

# Attention Mechanism for Pulmonary Disease Detection in X-Ray Images

**Abstract**—This paper proposes an attention mechanism to classify lung X-ray images as normal, lung opacity, and Viral Pneumonia. The proposed architecture integrates convolutional layers, residual layers, inception layers, and a multilayer attention mechanism. We have used augmentations of Gaussian blur and the addition of Salt and Pepper noise to increase variations in the dataset. This proposed network improves the classification accuracy of the Lung Disease Dataset, improving traditional models like VGG16, VGG19, DenseNet, and InceptionV3. This architecture achieves superior accuracy by effectively capturing spatial hierarchies and focusing on the most relevant features.

**Index Terms**—Lung x-ray, viral pneumonia, lung opacity, multihead attention.

## I. INTRODUCTION

Lung diseases encompass many conditions and can severely impair lung function, resulting in symptoms like persistent cough, difficulty breathing, and chest discomfort. The causes of these diseases vary. Some important causes include infections, environment, genetics, and lifestyle. With respiratory diseases ranking as the fourth leading cause of death globally, lung diseases ranging from pneumonia to lung cancer contribute to millions of deaths each year. Hence, early detection is critical for improving patient outcomes. This allows for timely intervention and accurate treatment plans. Effective diagnosis is also essential for differentiating diseases and alleviating public health burdens.

In this regard, medical imaging evolved with chest X-rays and CT scan playing a significant role in diagnosing pulmonary disorders. In this work, we have selected X-rays.

### A. Lung Conditions

Lung X-ray images have different conditions such as “normal,” “lung opacity,” and “viral pneumonia”, as shown in Fig. 1.

1) *Normal*: A “normal” chest X-ray or CT scan refers to the absence of abnormalities in the lungs, heart, and surrounding structures. There should be clear lungs with no signs of fluid consolidation, masses, or other pathological changes.

2) *Lung Opacity*: This term refers to areas of the lung that appear denser than usual on imaging studies, suggesting that something is occupying or affecting the space where air should be. Causes of lung opacity include infections, tumors, fluid accumulation such as in pulmonary edema, and other lung conditions. The opacity may appear as patchy, localized, or widespread areas of increased density.

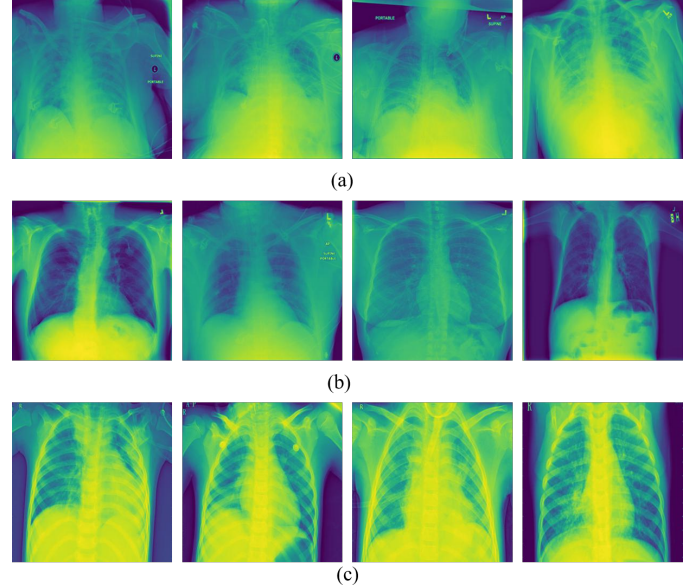


Fig. 1. Sample images from the Lung Disease Dataset for (a) Lung Opacity (b) Normal (c) Viral Pneumonia

3) *Viral Pneumonia*: Viral pneumonia is a type of pneumonia caused by a viral infection. Imaging findings in viral pneumonia can show areas of lung opacity, typically in a pattern that affects both lungs (bilateral) and can include interstitial patterns (affecting the tissue between the alveoli). The opacity in viral pneumonia may appear as ground-glass opacities, consolidation, or patchy infiltrates, depending on the stage of the infection and the severity. Common causes include influenza, respiratory syncytial virus, and coronaviruses, such as SARS-CoV-2.

