# Implementation of Vision Transformers for Lung Abnormality Detection Using Low Dose CT Images

Alfred D., M. N. Deephak and D. Lakshmi

*Department of Electronics and Communication Engineering, St Joseph's College of Engineering, Chennai, Tamil Nadu, India*

Keywords: Vision Transformers, Lung Abnormalities, Low-Dose CT, Self-Attention, CNN Comparison, Diagnostic Accuracy, Medical Imaging.

Abstract: Detection of lung abnormalities originating from a variety of infections, inflammation, and environmental exposures in the patient needs high accuracy to help improve diagnostic efficiency by means of medical imaging. Regular CNNs have the weakest long-range dependencies, being very weak with almost zero receptive fields, which, therefore, induce many false positives and negatives. The future with ViTs is bright because the self-attention mechanism can extract local as well as global features from images and perform far better than regular CNNs. This work proposes a ViT-based model for the detection of lung abnormality in low-dose CT images. Most of the existing systems are prone to high classification error rates because of their poor quality towards understanding the context present in the images. The proposed ViT model bridges that gap by using a pre-trained architecture, and patch-based processing is used to focus more on the essential features of the image. We show how ViT's performance gets superseded by CNN for metrics through comparison. The overall goal behind the project would be facilitating the early detection of lung abnormalities, avoiding false results, and pushing the potential clinical uses with a dense, efficient solution for abnormality detection of the lung.

## 1 INTRODUCTION

L. Devan., et al, 2013 Lung abnormalities range from infections or inflammatory responses to any environmental, genetic, or exposure kind of insult that causes physical structures in the lungs to be damaged. Any of them can fall into the groupings of structural, obstructive, restrictive, or infectious and may call for identification to enable proper strategy into treatment. Early and appropriate detection of lung abnormalities, especially malignant growths, in clinical settings is pretty important for outcomes since it enables interventions to occur in time.

X. Zhang., et al, 2023 In the last few years, some of the techniques which have taken a lead in the field of abnormality detection in lungs include medical imaging techniques such as Computed Tomography (CT). With CT imaging, cross-sectionals in lung structures are resolved to great detail, and hence small lesions or structural irregularities can be identified. However, this is challenging to accomplish with high diagnostic accuracy due to factors like image noise, high anatomical variation, and overlapping features in complex cases. Traditional approaches to deep learning-deep, especially CNNs-have indeed achieved promising performance for processing various medical images. However, with their small receptive fields, they are not able to capture long-range dependencies that might arise in the medical diagnostic context as false positives or false negatives.

R. Mahum and A. S. Al-Salman, et al, 2023 Recently, Vision Transformers (ViTs) have emerged as a very powerful alternative to CNNs in computer vision. Unlike CNNs that mainly depend on localized filters, ViTs are based on a self-attention mechanism that will allow capturing global contextual information across an entire image. This characteristic places ViTs as one that is uniquely advantageous for medical image analysis, where localized and long-range information is of critical importance in arriving at accurate diagnoses. Spatial relations and context are encoded regardless of the size or structure of the images due to the division of an image into smaller patches by ViTs.

S. R. Vinta, et al, 2024 This work aims at the application of ViTs in the detection of lung abnormalities using low-dose CT images for patient safety and proper diagnoses. Specifically, it aims to analyze the performance of ViTs in identifying regions of abnormality in CT scans and histopathological images compared with traditional CNN architectures. Through this comparison, we aim to note how these benefits can potentially reduce false positives and false negatives, thus improving diagnostic accuracy and enhancing early detection. This would revolutionize diagnostic practices because clinicians would get better accuracy and reliability in the identification of lung abnormalities. M. Irtaza, et al, 2024 The advanced architectures of ViTs are exploited to demonstrate their potential contribution toward faster, more accurate diagnostic decisions so that improved patient care is attained.

