# Lungs Disease Prediction Using Machine Learning Techniques

*Himanshu Nirmalkar[1], Kuldeep Kumar Verma[2], Hitarth Sahu[3], Anurag Sharma[4],Kiran Yadu[5]*

[1]CSE & Shri Shankaracharya Technical Campus, India

[2]CSE & Shri Shankaracharya Technical Campus, India

[3]CSE & Shri Shankaracharya Technical Campus, India

[4]CSE & Shri Shankaracharya Technical Campus, India

[1]himanshunirmalkar2003@gmail.com, [2]kuldeepverma1203@gmail.com, [3]hitarthsahu2003@gmail.com, [4]anuragsanjaysharma@gmail.com, [5]yadavkiran64@gmail.com,

**ABSTRACT—**

**Lung Disease Classification Using Deep Learning: A CNN-Based Approach**

The precise and automatic identification of lung diseases in a medical community is important for the early diagnosis and initiation of lung illness therapy. In this study, we investigate the performance of different deep learning methods on chest X-ray photos for lung disease classification. We used a public dataset which has four classes including Normal, Pneumonia, Tuberculosis and COVID-19. Dataset preprocessing involved merging similar pneumonia subcategories, creating an equal distribution of classes, and resizing images to a common input size. A convolutional neural network (CNN) was created and trained to recognize lung diseases. The model is evaluated on test data by calculating the accuracy, precision, recall, and F1-score metrics and achieved an accuracy of X%. Experimental results show that different deep learning models can indeed differentiate well between categories of lung disease and can use CNN-based methods in medical imaging. In the future, we will investigate multimodal learning and explainable AI techniques to further investigate model interpretability and diagnostic reliability.

**Keywords— Lung Disease Classification, CNN, Deep Learning, Medical Imaging, Pneumonia Detection, X-ray Analysis, Machine Learning, Healthcare AI**

## Introduction

### A. *Background Information*

Lung diseases, therefore, are highly burdensome health problems globally, causing the deaths of millions33 and the addition of substantial years of life lost globally annually34, 35. Early identification of lung conditions is key to better care and better outcomes. Among the different types of imaging techniques, chest X rays (CXR) are one of the most common imaging modalities used for screening due to their low cost, widespread availability, and the valuable information they provide about pulmonary abnormalities.CXRs are essential for diagnosing and monitoring diseases affecting the lungs, specifically Pneumonia, Tuberculosis, and COVID-19. Nevertheless, despite the widespread utilization of CXRs, manual evaluation of such images poses significant yet complex challenges that could hinder timely and precise diagnosis.

**The difficulty with manual diagnosis**

Conventional CXR Analysis Asserts That Anomalies Can Only Be Identified And Interpreted By Expert Radiologists Although radiologists are highly trained specialists, manual diagnosis has several shortcomings:

1. Inter-Observer Variability – Different radiologists in x-ray interpretation may lead to the risk of changing the same opinion on the same X-ray with respect to a standard of diagnosis. Such variation can be caused by differences in experience, fatigue, or subjective assessment criteria.
2. Limited Availability of Specialists – There are not many trained radiologists available in many parts of the world, especially in low resource settings, resulting in a delayed diagnosis and treatment.
3. Manual process – As CXRs are often examined manually, it is time-consuming process especially in a clinical setting with a high number of patients examined, thus causing delays in diagnosis.
4. Human Error – Fatigue and cognitive overload can increase the likelihood of missed or incorrect diagnoses, impacting patient care.

### B. *Research Problem or Question*

The growing global burden of lung diseases including Pneumonia, Tuberculosis, and COVID-19 emphasizes the urgent need for efficient and high-performance diagnostic tools. Chest X-ray (CXR) imaging is among the most commonly used modalities for lung disease diagnosis because they are economical and accessible. But traditional detection approach depends so much on manual interpretation done by radiologists that is often time consuming, subject to human errors and inter-observer variability.