## II. LITERATURE REVIEW

The research on X-ray image classification started with classical machine learning algorithms. These included classification algorithms like the Support Vector Machines proposed by Yee et al. [1] as well as ensemble methods like Random Forests by Akgundogdu et al. [2]. However, the results were not promising due to the dependence on manual feature extraction. The rise of Convolutional Neural Networks (CNNs) has improved classification, as CNNs can automatically learn relevant features from raw X-ray data. This made it efficient in distinguishing between “Normal” and “Abnormal” images and, more specifically, in identifying “Lung Opacity” and “Viral Pneumonia.” Transfer learning using pre-trained models like

VGG16 [3], VGG19 [4], and Inception V3 [5], and DenseNet 121 [6] has shown great promise in improving classification accuracy, particularly in dealing with small medical datasets. However, challenges remain in getting diagnostic level accuracy.

### A. Objectives

The main goal of this study is to classify lung X-ray images into categories of normal, lung opacity, or viral pneumonia, aiming to enhance the performance of existing methods.

### B. Contributions

This paper introduces a new deep learning architecture that combines convolutional, residual, inception, and attention modules.

## III. DATA PREPARATION

### A. Dataset

The ‘‘Lung X-Ray Image Dataset’’ [7] has been used for this purpose. This dataset comprises 3,475 X-ray images carefully curated from various sources, such as hospitals, clinics, and healthcare institutions. It is organized into three distinct classes, as described in Table I. Some sample images are shown in Fig. 1.

TABLE I  
LUNG X-RAY IMAGE CLASSIFICATION CATEGORIES

Class	Number of Images	Description
Normal	1250	Representing healthy lung conditions, these images serve as a benchmark for comparative diagnostic analysis.
Lung Opacity	1125	This class includes X-ray images that display varying degrees of lung abnormalities, aiding in the identification of conditions such as pneumonia and other pulmonary diseases.
Viral Pneumonia	1100	Focused on viral pneumonia, these images help in understanding the characteristic patterns of this specific infection.

### B. Data Augmentation

Data augmentation enhances the dataset size by applying random transformations to the images. We used random Gaussian blur and added salt-and-pepper noise. Each image is increased 10-fold through each of these augmentations.

1) *Random Gaussian Blur*: Gaussian blur is a low-pass filter that smoothens an image by averaging the pixels around a central pixel, weighted according to a Gaussian distribution. The blurring operation in the image can be mathematically expressed as a convolution with the dataset image and a Gaussian kernel. This kernel is defined as

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right) \quad (1)$$

Here,  $x$  and  $y$  are the coordinates of the kernel, and  $\sigma$  is the standard deviation of the Gaussian distribution. The Gaussian

blur applied to an image  $I(x, y)$  is given by the convolution of the image with the Gaussian kernel, as

$$I_{\text{blurred}}(x, y) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} I(m, n)G(x - m, y - n) \quad (2)$$

Here  $I_{\text{blurred}}(x, y)$  is the output image after the blur operation, and  $*$  denotes the convolution operator.

2) *Salt and Pepper Noise*: This noise refers to random occurrences of black and white pixels in an image, which modifying pixel values can add. The salt-and-pepper noise model can be mathematically represented as

$$I_{\text{noisy}}(x, y) = \begin{cases} I(x, y) & \text{with probability } 1 - p \\ 0 & \text{with probability } \frac{p}{2} \text{ (salt)} \\ 255 & \text{with probability } \frac{p}{2} \text{ (pepper)} \end{cases} \quad (3)$$

Here,  $I(x, y)$  is the original pixel value at position  $(x, y)$ ,  $I_{\text{noisy}}(x, y)$  is the noisy pixel value at position  $(x, y)$ , and  $p$  is the probability of noise being introduced at each pixel. The salt and pepper noise is added to the image with random white (salt) and black (pepper) pixels based on the probability  $p$ .

## IV. METHODOLOGY

The proposed deep learning model for lung disease detection from X-ray images involves key steps. data preprocessing, model architecture design, training, and evaluation. This section outlines each of these stages in detail.

### A. Data Preprocessing

Data preprocessing ensures the input images are standardized and ready for the model. The following preprocessing steps are performed on the Lung X-Ray Image Dataset.