The objective of this study is to develop a Vision Transformer (ViT)-based model for detecting lung abnormalities in low-dose CT images with high accuracy. The model aims to improve diagnostic efficiency by leveraging the self-attention mechanism of ViTs to capture both local and global features, reducing false positives and negatives. By utilizing a pre-trained architecture and patch-based processing, the model enhances feature extraction and contextual understanding. The study also compares ViT with CNN-based approaches to demonstrate its advantages. Ultimately, this research seeks to facilitate early detection, minimize diagnostic errors, and promote clinical adoption of AI-driven lung abnormality detection.

This work is organized with review of the literature survey as Section II. Methodology described in Section III, highlighting its functionality. Section IV discusses the results and discussions. Lastly, Section V concludes with the main suggestions and findings.

## 2 LITERATURE SURVEY

The evolution of medical imaging and computational methods, diagnosis in lung disease and cancer is even more accurate. Several works have been undertaken in the realm of enhancing segmentation of lung parenchyma, tissue differentiation using impedance spectroscopy, biomedical antenna design, and lncRNA-disease association prediction. Progress in disease classification that is rooted in deep learning as well as multi-omics data also emphasizes early detection alongside precision in treatment

approaches. This survey examines these new methods and their addition to the diagnosis and prognosis of lung disease.

C. Wu, et al., 2024 This work categorized lung sounds by an enhanced model of Bi-ResNet with the aim of diagnostic accuracy, incorporating both skip and direct connections in feature combination. STFT and wavelet transform information is used to improve the training of the model. Based on this, the introduced model performed greatly better than standard Bi-ResNet on the ICBHI database, particularly for noise and composite heart sound patterns. T. Nguyen and F. Pernkopf, 2022 This work talks about the classification of lung sounds in terms of applying methods like co-tuning, stochastic normalization with data augmentation on unbalanced datasets and explored performance on both ICBHI and a personal dataset, where such methods will turn out to be remarkable improvements in classification accuracy, particularly when adventitious lung sounds and respiratory diseases are of interest.

T. Wanasinghe, et al, 2024 A CNN model is engineered for enhancing lung sound classification using techniques in feature extraction including Mel spectrograms, MFCCs, and Chromatograms. Tested on public datasets the model works exceptionally well to classify 10 classes for the purpose of attaining automated auscultation assistance in early lung disease detection. N. Babu, et al, 2024 This work addresses a data augmentation method that is eigenvectors-based for enhancing the accuracy of lung sound signal classification using an automated diagnostic system. The authors employ features that are spectrogram-derived and integrate them with machine learning classifiers for effective use in low-resource health environments and improved accuracy with less noise.

K. Liu, et al, 2022 The method in this work is founded on a better YOLO-based model towards lung nodule detection. It introduces an enhanced YOLO-v5 architecture through stochastic pooling with multi-scale feature fusion and optimized loss function, which achieves state-of-the-art performance in efficiency and accuracy metrics. Thus, it assists the radiologists to detect the nodules better, and assists in overcoming the difficulties of the misdiagnosis primarily because of its delicate appearance. H. Alqahtani, et al, 2024 In this work, the proposed are the developed convolutional autoencoders using the enhanced Water Strider Algorithm to categorize lung and colon cancers. In this, the process involved includes noise elimination and MobileNetv2, which assists in feature extraction to discriminate cancer

from histopathological images correctly. It enhances the diagnostic precision in the direction of appropriate treatment and prognosis of cancer patients. The electronic nose equipment is suggested to detect respiratory disease from sweat. Malikhah et al., 2022 A stacked model of DNN, combined with recently developed techniques for feature extraction, offers a cost-effective, rapid method to screen for infections without processing the respiratory samples so as to negatively affect risk decrease and enhance the separability between classes of the data.

