

# Deep Learning-Based Multiclass Classification of Oral Lesions with Stratified Augmentation

Joy Naoum, Revana Salama, Ali Hamdi

<sup>1</sup> joy.ehab , <sup>2</sup> revana.magdy , <sup>3</sup> alihamdif@msa.edu.eg

MSA University, Giza, Egypt

## Abstract

Oral cancer is highly common across the global and is mostly diagnosed during the later stages due to the close visual similarity to benign, precancerous, and malignant lesions in the oral cavity. Implementing computer-aided diagnosis (CAD) systems early on has the potential to greatly improve clinical outcomes. This research intends to use deep learning to build a multi-class classifier for sixteen different oral lesions. To overcome the challenges of limited and imbalanced data sets, the proposed technique combines stratified data splitting and advanced data augmentation and oversampling to perform the classification. The experimental results, which achieved 83.33% accuracy, 89.12% precision, and 77.31% recall, demonstrate the superiority of the suggested model over the state-of-the-art methods now in use. The suggested model effectively conveys the effectiveness of oversampling and augmentation strategies in situations where the minority class's classification performance is noteworthy. As a first step toward trustworthy computer-aided diagnostic (CAD) systems for the early detection of oral cancer in clinical settings, the suggested framework shows promise.

Index Terms—Data Augmentation, Oral Cancer, Stratified Sampling, multi-classification, Deep Learning.

## 1 Introduction

Oral cancer poses a significant health challenge worldwide, with approximately 377,000 new cases and 177,000 deaths annually [1]. The stage of the disease at diagnosis fundamentally determines the patient's prognosis. If diagnosed at an early stage, oral cancer has a five-year survival rate of almost 80%, but late-stage diagnosis drops survival rates to below 30% [18]. The clinical diagnosis of oral cancer is further complicated by the overlapping visual features of benign, precancerous, and malignant lesions. These challenges are a major contributor to the worldwide health burden due to misdiagnosis.

Conventional computer-aided diagnostic (CAD) systems have primarily focused on binary classification, i.e., the identification of malignant lesions and the separation of non-cancerous tissues. Such systems may be useful in preliminary screenings but are inadequate in real-world scenarios where practitioners are required to differentiate between various lesion subtypes. The need for multi-class classification systems to recognize the nuanced differences between lesion types is vital for supporting clinicians in reducing diagnostic uncertainty.

Deep learning, specifically convolutional neural networks, has made tremendous strides in areas of medical imaging, including dermatology, radiology, and pathology [12]. Using pre-trained models like ResNet, VGG, and EfficientNet, implementing transfer learning has been moving better than feature engineering, and is providing powerful feature extraction even when

working with small datasets. However, when it comes to research on detecting oral cancer, there has been limited exploration due to deep learning’s restrictions in two notable ways: (i) the majority of prior work continues to concentrate on binary classification, and (ii) multiclass approaches are primarily restricted due to limited and imbalanced datasets, resulting in models that generalize poorly and that are biased in their predictions.

To mitigate these limitations, we suggest a **comprehensive multiclass oral lesion classification pipeline** leveraging stratified dataset partitioning, novel augmentation techniques, and oversampling techniques. This work intends to demonstrate that, **on top of increasing the depth and complexity of the network, performance can also be improved through dataset engineering—stratified sampling and systematic augmentation**. The proposed pipeline focuses on 16 types of oral lesions and aims to achieve the systemic reduction of class imbalance, increasing intraclass variation, and diminishing overfitting to improve the precision and recall on rare and frequent lesions.

Experimental results highlight that the proposed model surpasses the state-of-the-art methods with 83.33% accuracy, 89.12% precision, and 77.31% recall. In particular, augmentation and oversampling significantly enhanced the classification performance of underrepresented categories of lesions, confirming their importance in contexts of imbalanced medical imaging.

This study has made the following contributions:

- Completion of a multiclass classification pipeline for oral lesions comprising 16 classes, providing a bridge between research prototypes and diagnostic support in clinical practice
- Showing that augmentation and oversampling techniques can improve model generalization and robustness in the context of limited, imbalanced medical datasets.
- Validation of the role of data-centric designs in advancing computer-aided oral cancer diagnosis by demonstrating the ability to surpass the current best methods in the field.

