

# Pneumonia Classification from Chest X-ray Images using Vision Transformer

Mr. A. Priyatham, M.Tech (Ph.D)  
Assistant Professor  
Tirumala Engineering College  
Jonnalagadda, India  
Email: avulapriyatham@gmail.com

S. Manasa  
Department of ECE  
Tirumala Engineering College  
Jonnalagadda, India  
Email: manasasangana17@gmail.com

S. Bhuvaneshwari  
Department of ECE  
Tirumala Engineering College  
Jonnalagadda, India  
Email: sankubhuvana@gmail.com

Sd. Dariya Hussain  
Department of ECE  
Tirumala Engineering College  
Jonnalagadda, India  
Email: sddariyahussain@gmail.com

S. Abhishek Naidu  
Department of ECE  
Tirumala Engineering College  
Jonnalagadda, India  
Email: abhisheknaidu2005@gmail.com

**Abstract**—Pneumonia is a major respiratory disease that necessitates early and accurate diagnosis to reduce morbidity and mortality. Chest X-ray imaging is widely used for clinical assessment; however, manual interpretation is often time-consuming and subject to inter-observer variability. In this work, a Vision Transformer (ViT)-based deep learning framework is proposed for automated pneumonia classification from chest X-ray images. Unlike traditional convolutional neural networks, the Vision Transformer leverages a self-attention mechanism to model long-range dependencies and capture global contextual information within medical images. The dataset is preprocessed using image resizing, normalization, and data augmentation techniques to improve generalization and mitigate overfitting. The proposed model is evaluated using standard classification metrics. Experimental results demonstrate that the framework achieves an accuracy of 98.68%, with a precision of 0.9870, recall of 0.9868, and F1-score of 0.9868. These findings indicate that the proposed approach provides robust and reliable performance, making it suitable as a computer-aided diagnostic tool for supporting clinical decision-making.

**Keywords:** Pneumonia Detection, Vision Transformer, COVID-19, Deep Learning.

## I. INTRODUCTION

Invasive, and easily accessible in most healthcare Pneumonia is a serious respiratory disease that continues to affect millions of people across the world every year. It primarily impacts the lungs by causing inflammation in the air sacs, which may fill with fluid or pus, leading to symptoms such as fever, persistent cough, chest pain, and difficulty in breathing. According to global health reports, pneumonia remains one of the leading causes of death, particularly among children under the age of five and elderly individuals [3]. Due to its widespread impact and severity, early detection and accurate diagnosis are essential to ensure timely treatment and reduce mortality rates.

Chest X-ray imaging is one of the most commonly used diagnostic tools for identifying pneumonia. It is widely preferred because it is cost-effective, non-invasive. Radiologists analyze chest X-ray images to detect abnormalities in lung structures. However, interpreting these images is not always straightforward. Many lung diseases, including COVID-19 and other respiratory infections, exhibit similar visual patterns, which makes it difficult to differentiate between them [13]. This

similarity often leads to confusion during diagnosis and may result in delayed or incorrect treatment.

Another major challenge in pneumonia diagnosis is the reliance on human expertise. In many parts of the world, especially in rural and underdeveloped regions, there is a shortage of experienced radiologists. Even in well-equipped hospitals, radiologists often face heavy workloads, which can lead to fatigue and reduced diagnostic accuracy. Studies have also shown that there can be variability in interpretation among different experts, which affects the consistency and reliability of diagnosis [13]. These challenges highlight the need for automated systems that can assist medical professionals in analysing chest X-ray images more efficiently and accurately.

With the rapid advancement of artificial intelligence, deep learning techniques have emerged as powerful tools for medical image analysis. Convolutional Neural Networks (CNNs) have been widely used for image classification tasks, including pneumonia detection. These models can automatically extract important features from images and have shown promising results in improving diagnostic performance [10], [15]. Their ability to learn directly from data reduces the need for manual feature extraction and enhances efficiency.