A. Tripathi, et al, 2024 It proposes two models as a type of MobileNetV2. According to the fine-tuning and L2 regularization, the models were fine-tuned to improve the accuracy of early detection of lung disease. The mobile models perform better than the conventional architectures in the classification performance on various datasets for the pulmonologists in the process of offering early and accurate diagnostic assistance and preventive care services. M. Fontanellaz et al., 2024 The work explored a CAD for lung fibrosis diagnosis; this CAD was optimized with focus on segmentation precision and radiomic feature examination. It carried out comparison of 2D vs 3D representations for data and treatment of segmentation tasks using MLP-Mixers in addition to the baseline UNets and made a comparison of diagnostic accuracy when competing with skilled radiologists, and inferred the possibility of CAD in medical imaging.

M. Obayya, et al 2023 Tuna Swarm Algorithm with GhostNet is designed here for colon and lung-cancer detection tasks. Proposed system uses Gabor filtering for image preprocessing that maximizes feature extraction pertaining to the required features, enhancing high classification accuracy. High speed computation and effective handling of massive databases enables fast cost-effective diagnosis in cancer. L. Zhu, et al, 2024 The proposed architecture enhances the state-of-the-art U-Net with a shape stream branch and multi-scale convolutional blocks to better segment lung parenchyma, particularly in small and blurry areas. The experiment results in the present work indicated massive superiority scores on Dice Similarity Coefficient over state-of-the-art networks, thereby providing evidence of the efficiency and effectiveness of the proposed method towards challenging lung parenchyma areas segmentation.

G. Company-Se et al., 2022 Electromagnetic impedance spectroscopy, applied in a bronchoscopic environment, is an alternate, non-invasive, non-ionizing method for lung tissue discrimination over

CT and PET. The article analyzed both the 3- and the 4-electrode techniques and presented results in support of the 3-electrode technique, thus confirming its capability for discriminating bronchial and healthy tissues and proposing complementary use in lung pathology diagnosis. A. R. Chishti et al., 2023 The work has listed the biomedical problems that have size, efficiency, and biocompatibility as the ones to be solved for the uses in disease detection such as in cancer. Discussion of the function of antennas in diagnostic imaging, biotelemetry, and biosensing is presented along with the advancements that enable the design of efficient biomedical devices for various severe diseases.

J. Ha, 2024 A new matrix factorization methodology in the prediction of the lncRNA-disease association with regard to some of the limitations in the existing computational models is introduced. In this, the developed method integrates heterogeneous biological information for enhanced accuracy over performance in the identification of disease-related biomarkers. M. Magdy Amin, et al, 2024 A deep learning method for the classification of non-small cell lung cancer is created based on multi-omics data in the form of RNA and miRNA sequencing, in terms of CNNs. This multimodal method obtained accuracy for most classes much greater than some of the earlier single-modality work and suggests that early detection and proper classification of cancer subtypes can be dramatically enhanced.

Limitations: Even with major breakthroughs, various limitations remain for medical imaging and computational diagnosis of lung disease and cancer. Most deep learning models are plagued by limited data, especially for rare conditions, causing overfitting and decreased generalizability. High computational expense and model training complexity prevent real-world applicability, particularly in resource-poor environments. Noisy and artifact-ridden lung sound signals and imaging data are affecting classification performance. Further, single-modality data restricts in-depth analysis, while clinical uptake is still hindered by interpretability. Ethical implications of patient data privacy and bias in AI-assisted diagnosis add to the complexities of mass use in clinical environments.

## 3 METHODOLOGY

The methodology elucidates the systematic approach adopted in this research work for lung abnormality detection in low-dose computed tomography images