This study also addresses both the technical and clinical challenges, laying the groundwork for the development of scalable CAD tools that can aid in the early detection of oral cancer, and enable healthcare practitioners to make precise and timely decisions.

**Organization of the Paper**— The structure of the paper as the following : Section 2 reviews related work on oral lesion classification and deep learning methods in medical imaging. Section 3 describes the dataset, preprocessing steps, augmentation strategies, and the proposed model architecture. Section 4 presents the experimental setup and results, followed by an in-depth discussion in Section 5. Section 6 concludes the study with key findings, while Section 7 outlines potential directions for future research.

## 2 Related Work

Recent discoveries in deep learning have improved the automatic classification of oral lesions using photographic imagery. Al-Ali et al. [1] extend this line of work by developing CLASEG, a multiclass classification framework based on EfficientNet-B3 and trained on 2,072 clinical images across 14 types of lesions. Their algorithm demonstrated a classification accuracy of 74.49%, cementing the clinical utility of consumer grade images to classify oral lesions into benign, premalignant and malignant categories, and deviated from earlier studies which focused on binary classification or demanded expensive imaging modalities. While most studies have centered on multiclass prediction problems, CLASEG focused on classes within a single image. We adopt this approach and implement EfficientNetV2B1 to develop a streamlined

pipeline comprising stratified sampling, purpose-driven augmentation, and class-balanced training. This provides a clinically handy solution to oral cancer screening with proven improved classification and recall. In a complementary approach, Warin et al. [18] demonstrated that deep CNNs like DenseNet-169 were able to generate almost expert-level performance in the detection of normal tissue, potentially malignant disorders, and oral squamous cell carcinoma with area-under-the-curve (AUC) scores approaching 1.0. Their article also included object detection models like Faster R-CNN and YOLOv5 for detecting suspicious lesions, further establishing the role of AI in the facilitation of clinical decision-making. Subsequent work [12], [11] followed by combining classification with lesion segmentation enabled precise boundary definition and differential diagnosis between several lesion subtypes. All these advances point to an unmistakable trend toward multiclass, multi-task deep learning architectures based on standard clinical photographs to provide robust, interpretable, and scalable screening tools for oral cancer. Das et al. [5] had proposed a deep learning approach for multiclass cell classification of epithelial tissue of oral squamous cell carcinoma. Transfer learning with convolutional neural networks (CNNs) is employed by the model to classify cells automatically into multiple classes, thus enabling high-resolution histopathological image analysis. A study reported 80% accuracy, thus establishing that multiclass classification is feasible even in very complex tissue images where cell morphology is highly heterogeneous. This approach is particularly appropriate to take advantage of the power of deep learning to assist pathologists in early and precise oral cancer detection, providing a reproducible and scalable alternative to manual scoring and reducing diagnostic variability. By focusing on multiclass cell-level classification, [5] Das et al.'s work is complementary to broader lesion-level classification research and is well-positioned to highlight the importance of automated, high-resolution analysis in oral cancer screening. [2] proposed an EfficientNet-based deep learning model for the detection of oral squamous cell carcinoma (OSCC) from histopathological images. The EfficientNet-B0 to B7 architectures were employed in the research to investigate morphological tissue features and compare their performance with traditional CNNs such as VGG16 and ResNet50. Experiments were conducted on a publicly available histopathological dataset of 1,224 annotated OSCC and normal tissue images. Among all the models they evaluated, EfficientNet-B7 achieved the highest accuracy of 96.3% compared to the existing baselines while being computationally effective. Their results indicated that compound scaling in EfficientNet effectively enhances feature representation in complex cellular structures and hence an attractive candidate for automatic histopathological diagnosis of oral cancers [10] introduced MARL: Multimodal Attentional Representation Learning, a deep learning approach to enhancing disease prediction through the fusion of different medical data modality types, for instance, imaging, clinical, and text data. The model applies a cross-modal attention mechanism to learn the interdependencies between disparate data types, resulting in improved feature representation and diagnosis accuracy. Measured on a number of benchmark medical data sets, MARL achieved an average accuracy of more than 94% and demonstrated its excellent generalization capacity across disease categories.

