# A Hybrid Deep Learning Model for Accurate Chest X-Ray Disease Classification

Abu Mukaddim Rahi
*Electrical and Computer Engineering*
*North South University*
*Bashundhara R/A, Dhaka, Bangladesh*
mariam.binte@northsouth.edu

Mariam Binte Bashir
*Electrical and Computer Engineering*
*North South University*
*Bashundhara R/A, Dhaka, Bangladesh*
mariam.binte@northsouth.edu

Maher Ali Rusho*
*Computational Material and Data Analytics*
*Senior Scientist, Mr. R Business Corporation (NGO)*
*Chennai, Tamil Nadu, India*
maher.rusho@colorado.edu

Md. Khurshid Jahan
*Electrical and Computer Engineering*
*North South University*
*Bashundhara R/A, Dhaka, Bangladesh*
khurshid.jahan@northsouth.edu

Saber Hossain
*Electrical and Computer Engineering*
*North South University*
*Bashundhara R/A, Dhaka, Bangladesh*
saber.hossain@northsouth.edu

*Abstract*—To begin with, the classification of chest X-rays is vital for automatically diagnosing respiratory diseases such as COVID-19, pneumonia, and other abnormalities. In our system, we operated custom convolutional neural networks (CNNs) and Vision Transformers (ViTs), including Tiny-ViT, Swin-Transformer, FocalNet, and VOLO, to compare their effectiveness and for the classification task under limited data constraints in the Covid19, Pneumonia, and Normal Chest X-Ray Images dataset. Moreover, advanced data augmentation and regularization techniques have been used to improve the strength and inference of our system. Consequently, we attained an accuracy of 98.29% on the custom CNN and 98% on the Vision transformer model VOLO, outperforming all other models. In addition, our selection of custom CNN and ViT models (VOLO-D1) was based on their intense feature extraction abilities and usefulness for transfer learning. Subsequently, we used FocalNet to replace self-attention (SA) with a focal modulation mechanism in our system's vision. Furthermore, we reached 98%, surpassing the state of art images.

*Index Terms*—disease, classification, CNN, ViT, model

Fig. 1. Chest X-ray images from dataset

## I. INTRODUCTION

Respiratory diseases like COVID-19 and pneumonia pose a substantial global health burden, for which early and proper diagnosis is essential. The World Health Organization stated about the outbreak of the public health emergency on 30 January 2020 [1], and levied the attack as a pandemic on 11 March [2]. Its symptoms range from symptomless to deadly, including fever, painful throat, night cough, and tiredness. [3]. Convolutional neural networks (CNNs) have shown outstanding results in fixing various machine-learning topics with multiple layers of architecture. In our system, CNNs and vision transformer models play a crucial role. Consequently, this process emphasizes the need for computational aids and strategies to expand artificial intelligence's image recognition and classification field [4]. Meanwhile, vision transformers offer self-attention tools,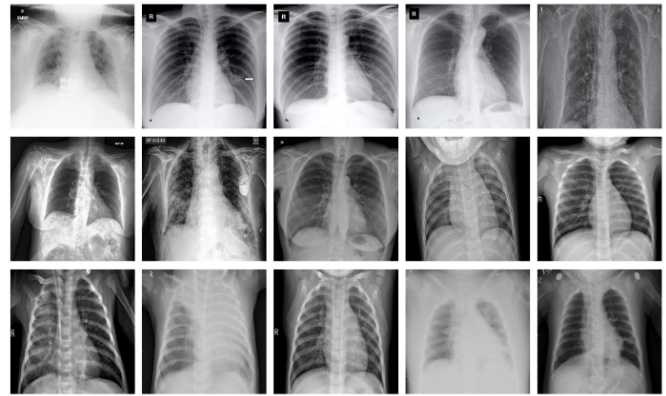 a novel approach to our image analysis by capturing international contextual information. Figure 1 shows the image classification of Normal, COVID-19-infected, and pneumonia-infected with the X-ray image of the Chest.

In this project, we proposed an AI-powered system for classifying chest X-ray images into three categories, "COVID," "PNEUMONIA," and "NORMAL," using the custom CNN and transformer models. The following are the contributions of our project :

- **Integrating CNNs and Vision Transformers in classifying medical images:** Our project merges with CNNs and vision transformers to achieve high classification accuracy and strengthen state-of-the-art architecture.
- **Ensuring Reliability with Data Preprocessing and Augmentation:** Our system provides reliable performance across diverse datasets by applying preprocessing and data augmentation.
- **Bridging Investigation and Real-World Applications**

**in Healthcare:** By bridging the gap between academic research and practical life applications in healthcare.

Our research is organized as Section II with a Literature review, Section III with Proposed System, Section IV with Model Architecture, and Section V with Results and Analysis. and Section VI, with the conclusion and future scope