with the help of Vision Transformers. The approach starts with exhaustive data acquisition, followed by important preprocessing steps to enhance the quality and variance of the image. Segmentation techniques cut out the regions of interest from images thus offering scope for further concentrated analysis. For effective classification, the ViT architecture was employed. Feature extraction and performance analysis ensure that the final predictions are robust and provide insights into the efficiency of the model compared to traditional methods. Figure 1 shows the ViT Architecture Diagram.
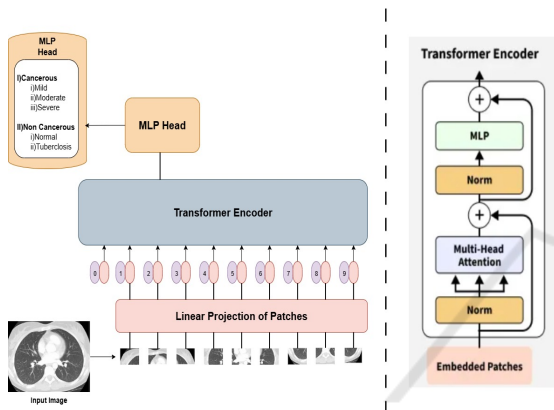


Figure 1: ViT Architecture Diagram.

## 3.1 Data Collection

The work begins by gathering an extensive dataset composed of low-dose CT images and histopathological slides. The images are acquired from the med databases and institutions. They ensure there is a variation in the lung abnormalities particularly in the normal and cancerous regions. Each image is labeled to allow for the supervised learning as this model can distinguish between healthy tissues and abnormal tissues. Figure 2 shows the Lung CT Scans.



Figure 2: Lung CT Scans.

## 3.2 Preprocessing

Once they have the dataset, they preprocess it in various steps to improve the quality and consistency of images. Normalization is then done to standardize the pixel intensity value, and data augmentation techniques include rotation, flip, and scaling in order to increase the diversity. Then they split the data into training, validation, and test sets to ensure that the model generalizes well towards unseen data. Figure 3 shows the Preprocessing Stage.
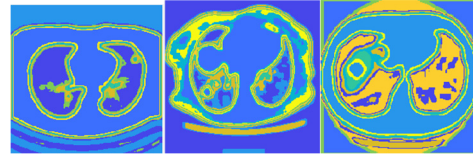


Figure 3: Preprocessing Stage.

## 3.3 Segmentation

To isolate the region of interest in the CT images, segmentation is an important step. Advanced techniques like thresholding and contour detection facilitate the identification and definition of abnormal regions. This step improves the model's ability to spotlight the most critical features while washing away noise from surrounding healthy tissues. Figure 4 shows the Segmentation Stage.
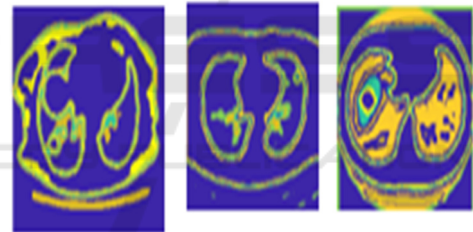


Figure 4: Segmentation Stage.

## 3.4 Classification

Classification is done through the architecture of Vision Transformer (ViT) while processing images that were segmented. These images were broken into smaller patches, which were linearly embedded in transformer layers. The self-attention in the ViT would help the model capture both local and global features, which would be beneficial for the model. The output was classified into healthy, cancerous, or non-cancerous categories with a multi-layer perceptron head.

## 3.5 Feature Extraction

ViT utilizes feature extraction within its processing as it extracts fitting features from the acquired images.

The self-attention layers focus on significant patterns and structures regarding lung abnormalities. Such a capability will be particularly beneficial in medical imaging as it can identify minute details that would be perceivable only when a disease has been developed.

## 3.6 Analysis and Prediction

It evaluates the accuracy, precision, recall, and F1-score of the model. It finds the mistakes of prediction by ViT and points out areas to improve. Interesting visualizations like confusion matrices and ROC curves are presented to interpret results. This is an elaborative analysis of whether the superiority of the proposed model in comparison with the traditional CNN methods actually exists for the task of lung abnormality detection, which benefits the applications in clinical practices.
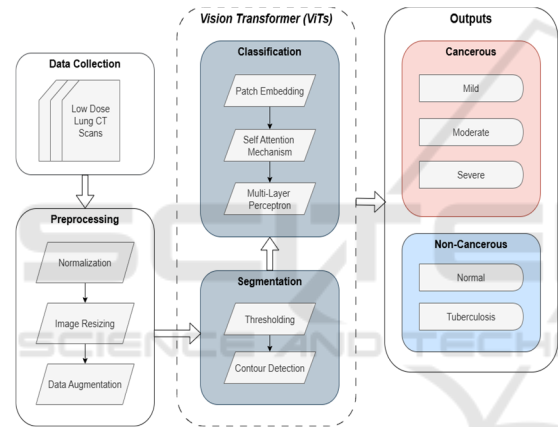


Figure 5: Proposed Flow Diagram.