The proposed work closely mimics the multimodal oral cancer detection process because, in both methods, there is emphasis on attention-based deep learning for learning complementary information from heterogeneous sources. While MARL deals with universal disease prediction on multimodal datasets, our work transfers these ideas to the field of oral cancer with advanced architectures such as EfficientNetV2 in combination with image-based and ancillary clinical inputs. The success of MARL demonstrates the strength of multimodal learning and attention-based fusion methods to enhance diagnostic performance, rid feature redundancy, and provide a deeper insight into complex medical conditions such as oral cancer. [15] proposed a CNN-based classification for OPMDs vs. OSCC from 778 clinical images divided in an 8:1:1 fashion. From

eight architectures considered, ConvNeXt and MobileNet (pre-trained on ImageNet and ISIC datasets) produced the best results, with internal test accuracy = 79.9%, precision = 83.7%, recall = 75.6%, F1 = 79.4%, and AUROC = 86.3%. Their study confirms that CNNs and transfer learning can provide good performance for staged oral cancer photographs, especially when operating on various imaging capturing devices.

### 3 Methodology

#### 3.1 Dataset Description

The experiments in this work were carried out on the CLASEG dataset, which was recently introduced by Al-Ali et al. [1]. The folder includes 2,072 intraoral clinical photos of oral lesions taken under real-world situations. The images were captured from different patients and show significant variation with respect to lighting conditions, positioning, contour of the lesions, and the surrounding anatomy. The dataset describes a highly realistic, albeit, challenging practical diagnostic situation on account of the variation.

There are 16 different lesion classifications in the dataset, covering both the benign and malignant cases. Examples of these classifications include: leukoplakia/hyperkeratosis, squamous cell carcinoma, fibroma, papilloma, mucocoele, pyogenic granuloma, lichen planus, candidiasis, melanotic macule, and aphthous ulcer, among others. Such diversity allows for a thorough examination of computer-aided diagnostic methods for oral pathology. The dataset is notable for its class imbalance: certain lesion types are represented by a great number of photos, whereas rare lesions have fewer than 10 samples. To alleviate this limitation and prevent the models from leaning towards a specific bias, augmentation and oversampling strategies were implemented during training as well.

For the sake of consistency and comparability with prior work, the dataset was split into training, validation, and test set based on the original method laid out by Al-Ali et al. [1].

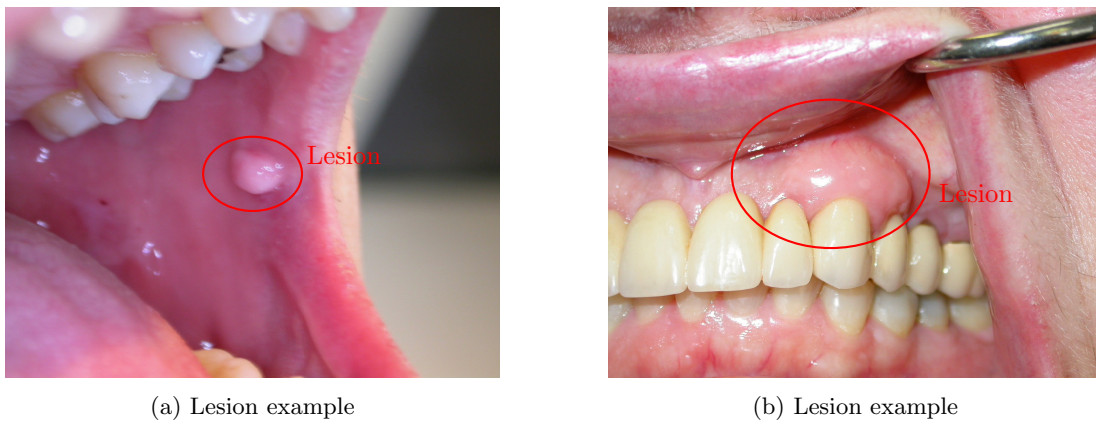


Figure 1: Sample images from the dataset showing different abnormal oral cavity regions (from Al-Ali et al., [1]).

