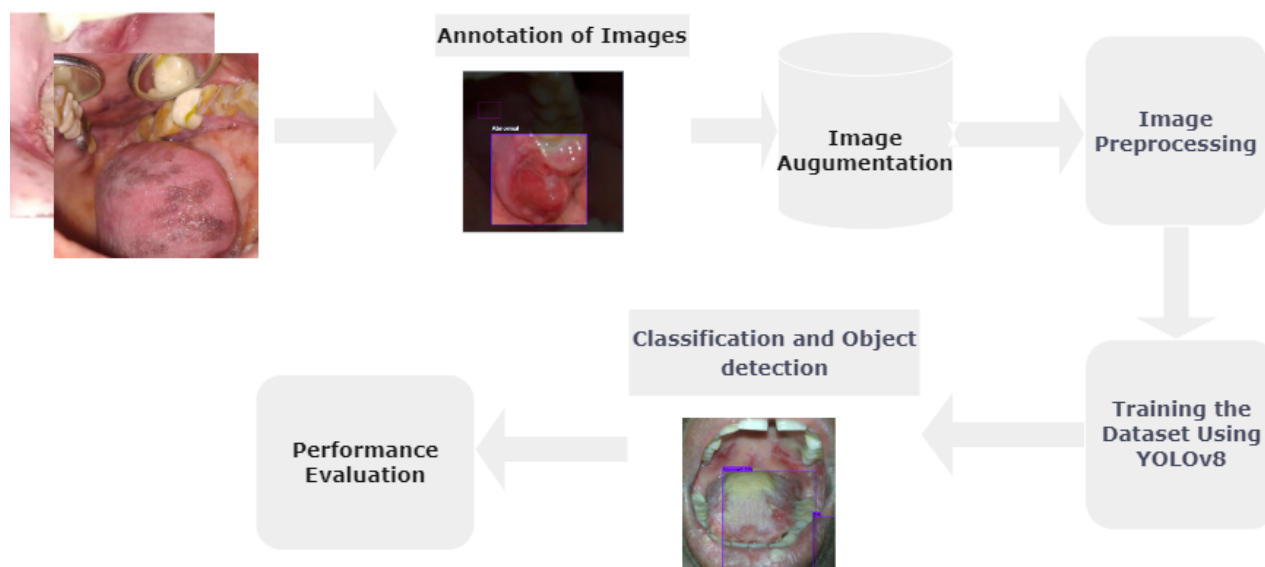


Figure 2. Proposed Design Methodology

the backbone network. EfficientNet-B0 was specifically selected for its lightweight design, while still delivering impressive accuracy at a lower computational expense. EfficientNet utilizes depth-wise separable convolutions as part of its design to maintain simplicity. This aids in lowering the computational expenses while maintaining the same level of accuracy. Accurate feature extraction is crucial for health care needs, and our research with limited resources benefited greatly from the use of EfficientNet-B0. The major architectural change made in YOLOv8 is the use of EfficientNet-B0 as a backbone network to extract the features at the multiple depths. This helps YOLOv8 to use the lower level as well as higher level semantics, thus the model has the ability to detect the objects of different size and enabling it to accept input and process it through the YOLOv8 head.

The neck merges the characteristics derived from the backbone and transmits them to the head for the ultimate detection and classification. The proposed methodology uses Feature pyramid network (FPN) over SPPF in YOLOv8 that aggregates the feature from different levels. This boosts the model's capability to recognize objects of diverse sizes. FPN enhances higher-level features through up sampling, integrates them with lower-level features, and subsequently forwards the combined information to the next stage. FPN is utilized in conjunction with convolutional layers to enhance the model's performance. The neck must possess sufficient efficiency to support the functions of the spine without impeding its performance. In our utilization of YOLOv8 for the analysis of medical images, we prioritize accuracy above speed and complexity.

The head of the network is responsible for producing

the final output. It processes the features extracted and aggregated by the backbone and neck components to generate predictions such as bounding boxes, class scores, and objectness scores. The head plays a vital role in transforming the feature maps into significant detection outcomes. YOLOv8 employs a regression-based method to anticipate the coordinates of the bounding box surrounding the identified objects. The Objectness Score is a measure of the probability that a bounding box encloses an object, allowing the distinction between background and object. The classification layers determine the likelihood of each class for the identified object.

4. EXPERIMENTAL SETUP

Figure 2 shows the flow of the study from the image acquisition to the performance evaluation.