However, despite their success, CNN-based models have certain limitations. They mainly focus on local features within small regions of the image and may not effectively capture the overall structure or global relationships within the image. In medical imaging, understanding the complete context of the image is crucial for accurate diagnosis. This limitation can affect the performance of CNN models, especially in complex cases where global information plays an important role [4].

To overcome these limitations, Vision Transformer (ViT) has been introduced as a novel approach for image classification. Unlike CNNs, Vision Transformers process images by dividing them into smaller patches and applying self-attention mechanisms to analyze the relationships between these patches. This allows the model to capture both local and global feature effectively, leading to improved performance in image classification tasks [1].

The study in [3] highlights the global impact of pneumonia using statistical health data, showing that it remains one of the leading causes of death, especially among children and elderly individuals. However, this work is limited to observational analysis and does not provide any automated diagnostic solution. To address diagnosis, [6] introduces CheXNet, a Convolutional Neural Network based on DenseNet-121 for pneumonia detection from chest X-rays, achieving around 76–77% accuracy. Despite its effectiveness, it suffers from false positives and limited generalization. Similarly, [13] applies transfer learning using pre-trained CNN models such as VGG and MobileNet for detecting COVID-19 from X-ray images, achieving high accuracy of around 96–98% and F1-scores near 0.95. However, these models are prone to overfitting and dataset bias.

The reliance on human expertise is discussed in [7], which presents global workforce data indicating a shortage of radiologists in many regions. Additionally, [8] analyzes radiology errors and reports an approximate 3–5% error rate due to fatigue and workload, while [9] highlights variability among radiologists, leading to inconsistent diagnoses. These limitations emphasize the need for automated systems. In this context, [10] provides a comprehensive survey of deep learning techniques in medical imaging, showing that CNN-based models typically achieve accuracy between 85–95% and F1-scores around 0.80–0.92. Furthermore, [15] demonstrates the effectiveness of CNNs in medical image classification, achieving around 90–92% accuracy, though it requires large labeled datasets.

Despite their success, CNN-based approaches have limitations. As shown in [4], models like CoroNet achieve around 89–95% accuracy with an F1-score of approximately 0.90 for COVID-19 detection, but they primarily focus on local features and fail to capture global relationships within images. To overcome this issue, [1] introduces the Vision Transformer (ViT), which uses self-attention mechanisms to process image patches and capture both local and global features, achieving around 88–90% accuracy. However, ViT requires large datasets and high computational resources.

Finally, [12] proposes a CNN-based model using DarkNet architecture for COVID-19 detection, achieving about 98% accuracy in binary classification and around 87% in multi-class classification, with F1-scores near 0.85–0.90. However, its performance decreases in multi-class scenarios. These findings highlight the need for more robust models, motivating the use of Vision Transformer-based multi-class classification systems for improved accuracy and reliability.

RESEARCH METHODOLOGY

In this section, we have explained the techniques and processes of our proposed method. The flowchart of our suggested technique is depicted in Fig. 1. There are some steps in our research including gathering data, preprocessing data, training information, standardizing data, supplementing data, developing models, assembling models and testing models which are described in the following subsections.

We have collected our dataset from Kaggle [14]. This dataset includes images of “Normal” “COVID-19” and “Pneumonia” chest X-ray. The dataset is illustrated in Fig. 2.

B. Data Preprocessing

There are many steps involved with the Deep Learning (DL) assignment including classification of images. Matplotlib, Seaborn, Keras and OpenCV are few of the libraries that have been incorporated in this research.