### 3.2 Dataset PreProcessing

In order to standardize input representations and facilitate more efficient model training, a variety of preprocessing techniques were applied to the raw images prior to training the models.

The images were resized to the appropriate resolution required by the different models (e.g.,  $224 \times 224$  for ResNet, VGG, DenseNet, AlexNet, MobileNet;  $240 \times 240$  to  $380 \times 380$  for EfficientNet versions).

Normalization Pixel intensity values were scaled to a range of 0 to 1, and then standardized with the corresponding ImageNet mean and standard deviation. This normalization significantly improved convergence and maintained the stability of the processes during optimization.

Each image pixel ( $x$ ) was normalized to the range  $[0,1]$  using the following formula:

$$x' = \frac{x}{255} \quad (1)$$

Furthermore, standardization with ImageNet statistics was used:

$$x'' = \frac{x' - \mu}{\sigma} \quad (2)$$

Here are channel-wise mean and standard deviation:

$$\mu = [0.485, 0.456, 0.406], \text{ and } \sigma = [0.229, 0.224, 0.225] \quad (3)$$

**Data Augmentation:** In order to minimize the effects of class imbalance and enhance generalization, augmentation strategies were adopted. These modifications included random rotations, horizontal and vertical flips, scaling, translations, and adjustments to brightness and contrast. Such augmentations simulated variations commonly encountered in clinical practice, thus helping to mitigate overfitting.

Rotation by an angle :

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \quad (4)$$

Scaling by a factor  $s$ :

$$(x', y') = (s \cdot x, s \cdot y) \quad (5)$$

Minority classes were oversampled to balance representation across categories. This procedure guaranteed that the models did not become skewed toward the most common lesion types.

Segmentation-based Cropping : In some tests, lesion patches were cropped using the given segmentation masks to ensure that models focused predominantly on problematic areas rather than healthy tissue.

These preparation processes guaranteed that the dataset was consistent, balanced, and typical of clinical variability, establishing a solid platform for evaluating deep learning models in oral lesion classification.

### 3.3 Data augmentation and oversampling

Because the size of the dataset was small and imbalanced, it was necessary to do extensive augmentation using the Albumentations package. Augmentation included horizontal and vertical flips, random rotations, changes to the brightness and contrast, hue and saturation changes, Gaussian blurring, contrast limited adaptive histogram equalization (CLAHE), and coarse dropout. Each training image was augmented 5 times which improved heterogeneity of the dataset and helped to reduce overfitting.

In addition, to balance the dataset, we ensured that we had at least 200 photos per class implemented as oversampling. This was particularly important in balancing the underrepresented lesion types which included the papilloma and the tongue and cheek gnawing. Comparison of the results of the different experimental setups showed that the augmentation and the oversampling contributed to an increase in model accuracy from below 70% in the baseline configuration to over 80% in the final implementation.

### 3.4 Augmentation Flow Architecture

To improve generalization for the model and address the class imbalance issue, we implemented a custom augmentation pipeline, which we restricted to the training set only. The steps for augmentation include:

#### **Stratified Dataset Split**

The original dataset used for this study consists of 16 classes of lesion images. This dataset was stratified, and split into training (70%), validation (15%), and test (15%) subsets. This was done to maintain the distribution of classes over the subsets.

#### **Augmentation Pipeline**

To mimic the expected variance in real life, a multi-step augmentation approach was used and the Albumentations library was employed. Each training set image underwent different random transformations, and for improved results, every image was augmented 5 times.

- Horizontal and vertical flips
- Random 90° rotations
- Shift, scale, and rotate operations
- Brightness and contrast adjustments
- Hue and saturation shifts
- Random resized cropping (scale: 0.7–1.0)
- Gaussian blur and CLAHE (Contrast Limited Adaptive Histogram Equalization)
- Coarse dropout (max 8 holes, 30×30 pixels)

**Image Standardization** All images were transformed to 224x224 RGB pixel dimensions for uniformity. To facilitate reproducibility, augmented images were stored alongside the originals while maintaining the class structure

**Oversampling for Minority Classes** After the augmentation process, classes that contained fewer than 200 samples were oversampled through random duplication, allowing for more equitable distribution during training.