A. Dataset Preparation

Dataset preparation is crucial in the proposed algorithm as it is founded on the deep neural network. Oral cancer photographic image dataset is taken from public repository from Roboflow platform [30] that contains 323 images. The annotation of these images are crucial that greatly helps in the detection. The annotation of the images is done with the help of experts from Suraksha Speciality Dental Care, Hoskote, India. The team contributed in understanding of oral cancer screening, which facilitated the annotation of the images. The labels used for annotation are Abnormal and Normal regions. The precise annotation leads to better detection of objects and improves the overall performance. The major challenge encountered in annotation is the overlapping of Normal and Abnormal regions. However we have tried to focus on labeling Abnormal regions as the objective of the study is to identify the cancer affected region. Figure

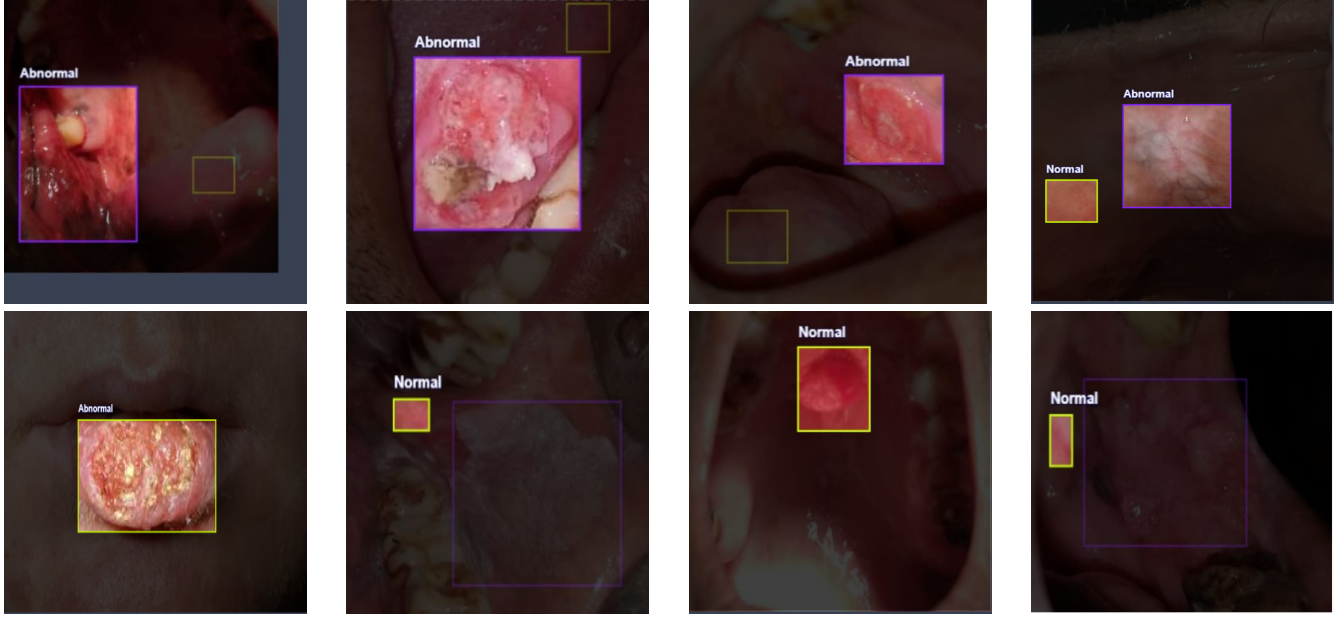


Figure 3. Images annotated with Normal and Abnormal labels

3 shows the annotation of oral cancer images as Normal and Abnormal objects. Multiple objects are present in each image.

B. Data Augmentation and Preprocessing

Utilizing balanced datasets can lead to more accurate outcomes and mitigate the risk of overfitting the model. The dimension of the dataset is increased with augmentation through horizontal flipping, rotating the images and blurring upto 10px. Flipping assists in recognizing images from various orientations by reflecting an image across both the horizontal and vertical axes. Rotation fulfills a similar function, enhancing robustness and reducing sensitivity to orientation. This is essential for real-time analysis, as the images are not accurately aligned. In our research, a 90° rotation is employed; however, it is compatible with any angle. The purpose of blurring is to eliminate noise from images, thereby enhancing the model's learning capabilities. The integration of these augmentation techniques enhances the dataset's size, facilitates learning from features, and aids in addressing issues related to class imbalance. The other augmentation techniques, such as cropping, introducing noise, and adjusting the brightness of the images, were not employed in our study, as the quality of our images is deemed sufficient for use. The augmentation process aimed to increase the dataset's size while also ensuring its balance. The data set is significantly increased to 689 images. The images are preprocessed by resizing it to 640*640 as required for YOLOv8 architecture. The augmented and preprocessed dataset is divided with 80% training set, 12% testing set and 8% for validation.