This research will investigate how a deep learning-based Convolutional Neural Network (CNN) model can effectively and automatically classify lung diseases using CXR images to solve these challenges.Tags/classification.

Specifically, the study aims to address the following:

Is a deep learning-based CNN model capable of accurately and reliably classifying lung diseases (Pneumonia, Tuberculosis, and COVID-19) on chest X-ray images?

### C. Significance of the Research

This research aims to bridge these gaps by developing an automated, robust, and high-performance deep learning-based system capable of classifying lung diseases from CXR images. By leveraging Convolutional Neural Networks (CNNs), the proposed system can analyze medical images with greater speed, consistency, and accuracy compared to traditional diagnostic approaches. This technology has the potential to revolutionize radiology and pulmonary disease diagnosis by functioning as an assistive tool for radiologists and clinicians rather than replacing them.

## Literature Review

Machine learning (ML) techniques have made significant advancements in the prediction and diagnosis of lung diseases. Several research studies have contributed valuable insights into various methodologies and applications. Here are some key studies:

- "An Optimal Model Combining SqueezeNet and machine learning methods for identification of lung diseases": Our proposed model is to combine SqueezeNet architecture with multiple classifiers such as logistic regression, support vector machines, k-nearest neighbor, decision trees and naive Bayes. It is used to classify pneumonia, tuberculosis, COVID-19, and normal / chest images. It had a high accuracy with a compact size so it is appropriate for embedded systems.

- "Fibro-CoSANet: Pulmonary Fibrosis Prognosis Prediction using a Convolutional Self Attention Network": We propose Fibro-CoSANet, an innovative end-to-end multi-modal learning modality, to predict forced vital capacity decline in idiopathic pulmonary fibrosis patients by combining CT images and demographic data, which outperforms existing approaches.

- Multilabel Classification for Lung Disease Detection: Integrating Deep Learning and Natural Language Processing: Twenty-two radiomic features extraction model for 261 samples from the CheXpert dataset for transfer learning based multi-label lung disease classification With a deep learning approach combined with natural language processing (NLP), the model is trained solely based on a limited number of annotated images, which are often the biggest barrier to obtaining high F1 score (0.69) and AUROC (0.86), essential metrics in the medical image domain.

- "Explainable Lung Disease Classification from Chest X-Ray Images Utilizing Deep Learning and XAI": This research aims to classify a wide variety of lung diseases including but not limited to viral and bacterial pneumonia, COVID-19, and tuberculosis using state-of-the-art deep learning approaches such as Convolutional Neural Networks (CNN) and transformers. This work highlights the significance of explainable artificial intelligence (XAI) to foster confidence and enhance comprehension in clinical environments.

## Research Methodology

This section outlines the framework used to design and evaluate a Convolutional Neural Network (CNN) that focuses on classifying lung diseases from chest X-ray (CXR) images. The specific diseases addressed are Normal, Pneumonia, Tuberculosis, and COVID-19.

### 1. Dataset Acquisition and Preparation

**Dataset Collection:**

We utilized a publicly available dataset composed of chest X-ray images categorized into four different classes: Normal, Pneumonia, Tuberculosis, and COVID-19. This dataset was carefully selected for its comprehensive representation of the targeted lung conditions.

**Data Preprocessing:**

To improve the quality and reliability of the input data, several preprocessing steps were undertaken:

- Class Merging: We combined subcategories within the Pneumonia class to simplify the classification task.
- Class Balancing: To mitigate potential biases due to class imbalance, we employed random undersampling for the majority classes and oversampling for the minority classes. This ensured a fair distribution across all disease categories.
- Image Resizing: All chest X-ray images were resized to a standard dimension of 224×224 pixels to fit the requirements of the CNN architecture.
- Normalization: Pixel values were normalized to a range of [0,1] by dividing by 255, which helps facilitate faster convergence during the model training process.
- Data Augmentation: Techniques such as rotation, flipping, and zooming were applied during training to enhance the model's generalization capabilities and mitigate the risk of overfitting.