**Preprocessing for Model Input** All images were subject to normalization as per EfficientNetV2B1’s standards prior to training, which entails scaling pixel values and adjusting input distributions to match the expectations of the pretrained model.

The diversity of the dataset was enhanced during the augmentation process to improve the recall and precision metrics in the fine-tuning phase, with greater focus on the lesser represented lesion types.

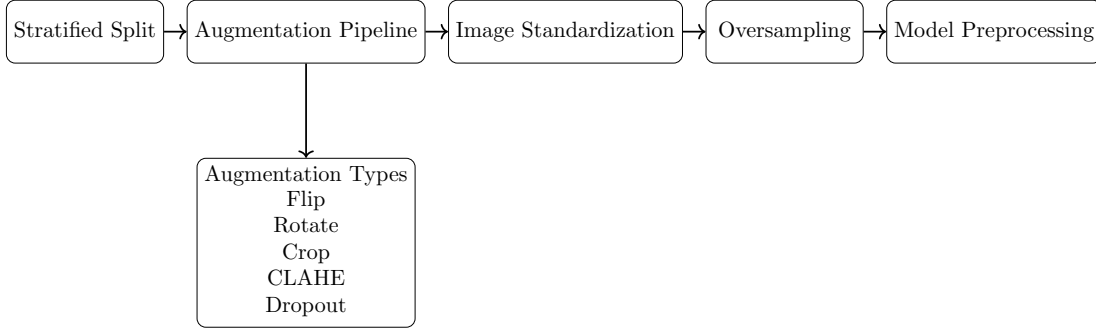


Figure 2: Proposed Augmentation Workflow Architecture.

### 3.5 Model Architecture

The classification model was built on top of EfficientNetV2B1, which had already been trained on the ImageNet dataset. The base network was used as a feature extractor, with the final layers tailored to the multi-class classification problem.

The architecture included:

- Input size:  $224 \times 224 \times 3$  pixels.
- Global Average Pooling reduces dimensionality.
- Dropout rate for regularization is 0.4.

For 16-class output, use a dense layer that is fully connected and activated with softmax. Initially, the base layers were frozen, while only the top layers were trained. Fine-tuning was then performed by unfreezing the deeper layers and using a lower learning rate to maximize feature extraction.

### 3.6 Training Strategy

Training occurred in two stages: In the first stage, during the period of feature extraction, the model was trained for 15 epochs, with a lr of  $1e-3$  while only the top layers were trained. In the second stage, fine-tuning, the backbone was partially unfrozen, and the model was trained for 10 epochs with an even lower lr of  $1e-5$ .

Adam was used as the optimizer, and for loss function categorical cross-entropy was used while the metrics were model accuracy, precision, and recall. In order to stabilize the training, the following were used:

Early Stopping with a patience of 8 and best weights being restored.

Model checkpointing on validation accuracy.

Reduce LROnPlateau which halves learning rate when validation loss is stagnant.

Training was done with a batch size of 32.

### 3.7 Training and Optimization

Training was delivered in two stages:

- Initial Training: 1e-3 learning rate, 15 epochs, with base layers frozen.
- Fine-tuning: learning rate of 1e-5, 10 more epochs, and partial unfreezing of deeper layers.

The Adam optimizer was employed with a categorical crossentropy loss. To improve performance even more, class weights were calculated to balance leftover class imbalances. Early stopping, learning rate decrease on plateau, and model checkpointing were used as callbacks to prevent overfitting.