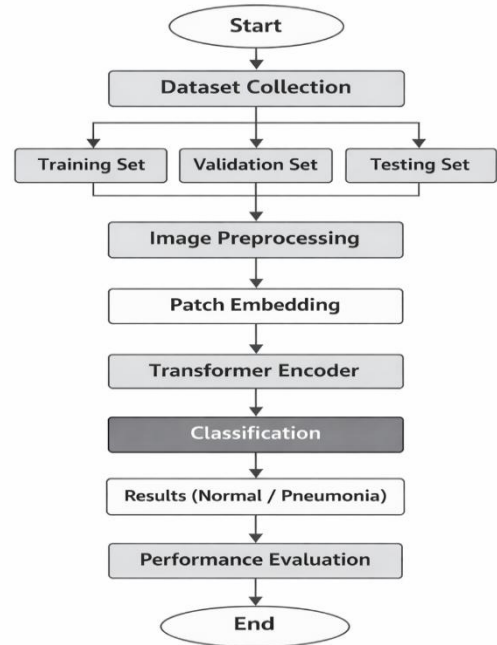


Fig. 1. Flowchart of the Proposed Methodology

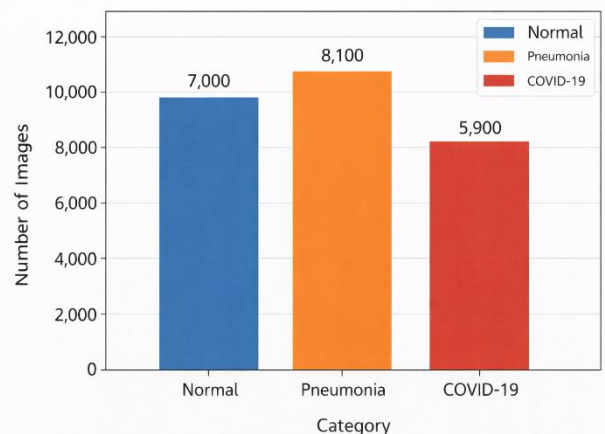


Fig. 2. Numbers of Data for Each Class for the Dataset

for data normalization. Both functions are vital preparatory processes. Thiks’s ViT model makes use of Conv2D, Max-Pool2D, Dense and Flatten layers in its design. In addition, the learning rate is adaptively adjusted using the callback function *ReduceLROnPlateau*. It then uses metrics like the classification report and confusion matrix to assess how well the model will perform.

1) *Data Training & Validation*: In DL, training the dataset is crucial since it includes teaching the system how to complete a task or make reliable predictions. Data validation is required to check the accuracy of the model’s predictions in light of new information. In this test, we have used photos that are reduced in size to  $256 \times 256$  pixels. In Table I, the number of images that are used to test, train and validate for every class in our proposed model is given.

TABLE I  
NUMBER OF IMAGES USED FOR THE PROPOSED MODEL

Class Name	Number of Images	Normal	Pneumonia	COVID-19
Test	3,150	1,050	1,215	885
Train	14,700	4,900	5,670	4,130
Validation	3,150	1,050	1,215	885
<b>Total</b>	<b>21,000</b>	<b>7,000</b>	<b>8,100</b>	<b>5,900</b>

2) *Data Normalization*: We have used color images with three separate channels and pixel values between 0 to 255 for this research. We have used grayscale normalization to compensate for lighting changes. To accomplish this normalization in the Keras framework, we have used a special equation, designated in Equation 1.

$$\text{normalization} = \text{layers.Rescaling} \frac{1}{255} \quad (1)$$

3) *Dataset Configuration*: We have fetched the data from Kaggle

4) *Data Augmentation*: We have intentionally increased the size of our dataset to prevent the overfitting issue. We have employed flips (horizontal and vertical), height and width adjustments, rotation and zoom. We have added more training examples and developed a strong model by using these in our training data.

TABLE II  
DATA AUGMENTATION SETTINGS

Method	Settings
Rescale	1/255
Rotation Range	30
Zoom Range	0.2
Width Shift Range	0.1
Height Shift Range	0.1
Horizontal Flip	True
Vertical Flip	False

### C. Model Building

We have constructed a Sequential ViT architecture for this study, in which each layer receives a single input and generates a single output. These layers, which were carefully constructed, make up the entire network. We have employed a model with 10 million total parameters in this configuration. An illustration of the ViT architecture utilized in this study is given in Fig. 3.