The YOLOv8 is used to train the model with the training data. It uses the sigmoid linear units (SiLU) [31] as an

activation function that replaces the leaky Relu only in the convolution and batch normalization layer of CNN. The outputs are taken from the sigmoid function.

$$SiLU(x) = x \frac{1}{1 + e^{-x}} \quad (1)$$

The intersection over union (IoU) quantifies the overlap between the predicted bounding box and the ground truth bounding box. Figure. 4 shows the intersection of the bounding boxes and extracting the bounding box from the background.

$$IoU \text{ Score} = \frac{\text{Area of Intersection}}{\text{Area of Union}} \quad (2)$$

A loss function helps in adjusting the weights of the neural network to reduce the cost. YOLOv8 employs C_{IoU} (Complete Intersection over Union) and the DFL (Distributed Focal Loss) for BBox (bounding Box) loss and Binary Cross Entropy (BCE) for Classification (cls) loss. BBox loss measures how closely do the cls loss evaluations and the expected ground truth bounding boxes match and how accurately each anticipated bounding box was classified. DFL loss takes care of the class imbalance issue in training the neural network and optimizes the bbox boundary distribution. The C_{IoU} loss considers three elements: the overlap area, the distance between the center points of the boxes, and the aspect ratio. The loss function for C_{IoU} is given by

$$\mathcal{L}_{CIU} = 1 - IoU + \frac{\rho^2(b, b^{gt})}{c^2} + av \quad (3)$$

where b , B^{gt} represents the intersection of background with

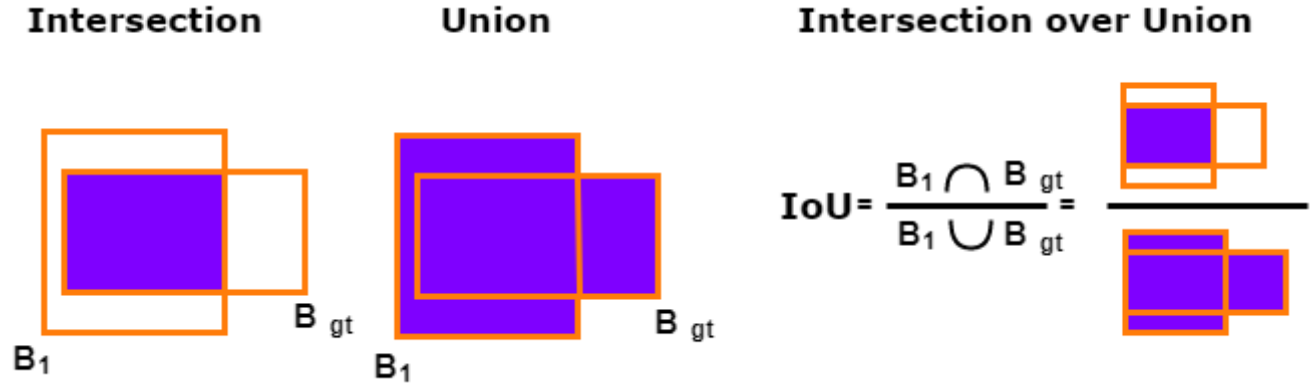


Figure 4. Intersection over Union

ground truth, ρ is the obtained euclidean distance. The smallest box that encloses the two boxes is represented by c , which is the diagonal measurement. α represents the positive trade-off criterion and v gauges aspect ratio's constancy. α can be given as

$$\alpha = \frac{a}{(1 - IoU) + a} \quad (4)$$

The specified a is

$$a = \frac{4}{\pi^2} \left(\arctan \frac{wt^{gt}}{ht^{gt}} - \arctan \frac{wt}{ht} \right)^2 \quad (5)$$

a defines the applying inverse tan function to the the bounding box's height and width.

BCE loss is utilized to measure the difference between the predicted category and the actual labels of the dataset. BCE loss is determined as

$$\mathcal{L}_{BCE} = \frac{1}{A} \sum_{i=1}^A -(m_i * \log(p_i)) + ((1 - m_i) * \log(1 - p_i)) \quad (6)$$

p_i represents probability of Normal class and $(1 - p_i)$ is the probability of class Abnormal. YOLO has originally trained using COCO dataset to identify 80 object classes like cars, books, handbags, phones and so on. This pretrained model has to be retrained to recognize normal or abnormal part in an oral cancer images by fine tuning the hyper parameters. The roboflow framework is used to custom train the model for oral cancer images. The model is custom trained using the YOLOv8 model for the dataset. The classification of an image as Normal and Abnormal is carried out using the test set. The bounded box is used to identify the cancerous part in an image. YOLOv8 is crafted to deliver enhanced speed and accuracy, but the performance of the model can be greatly enhanced by fine-tuning the hyperparameters. Table II shows hyper parameters considered for fine tuning the model:

TABLE II. Hyper parameters for Training the model

Hyper parameters	Value
Image Size	640*640
Batch size	32
Learning rate	1.00E-03
Epochs	25
Object threshold	0.5
NMS Threshold	0.4
Weight decay	0.0001