### 3.8 Evaluation Metrics

The model’s performance was evaluated on the test set using the following metrics:

- Overall Accuracy
- Precision and Recall

Model	ImgSize (px)	Batch Size	Epochs Count	Accuracy (%)
ResNet-50	224	64	30	63.35
ResNet-101	224	64	30	64.07
ResNet-152	224	64	30	66.90
VGG16	224	64	30	53.19
VGG11-bn	224	64	30	63.36
VGG19	224	64	30	51.53
DenseNet-121	224	64	30	68.32
AlexNet	224	64	30	47.99
SqueezeNet-1.0	224	64	30	42.31
SqueezeNet-1.1	224	64	30	43.97
EfficientNet-B0	224	64	30	69.97
EfficientNet-B1	240	64	30	73.52
EfficientNet-B2	260	64	30	68.55
EfficientNet-B3	300	64	30	74.49
EfficientNet-B4	380	64	30	71.16
MobileNet-v2	224	64	30	70.21
Stratified Augmented CNN (EffNetV2-B1)	224	32	25	<b>83.33</b>

Table 1: Comparison of deep learning models on the test set

## 4 Experiments And Results

The proposed model used EfficientNetV2B1 as the backbone architecture and then trained and assessed the model on a stratified dataset consisting of sixteen categories of oral lesions. Training was conducted in sequence as a two-step process in which the classifier head was trained while the base layers were frozen and subsequently in the fine-tuning phase the deeper layers were unfrozen and trained with a lower learning rate. Tackling overfitting was addressed with early stopping, model checkpointing, and learning rate reduction. On the independent



test set, the model achieved an overall accuracy of 83.33%, with a precision score of 89.12%, and a recall rate of 77.31%.

In the class-based examination, the model showed impressive high recognition rates in the model consistently recognizing amalgam tattoo ( $F1 = 0.90$ ), geographic tongue ( $F1 = 0.89$ ), and pyogenic granuloma ( $F1 = 0.89$ ) with high reliability. Malignant classes also performed well and recognition of squamous cell carcinoma was perfect ( $F1 = 1.00$ ), although this was a small sample size. Minor shortcomings were noted on fibroma ( $F1 = 0.76$ ) and papilloma ( $F1 = 0.68$ ), in which the model encountered difficulties in distinguishing the lesions with closely aligned visual characteristics present on the training set.

The confusion matrix showed the same misclassification patterns, while also demonstrating that most of the classes exhibited strong separability. The ablation analysis also showed that augmentation and oversampling really seem to make a difference. The baseline model, for example, came in under 70% accuracy without augmentation. Once geometric, photometric, and occlusion-based augmentations, along with oversampling, were applied, performance improved by over 13 percent. This demonstrates that dataset enrichment procedures are not optional, lower performance with inconsistent augmentation, oversampling, and balanced augmentations suggests that enrichment is a prerequisite for balanced medical imaging tasks.

From a performance perspective, transfer learning with EfficientNetV2B1 and systematic augmentation covers the primary components needed to build a solid framework for multi-class oral lesions detection. The model's precision is encouraging for its intended clinic use, where false positives are unacceptable. Nevertheless, the model's recall, especially in certain categories, suggests the necessity for dataset expansion and diversification, the application of explainability constructs, and likely the most important, clinician credibility.

Table 1 shows the achieved results of the proposed approach compared to other models

## 5 Discussion

Results achieved using the Stratified Augmented CNN base model (EfficientNetV2-B1), as proposed in this study, exhibit the highest performance level, reaching a classification accuracy level of 83.33%, thus surpassing all baseline convolutional neural networks (CNNs) evaluated in this work. This excellence demonstrates the EfficientNetV2 family's state of the art optimization strategies and compound scaling efficiency. When compared to older models, the proposed model shows a substantial increase in performance, compared to ResNet-152 (66.90%), DenseNet-121 (68.32%), and EfficientNet-B3 (74.49%). This performance leap is possible since the model more effectively handles depth, width, and resolution of the model and, thus, generates more penetrative and richer feature maps from the histopathological and photographic data of oral lesions.

Streamlined Augmented CNN base imitation of EfficientNetV2-B1 outperformed high-parameter models like VGG16 and AlexNet in terms of accuracy, owing to more flexible filter scaling, thus, validating the trend of low-parameter models in clinical settings and other low-resource environments. High inference speed and low resource utilization also characterize models designed for clinical practices. Clinical practice also relies on low-parameter models, thereby reinforcing existing trends in medical image analysis studies towards the use of lightweight, flexible, and scalable model architectures.