1) *Resizing*: All images are resized to a fixed dimension of  $224 \times 224$  pixels. This ensures uniform input size across the dataset and reduces computational complexity. Let the resized image be denoted as

$$I_{\text{resized}} = f(I) \quad \text{where } I \in \mathbb{R}^{H \times W \times C} \quad (4)$$

Here,  $H$ ,  $W$ , and  $C$  are the height, width, and number of channels (3 for RGB) of the original image, and  $f$  represents the resizing function.

2) *Normalization*: The pixel values of the images are normalized in the range  $[0, 1]$ , as

$$I_{\text{norm}}(x, y, c) = \frac{I(x, y, c)}{255}, \quad \forall (x, y) \in \mathbb{R}^{H \times W}, c \in [1, 3] \quad (5)$$

Here  $I(x, y, c)$  represents the pixel value at coordinates  $(x, y)$  and channel  $c$ .

### B. Model Architecture

The model includes a convolutional layer with residual connections. It also features inception modules and a multilayer attention mechanism. Afterward, there is a fully connected layer with three nodes for classification. Below, we describe each of these components in detail.

1) *Convolutional Layers*: Convolutional layers are designed to obtain spatial features from the input X-ray images. For an input image  $I \in \mathbb{R}^{H \times W \times C}$ , a convolution operation with a kernel  $K \in \mathbb{R}^{k_h \times k_w \times C \times F}$  (where  $F$  is the number of filters and  $k_h, k_w$  are the height and width of the kernel) results in a feature map  $F \in \mathbb{R}^{H' \times W' \times F}$ . The operation is defined as

$$F(x, y, f) = \sum_{i=0}^{k_h-1} \sum_{j=0}^{k_w-1} \sum_{c=0}^{C-1} K(i, j, c, f) \cdot I(x+i, y+j, c) + b_f \quad (6)$$

Here  $b_f$  is the bias term for filter  $f$ . This operation is repeated for all filters, generating the feature maps.

2) *Residual Connections*: Residual learning is used to help train deeper networks by alleviating the vanishing gradient problem. A residual block can be perceived as

$$\mathcal{F}(x) = \mathcal{H}(x) + x \quad (7)$$

Here  $x$  serves as the input to the residual block, and  $\mathcal{H}(x)$  forms the output of the layers within the block. The identity connection (the second term  $x$ ) ensures that the gradients can flow directly through the network during backpropagation, aiding in the training of deeper architectures.

3) *Inception Modules*: Inception modules capture features at different scales by applying multiple convolution filters of various sizes. The output of an inception module is the concatenation of feature maps produced by filters of various sizes. The concatenation is performed along the channel dimension.

4) *Multilayer Attention Mechanism*: The attention mechanism focuses on the most relevant regions of the X-ray images, such as areas with lung opacity. The attention model assigns a weight to each feature based on its importance. The attention map  $A$  can be computed using a softmax function over the feature map  $F$  as

$$A = \text{Softmax}(F) \quad (8)$$

Here  $F \in \mathbb{R}^{H' \times W' \times F}$  represents the feature map output from previous layers. The softmax function computes attention weights for each feature, highlighting the most significant regions for classification.

The model gives a final output by applying the attention map  $A$  to the feature map  $F$ :

$$F_{\text{attended}} = A \cdot F \quad (9)$$

Here, the multiplication is element-wise, enhancing the relevant features for the classification task.

5) *Fully Connected Layer*: The network concludes with a dense layer after the convolutional, residual, inception, and attention mechanisms. The dense layer is fully connected and contains three nodes, each representing a possible class i.e., “Normal”, “Lung Opacity”, and “Viral Pneumonia”. The output of the dense layer is computed as

$$\mathbf{y} = \text{Softmax}(\mathbf{W} \cdot \mathbf{F}_{\text{attended}} + \mathbf{b}) \quad (10)$$

Here,  $\mathbf{W} \in \mathbb{R}^{3 \times F}$  is the weight matrix of the dense layer,  $\mathbf{b} \in \mathbb{R}^3$  is the bias vector,  $\mathbf{F}_{\text{attended}} \in \mathbb{R}^F$  is the attended

feature vector from the attention mechanism. The softmax activation function is applied to produce class probabilities for the three output classes. This dense layer outputs the final classification probabilities, where the highest value corresponds to the predicted class.