5. PERFORMANCE ANALYSIS

In this investigation, the performance metrics considered are mAP, precision and recall. The term precision is used to assess how accurate and dependable the experiment's measurement is. In applications such as medical diagnosis, high precision is essential to minimize the occurrence of false positives. Such inaccuracies may result in unwarranted treatments that could pose risks to patients. Recall measures the model's ability to recognize positive samples. Ensuring a high recall rate is critically important, as it may have life-threatening implications or result in ineffective treatment. The Average Precision for each class is determined by constructing the precision-recall curve and calculating the area under this curve (AUC). The mAP is calculated by averaging the Average Precision (AP) scores for all classes. This metric offers a consolidated score that reflects the overall effectiveness of the model. mAP helps in identifying multiple conditions by balancing the precision and recall. Upon obtaining the True Positive (TP), False Positive (FP), True Negative (TN) and False Negative (FN) values from the confusion matrix, the metrics can be formulated as following:

$$Precision = \frac{TP}{TP + FP} \quad (7)$$

$$Recall = \frac{TP}{TP + FN} \quad (8)$$

TABLE III. Table of the confusion matrix for the test data

Class	Images	Labels	Precision	Recall	mAp@.5	mAP@.5:.95
All	56	76	0.873	0.875	0.913	0.734
Normal	56	29	0.901	0.889	0.972	0.723
Abnormal	56	46	0.812	0.873	0.896	0.768