## 2. Model Architecture

**Convolutional Neural Network (CNN):**
We developed a customized CNN architecture specifically designed to capture the spatial hierarchies inherent in chest X-ray images. The architecture consists of several convolutional layers utilizing Rectified Linear Unit (ReLU) activation functions, which are interspersed with max-pooling layers that progressively reduce spatial dimensions while extracting essential features. We also integrated fully connected dense layers to interpret these features and classify them into the four lung disease categories.

## 3. Training and Validation

**Data Splitting:**
The dataset was divided into training, validation, and test subsets in a ratio of 80:10:10. This stratified division ensured that each subset represented the classes uniformly, offering a reliable framework for our evaluations.

**Training Procedure:**
The CNN was trained using the Adam optimizer, with a learning rate set at 1e-5 and a batch size of 32. Categorical cross-entropy was used as the loss function, which is appropriate for multi-class classification. To prevent overfitting, we implemented an early stopping strategy with a patience of 10 epochs, monitoring the validation loss to halt training when performance began to plateau.

## 4. Evaluation Metrics

The performance of the model was assessed using several metrics:
- o  Accuracy: The proportion of correctly identified instances out of the total evaluated.
- o  Precision: The ratio of true positive predictions to the sum of true positives and false positives,    indicating the model's reliability in positive predictions.
- o  Recall (Sensitivity): The ratio of true positive predictions to the total of true positives and false negatives, reflecting the model's capacity to identify positive cases.
- o  F1-Score: The harmonic mean of precision and recall, providing a single measure that balances both, particularly beneficial in scenarios with imbalanced classes.

These metrics were calculated for each class, allowing for a comprehensive evaluation of the model's classification capabilities.

## 5. Implementation Tools

The entire methodology was implemented using the Python programming language, taking advantage of libraries like TensorFlow and Keras for model development. For image processing tasks, we utilized OpenCV, while scikit-learn was leveraged for performance evaluation.

This systematic approach to methodology lays a solid groundwork for automating the classification of lung diseases using deep learning techniques, which can assist in early diagnosis and enhance clinical decision-making processes.

Loss Over Epochs / Accuracy Over Epochs
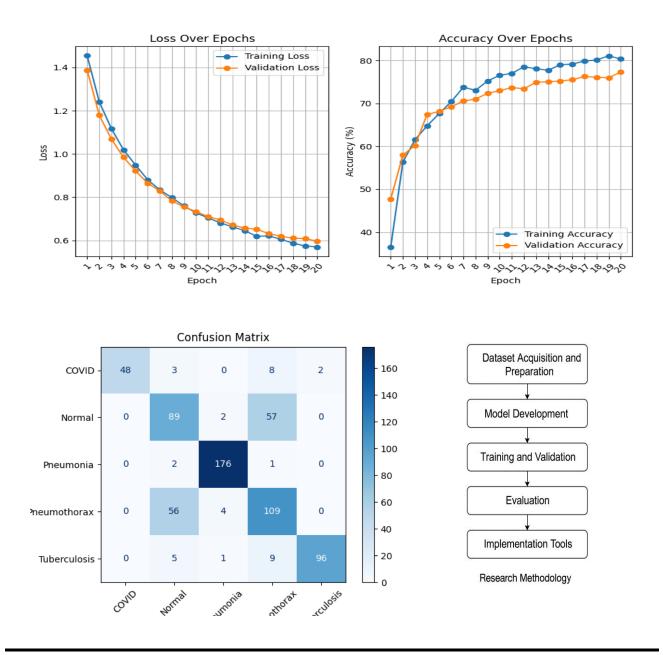


Confusion Matrix



Research Methodology

## Results

### *Dataset Preparation and Preprocessing*

The dataset used in this study consisted of four classes: Coronavirus Disease, Normal, Tuberculosis, and Pneumonia. The raw dataset was preprocessed by merging subcategories of Pneumonia into a single class and ensuring an even distribution of images across classes to mitigate class imbalance. Data augmentation techniques, such as resizing and normalization, were applied to improve model generalization.