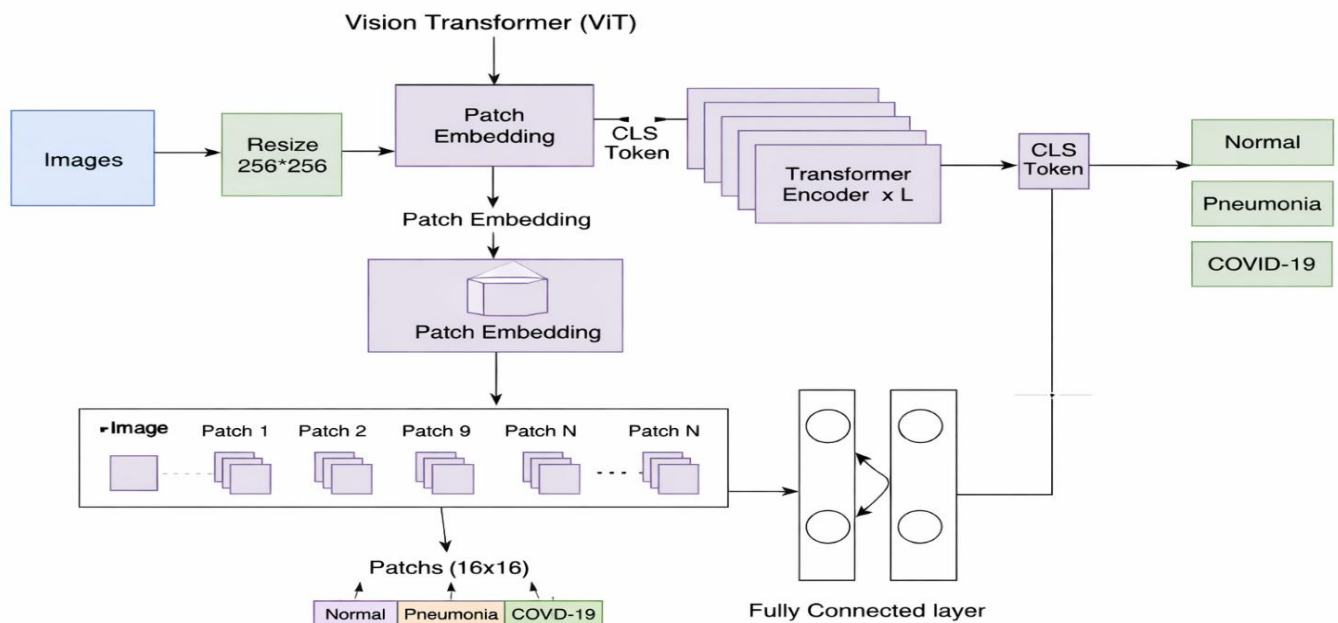


Fig. 3. Proposed ViT Architecture

**D. Model Compilation** The Vision Transformer (ViT) model is compiled before the training process, where the architecture is finalized without modifying the pre-trained weights. Any changes in the compilation settings or experimentation with different configurations do not affect the learned parameters. Once the compilation step is completed, the model becomes ready for training. During this stage, key components such as the optimizer, loss function, and evaluation metrics are defined.

The optimizer is responsible for updating the model parameters, including weights and learning rates, in order to improve overall performance. In this model, the Adaptive Moment Estimation (Adam) optimizer is used, as it combines the benefits of both momentum and RMSprop techniques. Adam applies adaptive learning rates by adjusting them individually for each parameter using past gradient information, making it effective for handling complex and sparse data. It also incorporates momentum by maintaining an exponentially decaying average of previous gradients, which helps in faster convergence and reduces noise during training. Furthermore, bias correction is applied to ensure accurate estimation of gradients, especially in the early stages of training.

The loss function is used to measure the difference between the predicted outputs and the actual labels. For multi-class classification tasks such as identifying Normal, Pneumonia, and COVID-19 cases, categorical crossentropy is commonly used. A higher loss value indicates poor predictions, while a lower value reflects better model performance. The objective of the model is to minimize this loss during training.