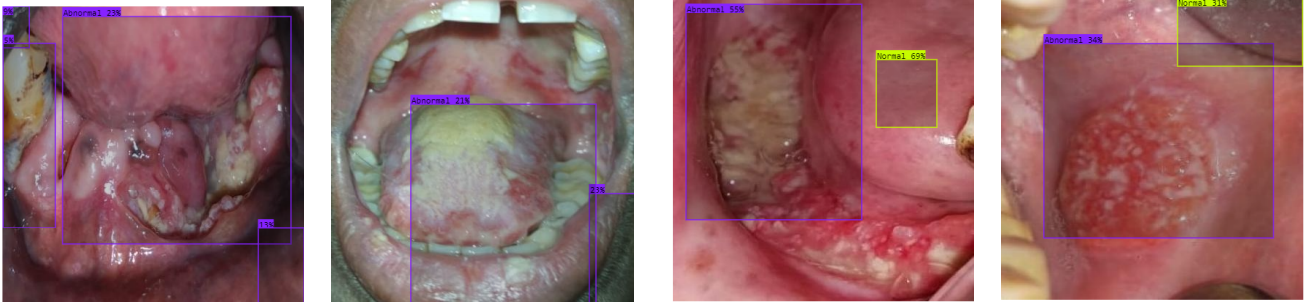


Figure 5. Identification of Normal and Abnormal part in the test images

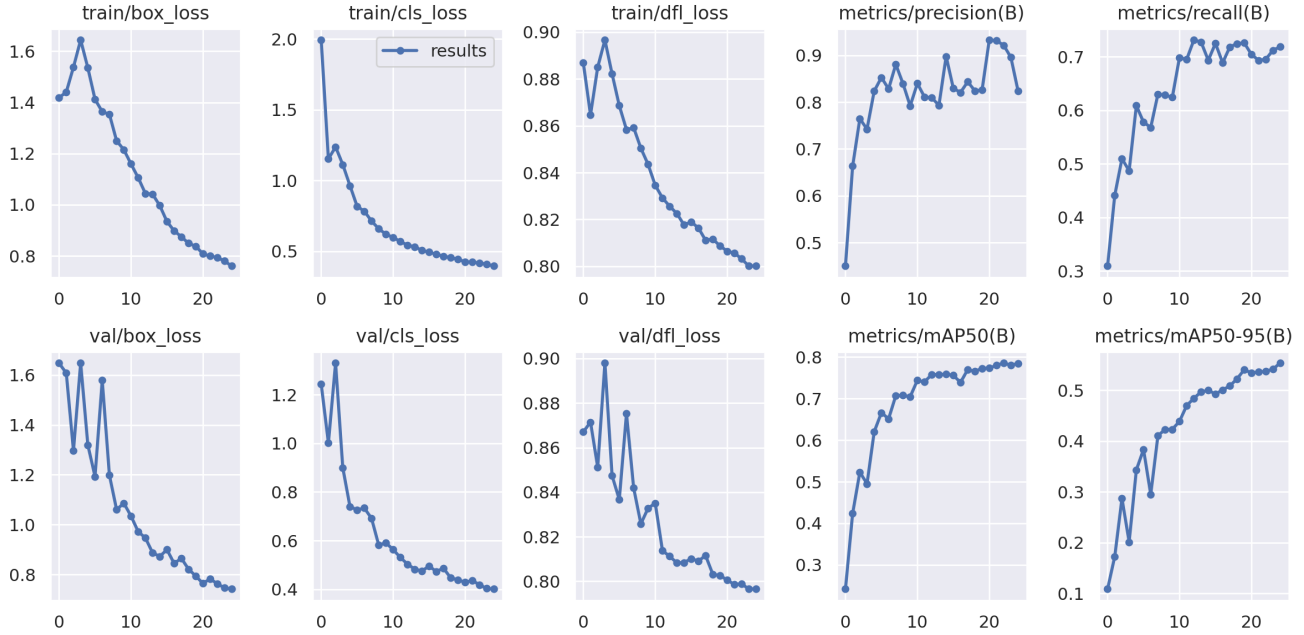


Figure 6. Analysis of metrics for 25 epochs for training and validation

$$mAP = \frac{1}{N} \sum_{i=1}^N AP_i \quad (9)$$

TABLE III shows confusion matrix table for the test data and shows the values of precision, recall and mAP for the normal and abnormal regions. The IoU threshold considered to calculate mAP is 0.5 and also till 0.95 with 0.05 increments. Both classes had good and noticeably better mAP ratings for precision and recall. Figure 5 shows detection of normal and abnormal class with level of confidence on test images. Figure 6 shows the various performance metrics

after the model runs for 25 epochs. It has both the training and validation values for reference. The loss is significantly reducing as the training progresses. The precision and recall are improved in both training and validation. Figure 7 shows the identification of Normal and Abnormal classes in the validation dataset. It also shows that the multiple objects are detected in an image with different sizes.

A. Comparative Analysis

The results obtained are evaluated against leading object detection methodologies, including Faster RCNN and Mask RCNN, both of which are part of the RCNN family, as well as the earlier version of the YOLO model, v7. YOLOv8