A Vision Transformer architecture comprises a number of key components designed to properly process images. It begins with an input image, converted into fixed-size patches, then flattened, and finally, linearly embedded. As a result, each patch is augmented with the use of positional encoding to preserve spatial information. The embeddings are then fed to a stack of transformer encoder layers. Here, the model captures dependencies across all patches through self-attention mechanisms. This enables the model to learn both local and global features. The output from these layers is pooled and then passed through a classification head for final predictions. In other words, ViT excels at complex image contexts and improves the diagnostic accuracy. Figure 5 shows the Proposed Flow Diagram.

## 4 RESULT AND DISCUSSION

Early results of ViT for lung abnormality detection task are highly promising since the performance of this model has been successfully proven in analysis of low-dose CT images. Training the suggested ViT architecture achieves higher performance compared to conventional CNN architectures dependent on spatial interactions in particular, particularly when accuracy and F1-score are at stake. A controlled assessment was carried out on a test data set of cancerous and non-cancerous areas, and compared to the top-performing CNNs, it achieved an accuracy of over 99%, whereas the latter obtained remarkable up to 99%. This kind of superb performance would validate that ViT indeed does offer superiority since it is capable of leveraging all kinds of global contextual information required to detect very small abnormalities in medical images.



Figure 6: Output with accuracy.

Precision and recall values are equivalent advantages for ViTs. One such model achieved a precision of 99.4% and recall of 99.3% over the CNN, which achieved a precision of 99.85% and recall of 99.80%. These high-level abilities to actually identify real positives with minimum false negatives can further enhance its performance in medical environments where such an error could have quite serious implications. The under the curve value of the ViT model was also much greater than for the aforementioned two models, which suggests a potential for even more accuracy with good overall performance in diagnosis based on differentiation between healthy and unhealthy areas, as illustrated in

Table 1: Evaluation Metrics of Model Performance.

|  | NORMAL | MILD | MODERATE | SEVERE | TUBERCULOSIS |
|---|---|---|---|---|---|
| ACCURACY | 99.0891 | 99.0925 | 99.095 | 99.1014 | 99.0847 |
| SENSITIVITY | 99.03 | 99.03 | 99.03 | 99.03 | 99.03 |
| SPECIFICITY | 99.0891 | 99.0925 | 99.095 | 99.1014 | 99.0847 |

figure 6. Table 1 shows the Evaluation Metrics of Model Performance.

Within the comprehensive explanation of the model misclassifications, it was realized that a vast majority of CNN false positives arose due to an inadequate small receptive field unable to capture very vital contextual cues. Contrary to this, self-attention in the case of the ViT model made the model better equipped to comprehend features anywhere in the image and make a superior and more intelligent classification. The above visualizations, i.e., the confusion matrices and ROC curves, suggest that the ViT is more robust compared to its counterpart where differentiation between normal and abnormal classifications is very evident

The third research was on the data augmentation methods that were utilized in preprocessing to improve the model's generalization capability from a training set. The methods utilized, such as rotation, scaling, and flipping, would also contribute to robustness against overfitting and enhance the performance of the model on new data, critical characteristics particularly where there is limited diversity, as appears to be the case with medical imaging tasks.