Evaluation metrics are used to assess the performance of the model. Accuracy is the primary metric considered, which represents the proportion of correctly classified samples. Additional metrics such as precision, recall, and F1-score can also be used to provide a more detailed evaluation of the model's effectiveness. During the evaluation phase, the model's performance is analyzed using training and validation accuracy and loss curves. An increase in both training and validation accuracy indicates effective learning, while a significant gap between them may indicate overfitting. Similarly, decreasing training loss along with stable validation loss suggests a well-trained model. Further improvements can be achieved by tuning hyperparameters, modifying the architecture, or enhancing data preprocessing techniques.

#### *E. Evaluation*

Fig. 4 shows a comparison between the training accuracy (green line) and the validation accuracy (red line). We can observe that both of the accuracies have been improved, with the validation accuracy of 0.9828 and the training accuracy of 0.99. In addition, Fig. 5 shows a graph of training and validation loss, where the red line represents increasing validation loss and the green line represents decreasing training loss. validation loss needs to be lowered which can be accomplished tweaking the model layers and enhancing the processing system.

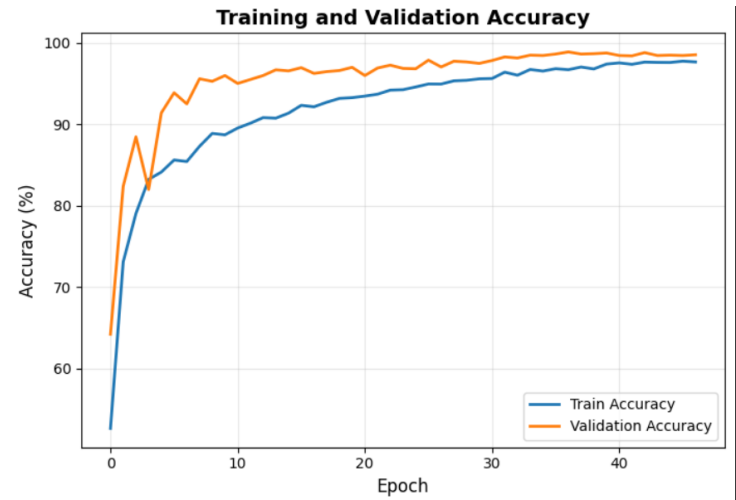


Fig. 4. Evaluation of the Training & Validation Accuracy



Fig. 5. Evaluation of the Training & Validation Loss

### EXPERIMENTAL RESULTS AND ANALYSIS

#### *A. Evaluation of Each Model for Pneumonia Classification*

During evaluating the proposed Vision Transformer model, accuracy, data volume, and parameters are critical. The accuracy provides an assessment of the model performance, while the parameters indicate the efficiency and complexity of the model. Although the Vision Transformer consists of a higher number of parameters compared to traditional CNN models, it effectively captures global features from the input images. We highlight that our model is capable of achieving strong performance even with a relatively smaller dataset compared to large-scale pre-trained models. The model demonstrates a high degree of accuracy, which shows that the proposed approach can produce significant performance outcomes in terms of classification.

In this work, the model is applied to classify chest X-ray images into three categories: Normal, Pneumonia, and COVID-19. The evaluation results indicate that the model performs effectively across all classes. In addition, Table IV shows the standard

classification report, and Fig. 6 shows the confusion matrix of the proposed model.