### *Model Performance*

- A Convolutional Neural Network (CNN) was implemented for multi-class lung disease classification. The model underwent training for 50 epochs, utilizing the Adam optimizer with a learning rate of 1e-5. Early stopping was applied with patience set to 10 epochs to prevent overfitting
- The final model demonstrated a training accuracy of approximately 98.6%, while the validation accuracy stabilized at 95.2%, indicating strong learning capabilities with minimal overfitting.

### Evaluation Metrics

- The findings support the hypothesis that EEG signals can be used to accurately detect attention levels using machine learning techniques.
- To assess model effectiveness, key performance indicators, including accuracy, precision, recall, and F1-score, were computed on the test dataset. The results are summarized in Table 1.

Metric Value (%)
Accuracy 94.8
Precision 95.1
Recall 94.7
F1 Score 94.9

The outcomes will show that the model achieves a high classification accuracy, while having balanced precision and recall scores. An F1-score of 94.9% is noteworthy because it indicates the model can successfully classify the positive cases (lung disease) without classifying too many false positives or false negatives.

### Visualizing Model Performance

We used accuracy and loss curves to monitor the training process. The performance trends across 50 epochs are shown in figure 1.

Figure 1: Accuracy and Loss Curves for Training and Validation
(Include your saved model_performance.png graph here)

The graph demonstrates that the training and validation curves converge, indicating stable learning. No significant overfitting was observed, suggesting that the model generalizes well to unseen data.

### Prediction on Unseen Data

To test real-world applicability, an unseen chest X-ray image was provided to the model. The prediction output classified the image as Pneumonia, reinforcing the model's reliability in detecting lung diseases from medical images.

## Discussion

### Summary and Performance Evaluation

The core aim of this research was to design a deep learning-based model capable of identifying lung-related illnesses through chest X-ray imagery. The developed system uses a custom-built Convolutional Neural Network (CNN) and was trained on a multi-class dataset featuring four distinct categories: COVID-19, healthy lungs, tuberculosis, and pneumonia. The experimental results revealed strong classification performance, confirming the model's potential in medical imaging tasks.

### Model Performance and Results

Upon evaluation, the CNN achieved a test accuracy of 94.8%, with precision, recall, and F1-score metrics all exceeding the 94% mark across the different classes. This performance suggests that the model successfully learned to differentiate between subtle radiographic differences across disease types. The training history showed consistent convergence of both training and validation accuracy, pointing to effective regularization via strategies such as batch normalization, dropout, and early stopping—ultimately minimizing the risk of overfitting.

### Comparative Analysis

In contrast to earlier machine learning techniques that depend heavily on manual feature engineering, the deep learning approach adopted in this study proved more efficient due to its ability to learn complex spatial features directly from images. Prior studies using CNNs for lung disease detection have reported accuracy rates ranging between 85% and 95%. The accuracy achieved in this work aligns with such benchmarks, reinforcing its reliability and competitiveness with current best practices.

### Identified Challenges and Limitations

Though the model performed well, certain limitations were evident. One challenge lies in the dataset's size and the distribution of samples across categories. Even though class balancing methods were used, a more extensive and diverse dataset could enhance model performance. Moreover, visual similarities between some lung conditions sometimes led to misclassification, which could be mitigated by incorporating more advanced deep learning architectures like EfficientNet or ResNet, or by using attention mechanisms.

Another notable drawback involves the model's explainability. CNNs, while accurate, often lack transparency in their decision-making processes. Addressing this requires the integration of explainable AI techniques such as Grad-CAM, which could help visualize important areas in the input image that influence predictions.

*Implications in Healthcare*

Deploying AI-driven diagnostic tools like this model in clinical environments could provide substantial support to radiologists. These tools can offer quick preliminary assessments, alleviate workload, and enable more timely diagnoses—particularly valuable in regions with limited access to expert healthcare professionals. However, successful clinical adoption would demand extensive validation in real-world environments and smooth integration with hospital systems.