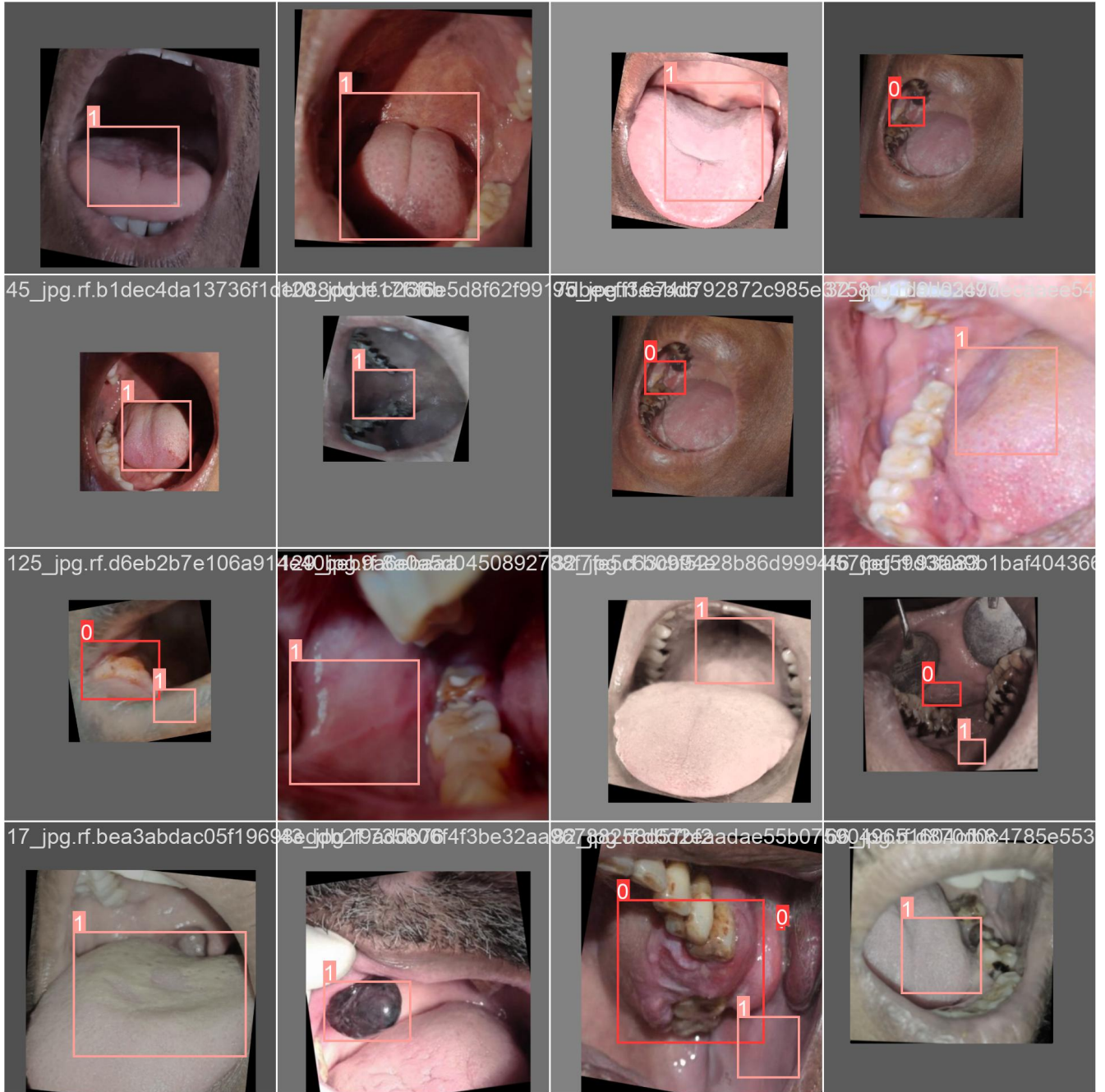
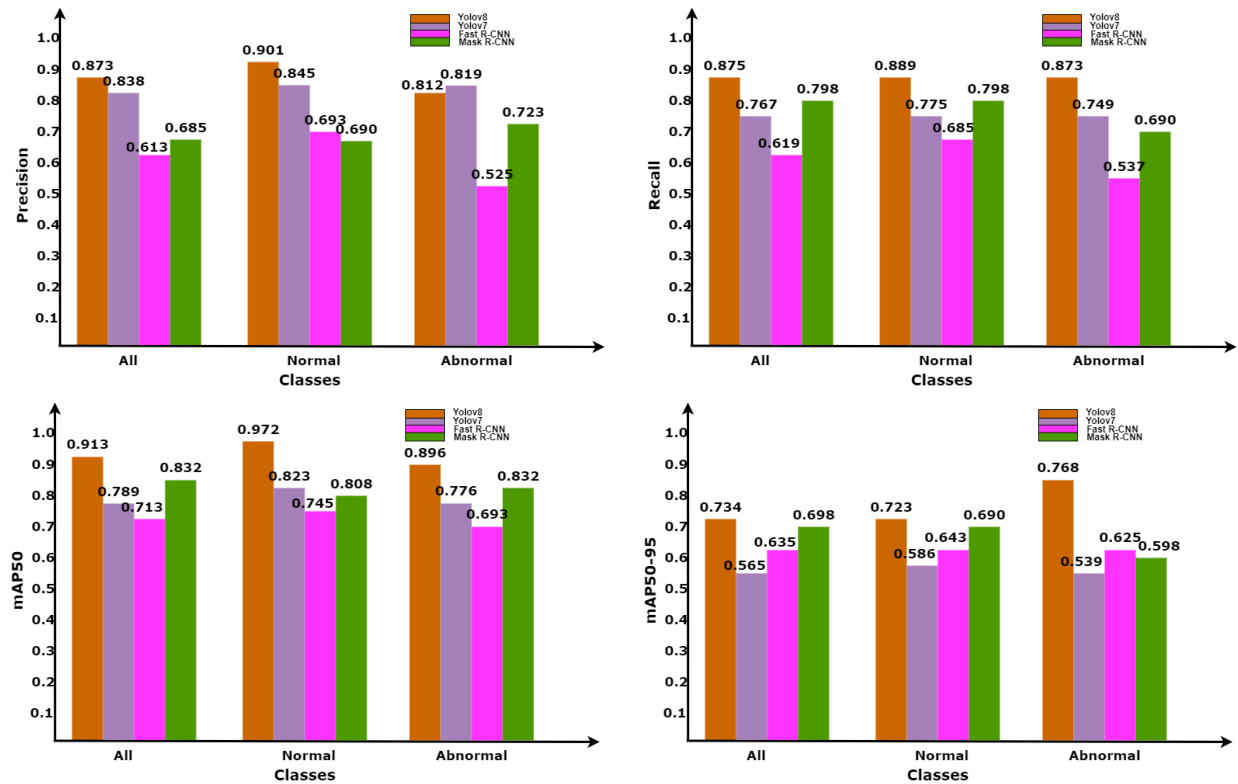


Figure 7. Detection of Normal and Abnormal regions for the Validation dataset

Figure 8. Comparison of YOLOv8 with YOLOv7, Mask RCNN and Faster RCNN



employs a unified architecture that enables object detection to occur in a single pass, resulting in improved speed by minimizing latency. YOLOv7 provides the better training strategies and performance similar to YOLOv8. Faster RCNN employs a two-stage methodology that encompasses both detection and classification processes. In contrast, Mask RCNN enhances this framework by incorporating a segmentation mask in addition to the features provided by Faster RCNN. The RCNN methodologies tend to be slower in comparison to YOLO models; however, they provide greater accuracy regarding precision and mAP when dealing with complex datasets. The selection of YOLOv8 over alternative models is attributed to its capability to identify smaller objects in real-time with exceptional accuracy and rapidity. This feature could prove advantageous in emergency scenarios within healthcare facilities.

Consequently, a comparative analysis offers a comprehensive overview of the comparison between object detection algorithms, highlighting the performance of YOLOv8 in relation to other leading models in the field. The same data set is applied to evaluate the performance of the different models. Figure 8 shows the comparative analysis of different object detection models in terms of precision, recall, mAP50 with the threshold of 50% overlapping with the ground truth and mAP50:5 with 0.05 increments for the classes Normal, Abnormal and all. From Figure 8, results shows that Yolov8 outperforms the other models in terms

of recall, mAP, whereas the precision is slightly less in identifying the Abnormal classes. This led to false positive but the good recall value ensures the false negatives are identified and hence suitable for health care applications. The performance of the Yolov8 can be enhanced by improving the quality and quantity of dataset.