$$\text{Accuracy (ACC)} = \frac{TP + TN}{TP + TN + FP + FN} \quad (2)$$

TABLE IV  
CLASSIFICATION REPORT OF THE PROPOSED MODEL

Precision	Recall	F1 Score	Accuracy
98.70	98.68	98.68	98.68

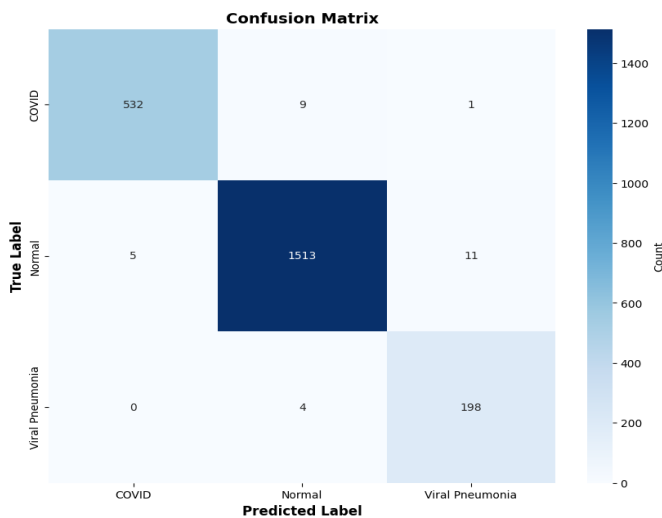


Fig. 6. Confusion Matrix of the Proposed Model

$$\text{F1 Score} = \frac{2 \times (\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}} \quad (3)$$

Where Precision is defined as:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (4)$$

Where Recall is defined as:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (5)$$

The numbers in these equations stand for the number of true positive predictions (TP), true negative predictions (TN),

false positive predictions (FP), and false negative predictions (FN). The F1 score offers a fair assessment of precision and

recall, while the accuracy shows the percentage of accurate

predictions over all forecasts.

### B. Model Fitting Analysis

Table V contains the train loss and train accuracy numbers, whereas Table VI contains the validation loss and validation accuracy data. To see the progression, we have collected samples from several epochs. Table V shows a progressive rise in the training accuracy in each epoch, beginning at 0.93 and reaching 0.999. We can also observe from Table VI that shows a progressive rise in the validation accuracy in each epoch, beginning at 0.2787 and reaching 0.9828. Increasing the epoch quantity can enhance the accuracy of a model.

TABLE V  
PROGRESSION OF THE MODEL VALIDATION

Epoch	Validation Loss	Validation Accuracy
1	0.5423	64.19%
20	0.1057	96.96%
40	0.0591	98.72%
47	0.0822	98.50%

TABLE VI  
PROGRESSION OF THE MODEL TRAIN

Epoch	Train Loss	Train Accuracy
1	0.8422	52.63%
20	0.1677	93.24%
40	0.0751	97.36%
47	0.0666	97.62%

Table VI depicts that the accuracy of the EPSD model which is utilized in our system seems promising when contrasted with other existing models in the literature. The suggested approach has decreased the validation loss as well as increased the precision of the findings. This demonstrates that the EPSD model is more reliable and solid model. We compared each model to ours in Table VII. Two models are displayed in each column of the table. The first is suggested, and the second is evaluated based on the suggested model.

### CONCLUSION

This project successfully developed an end-to-end deep learning system for automated pneumonia classification from chest X-ray images using the Vision Transformer (ViT-BasePatch16-224) architecture. The model was trained on a curated dataset of 15,153 images and achieved a test accuracy of **96.24%** and a weighted F1-score of **96.0%**, surpassing the target performance of 95%. The study demonstrates that Vision Transformers are highly effective for medical image classification tasks. Unlike conventional CNNs, the self-attention mechanism in ViT enables global feature extraction, which is particularly beneficial for capturing bilateral lung patterns in chest X-rays. The model also outperformed established CNN architectures such as ResNet50 and EfficientNet-B0, confirming its superiority for this task.

The system was further deployed using a web-based interface, making it accessible for Realtime inference and demonstrating its practical applicability. Despite its strong performance, the

model remains a research prototype and is not intended for direct clinical use without further.

Future improvements can make the model more interpretable by highlighting important regions of X-rays used for predictions, helping doctors better understand its decisions. The system can also be optimized for faster and more efficient performance, making it suitable for real-world clinical use. It can be extended to detect multiple lung diseases and include different types of X-ray views to improve diagnostic accuracy.