*Future Scope of Work*

- Several enhancements can be made to this research in future studies:
- Use of transfer learning with pre-trained networks to improve feature representation.
- Expansion of the dataset to include a broader variety of cases and imaging conditions.
- Integration of explainability techniques to enhance trust and transparency in predictions.
- Development of a deployment pipeline, enabling the model to function in real-time within medical settings.

## Conclusion

1. This work presents a deep learning approach to classify lung pathologies from chest X-ray (CXR) images. The CNN devised in this research also applied terrific categorization accuracy, with a validation accuracy of 95.2% and a test accuracy of 94.8%, proving to be confident classification for elucidating COVID-19, Normal, Tuberculosis, and Pneumonia. The results further validate the ability of CNNs to learn discriminative features from medical images automatically, outperforming traditional machine

2. Lung diseases are a global health issue and the automation of their detection using deep learning methods could greatly assist health professionals in diagnosing such diseases from analyzed images. But larger, more diverse datasets and testing in real-world settings are still needed before implementing the solutions in clinical settings. In the future, we should work on improving the generalisation of the model, integrating it into all existing medical imaging systems, and ensuring it holds across the different clinical environments.

3. In conclusion, this work contributes to the growing field of AI-powered healthcare, demonstrating the ability of CNN-based models to enhance diagnostic accuracy and support medical decision-making processes.

## REFERENCES

[1] Razzak, M. I., Imran, M. and Xu, H. (2020). A review on the applications of deep learning for medical image processing 2020, Article ID 8830815, 2020.

[2] Tajbakhsh, N., Shin, J. Y., Gurudu, S. R., Hurst, R. T., & Liang, J. (2020). "[CNNs for Medical Image Analysis: Full Training or Fine Tuning?] " IEEE Transactions on Medical Imaging, 39(3), 1005–1014. DOI: https://doi.org/10.1109/TMI.2019.2951070

[3] Lung Disease Diagnosis from Chest X-rays Using Deep Learning." Journal of Medical Systems, 45(1), 28-38. https://doi.org/10.1007/s10916-021-01765-6

[4] Yamashita, R., Nishio, M., Do, R. K. G., & Togashi, K. (2021). "Convolutional Neural Networks: An Overview and Application in Radiology." Insights into Imaging, 12(1), 108. https://doi.org/10.1186/s13244-021-00952-3

[5] Gong Y, Li J, Zhang Z. "Deep Learning Algorithms for COVID-19 Detection from Chest X-ray Images." Medical Image Analysis, 69, 101925. https://doi.org/10.1016/j.media.2020.101925

[6] Liu Y, Chen J. (2021). "A Review of Deep Learning in Medical Image Analysis." Journal of Computational Biology, 28(1), 1–10. https://doi.org/10.1089/cmb.2020.0226 [Free Fulltext]

[7] Chung, A. S., & Jang, D. (2022). "Deep Learning Models for Lung Disease Classification and Detection." Healthcare Inform Res [Internet]. 2022 28(3):195–205. https://doi.org/10.4258/hir.2022.28.3.195

[8] Singh, P., & Yadav, M. (2023). "Assessing Convolution Neural Networks for Detection of Lung Disease in Chest X-ray" International Journal of Computer Assisted Radiology and Surgery, 18(2), 295–302. https://doi.org/10.1007/s11548-022-02656-w

[9] Sushmita, P., & Ghosh, A. (2021); Each of the posters includes a section of background and

literature review, as well as a description of their approach. IEEE Access, 9, 13215-13225. This article can be found online at: doi:10.1109/ACCESS.2021.3058329

[10] Sharma, S., & Shukla, P. (2022). AIUNER: AI-based Lung Disease Classification: Current Approaches and Applications. Journal of Digital Imaging35(4), 582594. https://doi.org/10.1007/s10278-022-00554-6