6. CONCLUSION AND FUTURE WORK

In this work, a sophisticated approach for identifying and classifying oral cancer was created through the integration of YOLOv8 with EfficientNet-B0 and FPN. This collaboration optimizes feature extraction, enhances multi-scale representation, and preserves real-time processing capabilities, rendering it an essential asset in clinical practice. The roboflow provides a framework to perform the annotation of images. The augmentation of images improved the quality of the dataset. The pretrained YOLOv8 model is used to train 689 images. The model is validated using the test data and it is found that it could able to detect the normal and abnormal part in the given test image with a better confidence level. The system's outcomes show that it is possible to successfully identify and categorize normal and abnormal parts in a picture using deep learning method that is assisted by the YOLO object detection method. A comparison analysis is provided with the other object detection model with the YOLOv8 is providing the better results. In the future, we intend to generate a data set with approximately 10,000 images that are to be

collected from various hospitals to increase the system's resiliency. An ensemble model approach, which integrates various models, can improve detection capabilities while minimizing both false positives and false negatives. Hybrid architectures have the potential to significantly enhance this process, leading to improved patient outcomes across various medical applications.

REFERENCES

- [1] R. Serra, C. S. de Oliveira, S. Roque, F. Herrera, and H. Arco, "Oral hygiene care and the management of oral symptoms in patients with cancer in palliative care: a mixed methods systematic review protocol," *JBI Evidence Synthesis*, vol. 22, no. 4, pp. 673–680, 2024.
- [2] R. L. Siegel, A. N. Giaquinto, and A. Jemal, "Cancer statistics, 2024," *CA: a cancer journal for clinicians*, vol. 74, no. 1, 2024.
- [3] E. Saberian, A. Jenča, A. Petrášová, J. Jenčová, R. A. Jahromi, and R. Seiffadini, "Oral cancer at a glance," *Asian Pacific Journal of Cancer Biology*, vol. 8, no. 4, pp. 379–386, 2023.
- [4] H. A. Helaly, M. Badawy, and A. Y. Haikal, "A review of deep learning approaches in clinical and healthcare systems based on medical image analysis," *Multimedia Tools and Applications*, vol. 83, no. 12, pp. 36 039–36 080, 2024.
- [5] A. Vijayakumar and S. Vairavasundaram, "Yolo-based object detection models: A review and its applications," *Multimedia Tools and Applications*, pp. 1–40, 2024.
- [6] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [7] P. Jiang, D. Ergu, F. Liu, Y. Cai, and B. Ma, "A review of yolo algorithm developments," *Procedia computer science*, vol. 199, pp. 1066–1073, 2022.
- [8] T. Diwan, G. Anirudh, and J. V. Temburne, "Object detection using yolo: Challenges, architectural successors, datasets and applications," *Multimedia Tools and Applications*, vol. 82, no. 6, pp. 9243–9275, 2023.
- [9] F. Prinzi, M. Insalaco, A. Orlando, S. Gaglio, and S. Vitabile, "A yolo-based model for breast cancer detection in mammograms," *Cognitive Computation*, vol. 16, no. 1, pp. 107–120, 2024.
- [10] W. Brahmi and I. Jdey, "Automatic tooth instance segmentation and identification from panoramic x-ray images using deep cnn," *Multimedia Tools and Applications*, vol. 83, no. 18, pp. 55 565–55 585, 2024.
- [11] V. Patel, M. Kanojia, and V. Nair, "Exploring the potential of resnet50 and yolov8 in improving breast cancer diagnosis: A deep learning perspective," *International Journal of Computer Information Systems and Industrial Management Applications*, vol. 16, no. 3, pp. 16–16, 2024.
- [12] V. N. Jenipher and S. Radhika, "Lung tumor cell classification with lightweight mobilenetv2 and attention-based scam enhanced faster r-cnn," *Evolving Systems*, pp. 1–18, 2024.
- [13] G. H. Aly, M. Marey, S. A. El-Sayed, and M. F. Tolba, "Yolo based breast masses detection and classification in full-field digital mammograms," *Computer methods and programs in biomedicine*, vol. 200, p. 105823, 2021.