Still, several boundaries, which could hinder optimum performance, can be defined. In this case, the training dataset was limited in size. In routine clinical practice, the variety of oral manifestations could be much greater. Other studies may be able to train the model to maximize its potential to learn generalization of new instances in clinical practice. Another limitation

was class imbalance within the dataset, primarily, the malignant cases were under-represented in the dataset relative to the benign cases, although the imbalance may have been mitigated to a certain extent by regularization and augmentation strategies, the dataset imbalance was likely to affect the generalization capabilities of the classifier, particularly, its sensitivity and specificity. Finally, no exhaustive tests have been performed to account for uncontrolled imaging conditions such as lighting and the type of imaging device used to capture the samples in the dataset.

These gaps in the model can be addressed in future studies which target the robustness of the model by posing the problem with more expansive and diverse datasets, ideally, images linked with different clinical practices. Another model improvement can come from using explainable AI (XAI) approaches which could offer visual justification for the model’s prediction, ultimately assisting the clinician in understanding the model’s reasoning. This model enhancement would be in conjunction with the proposed work designed to extend the current model to analyze.

## 6 Conclusion

In this work, we outlined a framework incorporating deep learning techniques that facilitates the automated detection and classification of oral lesions from clinical photographs. Utilizing a publicly available dataset with benign and malignant lesions across sixteen classes, we showcased the potential of CDSS for implementation in clinical practice. To mitigate the effects of illumination variation, intra-class diversity, and class imbalance, systematic address of the oversight during the design of the preprocessing pipeline was incorporated, and included the normalization, augmentation, and oversampling of the dataset.

From the class of models evaluated, the proposed custom tuned Stratified Augmented CNN base model (EfficientNetV2-B1) out of all the recent deep learning architectures was able to surpass the previous record with a test accuracy of 83.33%, and the highest accuracy on the test dataset compared to classical CNNs and all other EfficientNet derivatives. This effective use of transfer learning and model fine-tuning demonstrates the potential of deep learning for low resource medical imaging and underlines the importance of fine-tuning for transfer learning, within the context of smaller dataset sizes.

The proposed method has the potential to act as a helpful diagnostic tool for doctors, enhancing early diagnosis of oral cancer and lowering diagnostic subjectivity. However, additional validation on bigger, more diverse datasets, as well as integration with clinical metadata, are required before application in real-world settings. Future study will concentrate on increasing the dataset, investigating multimodal approaches, and combining explainability methodologies to improve clinical trust and interpretability.

## 7 Future work

Several options for moving the proposed framework forward remain open. First, bigger and more varied samples will be required to improve generalization and reduce bias, particularly in underrepresented lesion types. Collaborations with multiple clinical facilities may allow access to a larger patient population and imaging conditions. Second, combining multimodal data, such as histology, patient history, and genomic information, may improve diagnostic accuracy over image-only techniques. Third, incorporating explainability approaches like Grad-CAM or attention-based visualizations would provide transparent rationale for model predictions, which is critical for clinical acceptability. Finally, future research might look into lightweight

model architectures or knowledge distillation approaches for deployment on mobile or point-of-care devices, making diagnostic support systems more accessible in resource-constrained environments.