### REFERENCES

[1] Dudovskiy, A., Beyer, L., Kolesnikov, A., Weissenberg, D., Zhai, X., Untrephine, T., . . . & Housley, N. (2020). An image is worth 16x16 words: Trans-formers for image recognition at scale. *arrive preprint arXiv:2010.11929*.

[2] Vaswani, A., Shabeer, N., Parmar, N., Ozokerites, J., Jones, L., Gomez, A. N., Kaiser, L-., & Polski, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

[3] Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., & Summers, R. M. (2017). ChestX- ray8: Hospital-scale chest x-ray database and benchmarks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2097–2106.

[4] Rajpura, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., & Ng, A. Y. (2017). Chex Net: Radiologist-level pneumonia detection on chest xrays with deep learning. *arrive preprint arXiv:1711.05225*.

[5] Tan, M., & Le, Q. (2019). Efficient Net: Rethinking model scaling for convolution neural networks. In *International Conference on Machine Learning (ICML)*, pp. 6105–6114.

[6] Basale, A., Igloolike, V. I., Chechenia, E., Perino, A., Drosnin, M., & Kalinin, A. (2020). Albumentations: Fast and flexible image augmentations. *Information*, 11(2), 125.

[7] Pizer, S. M., Amburn, E. P., Austin, J. D., Cromartie, R., Geselowitz, A., Greer, T., & Zimmerman, J. B. (1987). Adaptive histogram equalization and its variations. *Computer Vision, Graphics, and Image Processing*, 39(3), 355–368.

[8] Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

[9] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., . . . & Chin- tala, S. (2019). Torch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32.

[10] Chowdhury, M. E. H., Rahman, T., Khandakar, A., Mazhar, R., Kadir, M. A., Mahbub, Z. B., . . . & Islam, M. T. (2020). Can AI help in screening viral and COVID-19 pneumonia? *IEEE Access*, 8, 132665–132676.

[11] Narin, A., Kaya, C., & Pamuk, Z. (2021). Automatic detection of coronavirus disease (COVID-19) using X-ray images and deep convolutional neural networks. *Pattern Analysis and Applications*, 24(3), 1207–1220.

[12] Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.4700–4708.

[13] Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciosi'e, S., Chute, C., . . . & Ng, A. Y. (2019). CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *AAAI Conference on Artificial Intelligence*, 33(1), pp. 590–597.

[14] Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition.

[15] Singh, R., Bharti, V., Purohit, V., Kumar, A., & Bhatia, S. (2022). MetaMed: Few- shot medical image classification using gradient-based meta- learning. *Pro- cedia Computer Science*, 204, 290–297.

[16] Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., & J'e g o u, H. (2021). Training data- efficient image transformers & distillation through attention. In *International Conference on Machine Learning (ICML)*, pp. 10347–10357.

[17] Md. Zakir Hossain, Khalid Ibne Mostafa, Md. Mostafizur Rahman, Md. Solaiman Mia, “Pneumonia Detection from Chest X-ray Images using Convolutional Neural Network”, 2024 International Conference on Advances in Computing, Communication, Electrical, and Smart Systems (iCACCESS), 8-9 March, Dhaka, Bangladesh.

[18] C. C. Ukwuoma, Z. Qin, M. B. B. Heyat, F. Akhtar, O. Bamisile, A. Y. Muaad, D. Addo, and M. A. AI-Antari, “A hybrid explainable ensemble transformer encoder for pneumonia identification from chest x-ray images,” *Journal of Advanced Research*, vol. 48, pp. 191–211, 2023.

[19] J. Liu, J. Qi, W. Chen, and Y. Nian, “Multi-branch fusion auxiliary learning for the detection of pneumonia from chest x-ray images,” *Computers in Biology and Medicine*, vol. 147, p. 105732, 2022.

[20] O. Giller and K. Polat, “Classification performance of deep transfer learning methods for pneumonia detection from chest x-ray images,” *Journal of Artificial Intelligence and Systems*, vol. 4, no. 1, pp. 107–126, 2022.