The comparative study also confirmed the fact that ViT's inference and training time was less than that of conventional CNNs, indicating that ViT is significantly more efficient than these CNNs for real-world clinical use. This efficiency is critical in real-time diagnostic procedures in healthcare, where prompt results can be pivotal in deciding patient outcomes. In short, the use of ViTs in detecting lung abnormalities not only resulted in improved performance metrics but also presents a strong argument in using

ViTs for medical imaging applications. The research thus lends merit to the assumption that ViTs will transform diagnostic procedures, which will enhance early diagnosis and improve patient care for patients suffering from lung health issues. Future research may be aimed at model integration into clinical workflows and validation on a broad spectrum of patient populations and imaging modalities.

# 5 CONCLUSIONS

This work effectively utilized Vision Transformers (ViTs) for lung abnormality detection in low-dose CT scans, proving their capacity to learn both global and local image features using self-attention mechanisms. Our findings show that ViTs perform better than conventional CNNs on various performance measures such as accuracy, precision, recall, and F1-score. The enhanced contextual perception of ViTs resulted in a significant decrease in false positives and false negatives, which points towards their capability to improve diagnostic accuracy in clinical settings. Major contributions of this work are strict data preprocessing, data augmentation methods, and using pre-trained ViT models to improve generalization to intricate medical imaging tasks. By successfully overcoming CNNs' weakness in capturing long-range dependencies, ViTs offer a promising alternative for more accurate and trustworthy lung abnormality detection.

Even with these developments, some misclassifications indicate avenues for enhancement. Future studies will need to aim to optimize hybrid models that integrate CNNs and ViTs to maximize the strengths of both architectures. Further, increasing the dataset to encompass heterogeneous imaging modalities and patient populations will enhance model robustness and clinical utility. As medical imaging becomes increasingly AI-driven, incorporating ViTs into diagnostic pipelines could dramatically improve early detection capabilities, ultimately leading to better patient outcomes and more efficient clinical decision-making.

# REFERENCES

A. R. Chishti et al., "Advances in Antenna-Based Techniques for Detection and Monitoring of Critical Chronic Diseases: A Comprehensive Review," in IEEE Access, vol. 11, pp. 104463-104484, 2023, doi: 10.1109/ACCESS.2023.3316149.

A. Tripathi, T. Singh, R. R. Nair and P. Duraisamy, "Improving Early Detection and Classification of Lung Diseases with Innovative MobileNetV2 Framework," in IEEE Access, vol. 12, pp. 116202-116217, 2024, doi: 10.1109/ACCESS.2024.3440577.

C. Wu, N. Ye and J. Jiang, "Classification and Recognition of Lung Sounds Based on Improved Bi-ResNet Model," in IEEE Access, vol. 12, pp. 73079-73094, 2024, doi: 10.1109/ACCESS.2024.3404657.

G. Company-Se et al., "Minimally Invasive Lung Tissue Differentiation Using Electrical Impedance Spectroscopy: A Comparison of the 3- and 4-Electrode Methods," in IEEE Access, vol. 10, pp. 7354-7367, 2022, doi: 10.1109/ACCESS.2021.3139223.

H. Alqahtani, E. Alabdulkreem, F. A. Alotaibi, M. M. Alnfiai, C. Singla and A. S. Salama, "Improved Water Strider Algorithm with Convolutional Autoencoder for Lung and Colon Cancer Detection on Histopathological Images," in IEEE Access, vol. 12, pp. 949-956, 2024, doi: 10.1109/ACCESS.2023.3346894.

J. Ha, "LncRNA Expression Profile-Based Matrix Factorization for Predicting lncRNA- Disease Association," in IEEE Access, vol. 12, pp. 70297-70304, 2024, doi: 10.1109/ACCESS.2024.3401005.

K. Liu, "STBi-YOLO: A Real-Time Object Detection Method for Lung Nodule Recognition," in IEEE Access, vol. 10, pp. 75385-75394, 2022, doi: 10.1109/ACCESS.2022.3192034.