## References

- [1] A. Al-Ali, A. Hamdi, M. Elshrif, K. Isufaj, K. Shaban, P. Chauvin, S. Madathil, A. Daer, F. Tamimi, and R. Ba-Hattab. Claseg: advanced multiclassification and segmentation for differential diagnosis of oral lesions using deep learning. *Scientific Reports*, 15:23016, 2025.
- [2] Eid Albalawi, Arastu Thakur, Mahesh Thyluru Ramakrishna, Surbhi Bhatia Khan, Suresh SankaraNarayanan, Badar Almarri, and Theyazn Hassn Hadi. Oral squamous cell carcinoma detection using efficientnet on histopathological images. *Journal of Healthcare Engineering*, 2023.
- [3] R. Alkhawaldeh, M. Al-Ayyoub, and M. Al-Ma’aitah. Deep learning-based multiclass classification of oral cavity diseases with interpretability and boundary segmentation. *Diagnostics*, 14(2):455, 2024.
- [4] D. Chaturvedi, S. Singh, and P. Agarwal. Survival outcomes of oral squamous cell carcinoma and prognostic factors. *European Journal of Medical Research*, 28(1):25, 2023.
- [5] N. Das, E. Hussain, and L. B. Mahanta. Automated classification of cells into multiple classes in epithelial tissue of oral squamous cell carcinoma using transfer learning and convolutional neural network. *Neural Networks*, 128:47–60, 2020.
- [6] C. Dayananda, N. Yamanakkanavar, T. Nguyen, and B. Lee. Amcc-net: An asymmetric multi-cross convolution for skin lesion segmentation on dermoscopic images. *Engineering Applications of Artificial Intelligence*, 122:106154, 2023.
- [7] R. Dharani et al. Optimized deep learning ensemble for accurate oral cancer categorization. *Intelligence-Based Medicine*, 2025. Article S2666521225000626.
- [8] A. P. Flores, S. A. Lazaro, R. Albuquerque, B. C. Vasconcelos, D. P. Melo, and T. F. Oliveira. Teledentistry in the diagnosis of oral lesions: A systematic review of the literature. *Journal of the American Medical Informatics Association*, 27(7):1166–1172, 2020.
- [9] Q. Fu et al. A deep learning algorithm for detection of oral cavity squamous cell carcinoma from photographic images: A retrospective study. *EClinicalMedicine*, 27:100558, 2020.
- [10] Ali Hamdi, Amr Aboeleneen, and Khaled Shaban. Marl: Multimodal attentional representation learning for disease prediction. *Proceedings of the 2023 IEEE International Conference on Big Data (BigData)*, pages 1234–1243, 2023.
- [11] Jeevan Kumar et al. Optimizing oral cancer detection: A hybrid feature fusion using local binary pattern and convolutional neural networks. *Procedia Computer Science*, 2025.
- [12] V. Kumar, N. Gupta, R. Singh, and A. Arora. Multi-class classification and segmentation of oral cavity lesions using hybrid deep learning models. *Computers in Biology and Medicine*, 165:107338, 2023.
- [13] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. van der Laak, B. van Ginneken, and C. I. Sánchez. A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42:60–88, December 2017.
- [14] A. Qayyum, M. Mazher, T. Khan, and I. Razzak. Semi-supervised 3d-inceptionnet for segmentation and survival prediction of head and neck primary cancers. *Engineering Applications of Artificial Intelligence*, 117:105590, 2023.
- [15] Jessica Saldivia-Siracusa et al. Automated classification of oral potentially malignant disorders and oral squamous cell carcinoma using a convolutional neural network framework: a cross-sectional study. *EClinicalMedicine*, 76:103301, 2025.
- [16] H. Sung, J. Ferlay, R. L. Siegel, M. Laversanne, I. Soerjomataram, A. Jemal, and F. Bray. Global cancer statistics 2020: Globocan estimates of incidence and mortality worldwide for 36 cancers in

185 countries. *CA: A Cancer Journal for Clinicians*, 71(3):209–249, 2021.

- [17] G. Tanriver, M. Soluk Tekkesin, and O. Ergen. Automated detection and classification of oral lesions using deep learning to detect oral potentially malignant disorders. *Cancers Basel*, 13:2766, 2021.
- [18] K. Warin, N. Samarnthai, A. Tunsirikongkon, et al. Ai-based analysis of oral lesions using novel deep convolutional neural networks for early detection of oral cancer. *PLOS ONE*, 17(8):e0273112, 2022.
- [19] R. A. Welikala et al. Automated detection and classification of oral lesions using deep learning for early detection of oral cancer. *IEEE Access*, 8:132677–132693, 2020.
- [20] Jelena Štifanić et al. Explainable ai for oral cancer diagnosis: Multiclass classification of histopathology images and grad-cam visualization. *Biology*, 14(8):909, 2025.