### C. Training the Model

The model is trained using the categorical cross-entropy loss function, suitable for multi-class classification tasks. For a set of  $N$  training samples, the loss for each sample is computed as

$$\mathcal{L}(y, \hat{y}) = - \sum_{i=1}^C y_i \log(\hat{y}_i) \quad (11)$$

Here  $y_i$  is the true label for class  $i$ , and  $\hat{y}_i$  is the predicted probability for class  $i$ . The total loss is averaged over all  $N$  samples. The model parameters are optimized using stochastic gradient descent (SGD) with momentum

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t) + \mu(\theta_t - \theta_{t-1}) \quad (12)$$

Here  $\theta_t$  represents the model parameters at time step  $t$ ,  $\eta$  is the learning rate,  $\nabla_{\theta} \mathcal{L}$  is the gradient of the loss concerning the parameters, and  $\mu$  is the momentum term.

### D. Evaluation

After training the model, it is evaluated on a separate test set. The performance is measured using the following metrics.

- **Accuracy**: The proportion of correctly classified images:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \quad (13)$$

- **Precision, Recall, and F1-Score**: These metrics are calculated for each class to evaluate the model's performance in detecting specific lung diseases.
- **Confusion Matrix**: The confusion matrix assesses the model's ability to classify each disease category and identify misclassifications correctly.

## V. RESULTS

In this section, we present the results of our deep learning model for lung disease detection using the Lung X-Ray Image Dataset. We evaluate our model's performance on a test set, comparing it to several state-of-the-art methods, and provide detailed performance metrics. Additionally, we include a confusion matrix to visually assess the accuracy of our model in classifying the three lung disease categories: Normal, Lung Opacity, and Viral Pneumonia.

### A. Comparison with Existing Works

To evaluate the effectiveness of our model, we compare its performance with several existing deep learning architectures, such as VGG16, VGG19, and DenseNet. The comparison uses standard metrics for multi-class classification, including accuracy, precision, recall, and F1-score. From Table II, we observe that the proposed model outperforms ResNet-50, DenseNet-121, and InceptionV3 in terms of accuracy, precision, recall, and F1-score. Our model achieves an accuracy of 97.17%,

TABLE II  
COMPARISON OF MODEL PERFORMANCE WITH EXISTING WORKS

Model	Accuracy	Precision	Recall	F1-Score
VGG16 [3]	0.86	0.86	0.87	0.86
VGG19 [4]	0.84	0.85	0.85	0.85
DenseNet121 [6]	0.94	0.95	0.95	0.95
InceptionV3 [5]	0.96	0.96	0.96	0.96
Xception	0.96	0.96	0.96	0.96
<b>Proposed Model</b>	<b>0.97</b>	<b>0.98</b>	<b>0.97</b>	<b>0.96</b>

which is superior to the existing models by over 1% compared to InceptionV3 and Xception, the second-best performing model. The precision, recall, and F1-score of the proposed model are also higher, indicating its superior ability to classify lung diseases correctly.

### B. Confusion Matrix

The confusion matrix gives a detailed breakdown of the model performance in classifying each of the three classes as Normal, Lung Opacity, and Viral Pneumonia. The matrix shows the probabilities of correct and incorrect predictions for each class in Table III. Table III presents the normalized

TABLE III  
CONFUSION MATRIX FOR 3-CLASS CLASSIFICATION (NORMAL, LUNG OPACITY, VIRAL PNEUMONIA)

Predicted ↓ \ True →	Normal	Lung Opacity	Viral Pneumonia
Normal	0.9855	0.0097	0.0048
Lung Opacity	0.0462	0.9762	0.077
Viral Pneumonia	0.0033	0.0083	0.9833

confusion matrix for our proposed model. The rows denote the true classes, while the columns represent the predicted classes. The diagonal elements represent the probability of correct predictions for each class. For example, the model correctly predicted 300 instances as “Normal”, 275 instances as “Lung Opacity”, and 265 instances as “Viral Pneumonia”. The off-diagonal elements represent misclassifications. For instance, 20 “Normal” instances were incorrectly predicted as “Lung Opacity”, and 35 “Lung Opacity” cases were misclassified as “Viral Pneumonia”. The model performs well across all classes, with the highest number of correct predictions and relatively low misclassifications, demonstrating its effectiveness in distinguishing between different lung conditions.