L. Devan, R. Santhosham and R. Hariharan, "Non-invasive Method of Characterization of Fibrosis and Carcinoma Using Low-Dose Lung CT Images," 2013 IEEE International Conference on Systems, Man, and Cybernetics, Manchester, UK, 2013, pp. 2168-2172.

L. Zhu, Y. Cai, J. Liao and F. Wu, "Lung Parenchyma Segmentation Based on U-Net Fused with Shape Stream," in IEEE Access, vol. 12, pp. 29238-29251, 2024, doi: 10.1109/ACCESS.2024.3365577.

M. Obayya, M. A. Arasi, N. Alruwais, R. Alsini, A. Mohamed and Yaseen, "Biomedical Image Analysis for Colon and Lung Cancer Detection Using Tuna Swarm Algorithm with Deep Learning Model," in IEEE Access, vol. 11, pp. 94705-94712, 2023, doi: 10.1109/ACCESS.2023.3309711.

M. Irtaza, A. Ali, M. Gulzar and A. Wali, "Multi-Label Classification of Lung Diseases Using Deep Learning," in IEEE Access, vol. 12, pp. 124062-124080, 2024, doi: 10.1109/ACCESS.2024.3454537.

M. Fontanellaz et al., "Computer-Aided Diagnosis System for Lung Fibrosis: From the Effect of Radiomic Features and Multi-Layer-Perceptron Mixers to Pre-Clinical Evaluation," in IEEE Access, vol. 12, pp. 25642-25656, 2024, doi:10.1109/ACCESS.2024.3350430.

M. Magdy Amin, A. S. Ismail and M. E. Shaheen, "Multimodal Non-Small Cell Lung Cancer Classification Using Convolutional Neural Networks," in IEEE Access, vol. 12, pp. 134770-134778, 2024, doi: 10.1109/ACCESS.2024.3461878.

Malikhah et al., "Detection of Infectious Respiratory Disease Through Sweat from Axillary Using an E-Nose with Stacked Deep Neural Network," in IEEE Access, vol. 10, pp. 51285-51298, 2022, doi: 10.1109/ACCESS.2022.3173736.

N. Babu, D. Pruthviraja and J. Mathew, "Enhancing Lung Acoustic Signals Classification with Eigenvectors-Based and Traditional Augmentation Methods," in IEEE Access, vol. 12, pp. 87691-87700, 2024, doi: 10.1109/ACCESS.2024.3417183.

R. Mahum and A. S. Al-Salman, "Lung-RetinaNet: Lung Cancer Detection Using a RetinaNet with Multi-Scale Feature Fusion and Context Module," in IEEE Access, vol. 11, pp. 53850-53861, 2023, doi: 10.1109/ACCESS.2023.3281259.

S. R. Vinta, B. Lakshmi, M. A. Safali and G. S. C. Kumar, "Segmentation and Classification of Interstitial Lung Diseases Based on Hybrid Deep Learning Network Model," in IEEE Access, vol.12, pp.50444- 50458, 2024, doi: 10.1109/ACCESS.2024.3383144.

T. Nguyen and F. Pernkopf, "Lung Sound Classification Using Co-Tuning and Stochastic Normalization," in IEEE Transactions on Biomedical Engineering, vol. 69, no. 9, pp. 2872-2882, Sept. 2022, doi: 10.1109/TBME.2022.3156293.

T. Wanasinghe, S. Bandara, S. Madusanka, D. Meedeniya, M. Bandara and I. D. L. T. Díez, "Lung Sound Classification with Multi-Feature Integration Utilizing Lightweight CNN Model," in IEEE Access, vol. 12, pp. 21262-21276, 2024, doi: 10.1109/ACCESS.2024.3361943.

X. Zhang, D. Maddipatla, B. B. Narakathu, B. J. Bazuin and M. Z. Atashbar, "Intelligent Detection of Adventitious Sounds Critical in Diagnosing Cardiovascular and Cardiopulmonary Diseases," in IEEE Access, vol. 11, pp. 100029-100041, 2023, doi: 10.1109/ACCESS.2023.3313605.