- [14] Y. Hsu, C.-Y. Chou, Y.-C. Huang, Y.-C. Liu, Y.-L. Lin, Z.-P. Zhong, J.-K. Liao, J.-C. Lee, H.-Y. Chen, J.-J. Lee *et al.*, "Oral mucosal lesions triage via yolov7 models," *Journal of the Formosan Medical Association*, 2024.
- [15] P. A. Winata and I. Roysida, "Implementation of a faster r-cnn algorithm for identification of metastatic tissue using lymphoma histopathological images," *Journal of Soft Computing Exploration*, vol. 4, no. 2, 2023.
- [16] L. Gao, T. Xu, M. Liu, J. Jin, L. Peng, X. Zhao, J. Li, M. Yang, S. Li, and S. Liang, "Ai-aided diagnosis of oral x-ray images of periapical films based on deep learning," *Displays*, vol. 82, p. 102649, 2024.
- [17] S. Mammeri, M. Amroune, M.-Y. Haouam, I. Bendib, and A. Corrêa Silva, "Early detection and diagnosis of lung cancer using yolo v7, and transfer learning," *Multimedia Tools and Applications*, vol. 83, no. 10, pp. 30 965–30 980, 2024.
- [18] C.-K. Chou, R. Karmakar, Y.-M. Tsao, L. W. Jie, A. Mukundan, C.-W. Huang, T.-H. Chen, C.-Y. Ko, and H.-C. Wang, "Evaluation of spectrum-aided visual enhancer (save) in esophageal cancer detection using yolo frameworks," *Diagnostics*, vol. 14, no. 11, p. 1129, 2024.
- [19] M. E. Salman, G. Ç. Çakar, J. Azimjonov, M. Kösem, and İ. H. Cedimoğlu, "Automated prostate cancer grading and diagnosis system using deep learning-based yolo object detection algorithm," *Expert Systems with Applications*, vol. 201, p. 117148, 2022.
- [20] I. Pacal, A. Karaman, D. Karaboga, B. Akay, A. Basturk, U. Nalbantoglu, and S. Coskun, "An efficient real-time colonic polyp detection with yolo algorithms trained by using negative samples and large datasets," *Computers in biology and medicine*, vol. 141, p. 105031, 2022.
- [21] A. Baccouche, B. Garcia-Zapirain, C. C. Olea, and A. S. Elmaghraby, "Breast lesions detection and classification via yolo-based fusion models," *Computers, Materials & Continua*, vol. 69, no. 1, 2021.
- [22] A. Baccouche, B. Garcia-Zapirain, Y. Zheng, and A. S. Elmaghraby, "Early detection and classification of abnormality in prior mammograms using image-to-image translation and yolo techniques," *Computer Methods and Programs in Biomedicine*, vol. 221, p. 106884, 2022.
- [23] A. Karaman, I. Pacal, A. Basturk, B. Akay, U. Nalbantoglu, S. Coskun, O. Sahin, and D. Karaboga, "Robust real-time polyp detection system design based on yolo algorithms by optimizing activation functions and hyper-parameters with artificial bee colony (abc)," *Expert systems with applications*, vol. 221, p. 119741, 2023.
- [24] R. Nersisson, T. J. Iyer, A. N. Joseph Raj, and V. Rajangam, "A dermoscopic skin lesion classification technique using yolo-cnn and traditional feature model," *Arabian Journal for Science and Engineering*, vol. 46, no. 10, pp. 9797–9808, 2021.
- [25] G. Hamed, M. A. E.-R. Marey, S. E.-S. Amin, and M. F. Tolba, "Deep learning in breast cancer detection and classification," in *Proceedings of the International Conference on Artificial Intelligence and Computer Vision (AICV2020)*. Springer, 2020, pp. 322–333.
- [26] Y. Nie, P. Sommella, M. O'Nils, C. Liguori, and J. Lundgren,



- “Automatic detection of melanoma with yolo deep convolutional neural networks,” in *2019 E-Health and Bioengineering Conference (EHB)*. IEEE, 2019, pp. 1–4.
- [27] S. O. Manoj, K. R. Abirami, A. Victor, and M. Arya, “Automatic detection and categorization of skin lesions for early diagnosis of skin cancer using yolo-v3-dcnn architecture,” *Image Analysis and Stereology*, vol. 42, no. 2, pp. 101–117, 2023.
- [28] Z. Ji, J. Zhao, J. Liu, X. Zeng, H. Zhang, X. Zhang, and I. Ganchev, “Elct-yolo: an efficient one-stage model for automatic lung tumor detection based on ct images,” *Mathematics*, vol. 11, no. 10, p. 2344, 2023.
- [29] ultralytics, “Yolov8 model is used from <https://docs.ultralytics.com/>.” NA, 2023.
- [30] —, “Oral cancer image dataset is used from <https://universe.roboflow.com/sagari-vijay/oral-cancer-data/dataset/1>.” NA, 2023.
- [31] G. Jocher, A. Stoken, J. Borovec, L. Changyu, A. Hogan, A. Chaurasia, L. Diaconu, F. Ingham, A. Colmagro, H. Ye *et al.*, “ultralytics/yolov5: v4. 0-nn. silu () activations, weights & biases logging, pytorch hub integration,” *Zenodo*, 2021.
-