## II. Literature Review

Regarding the customized convolutional neural network (CNN), S. Ashwini et al. [5], propose a multi-type classification of some diseases, such as lung opacity, tuberculosis, pneumonia, and COVID-19. They used two models, Classification-1 and Classification-2, to detect lung diseases and different types of lung diseases. They achieved 99.82% accuracy in Classification-1 and 98.75% accuracy in Classification-2. This will help to recognize and treat patients more effectively. Consequently, Mohammad Mousavi et al. [6], offer an automatic detection model for COVID-19 operating the respiratory sound and the medical image based on the Internet of Health Things (IoHT). Notably, they used the sound of coughing to detect COVID-19-affected patients, which achieved an accuracy of 94.999%. Among the models, they achieved an amazing result of 99.414% from InceptionResNetV2. Though it can diagnose the initial status of COVID-19, their system might also include other lung diseases. Likewise, Ibnu Utomo Wahyu Mulyono et al. [7], provide different arrangements of convolutional neural networks (CNN) used for image classification tasks in COVID-19. They analyzed the performance of VGG, ResNet-50, and classic CNN architectures with various datasets. Among them, the ResNet-50 architecture achieves the highest performance with an accuracy of 96.63%. Puji Dwi Rinanda et al. [8] also presented a practical approach to automatically recognize and classify mango leaf diseases by the Convolutional Neural Network (CNN). However, they performed a comparative analysis of the accuracy between VGG16, CNN, and InceptionV3. Among the three models, VGG16 outperformed with an accuracy of 96.87%. Despite impressive results, they might improve their dataset to improve model accuracy. Saravanan Srinivasan et al. [9], highlighted a deep convolutional neural network (CNN) to improve the early detection of brain tumors. Again, for the different classes of classification tasks, they offered three separate CNN models with an accuracy of 99.53%, 93.81%, and 98.56%, respectively. Md Nurul Absur et al. [4] illustrates the need to analyze digital images comparable to digital media pictures by leveraging CNN for the computer system. They achieved 98.71% accuracy using the MNIST dataset without any bias. As a result, Deep CNNs require less prior work than other image-processing algorithms.

Based on the transformer model, Asmi Sriwastawa et al. [10] CNNs have been the most prevalent image classification mechanism. They achieved the finest transformer-based classifier, with 91.57% test accuracy on the BreakHis on MaxViT. Subsequently, Attiapo Acybah Morel Omer et al. [11], introduces an approach to classify images operating Vision Transformer (ViT) architecture. Additionally, ViT emerged as an ideal option for CNN for image analysis tasks with improved performance, which can process image patches instantly without depending on spatial orders and enhance computational efficiency. Likewise, Mouhamed Laid ABIMOULOUD et al. [12] highlighted three low-weight systems on attention and convolution techniques: ViT, MVIT, and CCT. They used the BreakHis dataset for binary and multi-classification of breast cancer subtypes, resulting in fewer parameters and lower training time while achieving accurate breast tumor classification. Among the models, VIT obtains the highest accuracy of 98.64% and was compared with state-of-the-art models using the same dataset, which can minimize computational training resources and decision time. In this paper, Verren Angelina Saputra et al. [13], compare ResNet152 as the best CNN model for classifying skin diseases with ViT. They used the HAM10000 dataset with 10,015 images, where ViT achieved 98.28% accuracy, more than ResNet152's 96.70%. Though ViT proved superior to Resnet50, it has major drawbacks in potential overfitting.

## III. Proposed System

The study offers an image classification framework by comparing the arrangement of custom convolutional neural networks (CNN) and transformer-based architectures for the X-ray images utilizing custom CNN and transformer-based models, differentiating between COVID-19, pneumonia, and typical chest X-ray images employed. The image has been resized into 256x256 pixels for the custom CNN model and 224x224 pixels for the vision transformer model, which has been strategically put between every pair of convolutional coatings with a max-pooling layer of 2×2 to enhance quality extraction and decrease dimensionality. Figure 2, shows the method from the pre-trained model to model result visualization with all the steps by random flipping, rotation, zoom, brightness, and contrast adjustments. Lastly, we train and evaluate our system with the custom CNN, FocalNet, and Vision transformer models

**Dataset:** The dataset, which contains chest X-ray images, has been classified into COVID-19, NORMAL, and PNEUMONIA. Firstly, the images are organized in class-specific directories, split into 1626 pictures for "COVID", 1802 images for "NORMAL", and 1800 for the "PNEUMONIA" dataset. In our system, all the images are preprocessed and resized to 256x256 in the PNG form. Consequently, to improve generality, to form a vigorous foundation, and to detect COVID-19 for the classification method, each image undergoes processing and augmentation [14].