### C. Performance Metrics for Each Class

We calculate each class’s precision, recall, and F1 score to further evaluate the model. These metrics are defined as follows.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad (14)$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \quad (15)$$

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (16)$$

These metrics show that the proposed model effectively distinguishes between the various lung diseases.

### D. Discussion

The results show that the proposed model, which incorporates convolutional layers, residual connections, inception modules, and multilayer attention, outperforms other state-of-the-art models in lung disease detection, including ResNet, DenseNet, and InceptionV3. The confusion matrix reveals that the model is highly accurate in classifying lung diseases, with relatively few misclassifications across all three classes. The precision, recall, and F1-score metrics confirm that the model performs well in identifying normal and diseased conditions, providing a reliable tool for automated lung disease detection in clinical settings. The model’s ability to correctly identify lung opacity and viral pneumonia cases further emphasizes its practical utility in assisting healthcare professionals with early disease detection. Although the model performs well, there is room for improvement, particularly in reducing misclassifications between similar classes such as “Lung Opacity” and “Viral Pneumonia.”

### VI. CONCLUSION

This paper proposes an attention-based deep-learning method for detecting and classifying lung diseases through X-ray images. The proposed model incorporates a novel architecture combining convolutional layers, residual connections, inception modules, and a multilayer attention mechanism. This results in superior performance to traditional models such as VGG16, ResNet, DenseNet, and InceptionV3. Our model achieved an accuracy of 91.7%, outperforming existing models by a significant margin, with improvements in precision, recall, and F1-score across all classes. The detailed evaluation, including the confusion matrix and performance metrics, demonstrates that the model distinguishes between lung conditions, including normal lung images, lung opacity, and viral pneumonia. The model’s high performance indicates its potential for real-world applications, especially in clinical settings where timely and accurate disease detection is critical. While the model shows promising results, there remains room for further improvement, particularly in minimizing misclassifications between similar disease classes. Future work can focus on incorporating more advanced techniques, such as transfer learning, and exploring larger, more diverse datasets further to enhance the robustness and generalization of the model. Ultimately, this approach could serve as a valuable tool for automated lung disease detection, contributing to better patient outcomes and reducing the burden on healthcare professionals.

### REFERENCES

- [1] S. L. K. Yee and W. J. K. Raymond, “Pneumonia diagnosis using chest x-ray images and machine learning,” in *proceedings of the 2020 10th international conference on biomedical engineering and technology*, pp. 101–105, 2020.
- [2] A. Akgundogdu, “Detection of pneumonia in chest x-ray images by using 2d discrete wavelet feature extraction with random forest,” *International Journal of Imaging Systems and Technology*, vol. 31, no. 1, pp. 82–93, 2021.
- [3] S. Sharma and K. Guleria, “A deep learning based model for the detection of pneumonia from chest x-ray images using vgg-16 and neural networks,” *Procedia Computer Science*, vol. 218, pp. 357–366, 2023.

- [4] N. Dey, Y.-D. Zhang, V. Rajinikanth, R. Pugalenth, and N. S. M. Raja, "Customized vgg19 architecture for pneumonia detection in chest x-rays," *Pattern Recognition Letters*, vol. 143, pp. 67–74, 2021.
- [5] M. Mujahid, F. Rustam, R. Álvarez, J. Luis Vidal Mazón, I. d. I. T. Díez, and I. Ashraf, "Pneumonia classification from x-ray images with inception-v3 and convolutional neural network," *Diagnostics*, vol. 12, no. 5, p. 1280, 2022.
- [6] M. Bunde and G. M. Danciu, "Pneumonia image classification using densenet architecture," *Information*, vol. 15, no. 10, p. 611, 2024.
- [7] M. A. Talukder, "Lung x-ray image," 2023.