**Image Preprocessing:** Both CNN and ViT frameworks have been implemented to confirm that the data provided in the models is clean, constant, and optimized for training and preprocessing. For the custom CNN model, images were resized to 256x256 pixels, and the pixel values were standardized to a range of [0,1] using rescaling processes. However, the images were resized to 224x224 pixels for Vision Transformer models to align with common pre-trained model measurements. Consequently, mean and standard deviation
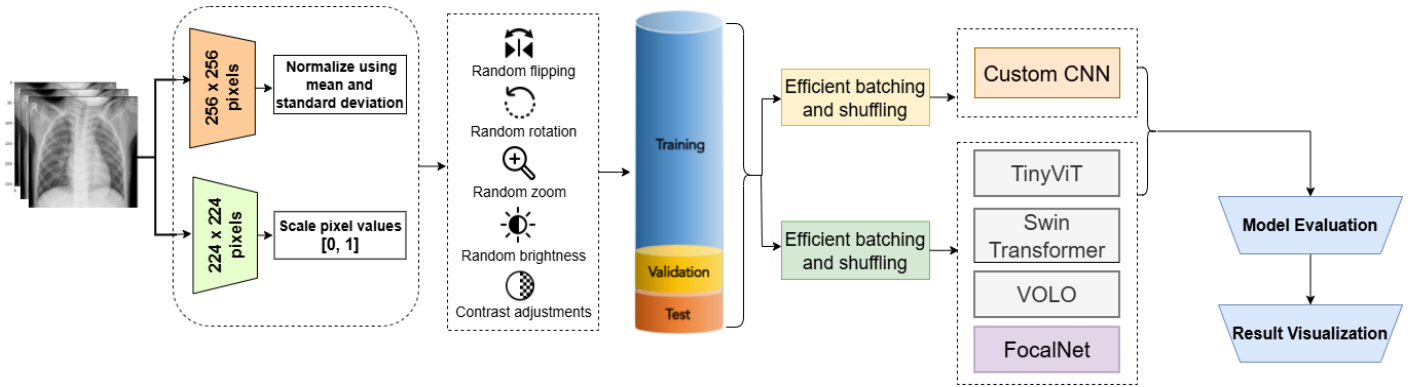
Fig. 2. Workflow of the proposed hate speech detection system

normalization was performed from the ImageNet dataset, and pictures were transformed into tensors to confirm compatibility and efficient GPU enumeration. As a result, both frameworks assured consistency across input data, enhanced model confluence.

**Image Loading:** Image loading was streamlined using the TensorFlow directory API for our custom CNN model, enabling automated labeling, efficient batch processing, and reduced training time by loading data asynchronously. In contrast, for the ViT model, a custom dataset class was implemented in PyTorch for manual image loading. Both methods support efficient data handling, batch processing, and prefetching to optimize training time. To ensure consistency in image dimensions and quality, images were loaded using the Python Imaging Library (PIL), converted to RGB format, and passed through resizing and normalization.

**Data Augmentation:** For instance, data augmentation plays an excellent role in increasing the diversity in the chest X-ray images. In our system, Random Flipping, Random Zoom, Random Brightness, and Random Contrast have been used by ±20%,±20%, ±10%, and ±10% to flip the images horizontally and vertically, to alter the magnification, to adjust the image's brightness and to handle differences in imaging equipment respectively. Table 2 shows the augmentation techniques that introduce random flipping, rotation, zoom, brightness, and contrast adjustments, with descriptions and parameters.

TABLE I
IMAGE AUGMENTATION TECHNIQUES AND PARAMETERS

| Technique | Description | Parameters |
|---|---|---|
| Random Flipping | Flips images horizontally and vertically to account for orientation | Horizontal and Vertical |
| Random Rotation | Rotates images within ±20% to handle misalignments | ±20° rotation angle |
| Random Zoom | Zooms in/out to simulate different magnification levels | Height and Width ±20% |
| Random Brightness | Adjusts brightness to mimic varying lighting conditions | ±10% brightness adjustment |
| Random Contrast | Changes contrast to address intensity variations | ±10% contrast adjustment |

**Data Set Allocation:** For instance, the dataset is separated into training, validation, and testing with an 80:10:10 ratio. Moreover, 80% of the training dataset's images were used to optimize the model's weight. During the training, Validation data, with 10% of the dataset, is employed to monitor overfitting and to fine-tune hyperparameters. To ensure that final evaluations were completed, the remaining 10% were reserved for testing.

## IV. MODEL ARCHITECTURES

**Custom CNN Model:**Our system was built on the custom CNN model with an attribute extraction and classification channel. For the input sizes with 256x256 pixels, our system resized the layer and then utilized the convolutional layers with kernel measures of (3x3), tracking batch normalization and ReLU activation processes.

Additionally, to improve the model's capability to generalize the unrecognized data, L2 regularization was used for the dense layers with 0.001. It employed the Adam optimizer, using data covering, shuffling, and prefetching with TensorFlow's autotune component to optimize the training further. Figure 3 shows the architecture diagram for the Custom CNN.

**VisionTransformers (ViT):** Our system used transformer models like TinyViT, Swin Transformer, VOLO, and FocalNet. Each brought individual strengths to the classification task as our dataset experienced extensive preprocessing, including image resizing to 224x224 pixels, normalization, and data augmentation techniques like random cropping, horizontal flipping, and brightness adjustments. We separated the pictures into smaller image patches using patch embedding coatings per transformer model. These patches were then linearly projected into token embeddings, comprising sequences provided into the transformer layers. In addition, our system allows the network to concentrate on different parts of the image simultaneously, featuring numerous self-attention heads. Finally, multi-layer perceptron (MLP) heads are employed with softmax activation to produce probabilities for the three target categories. Lastly, the AdamW optimizer with a learning speed of 1e-4 and a significant deterioration of 1e-4 was utilized for
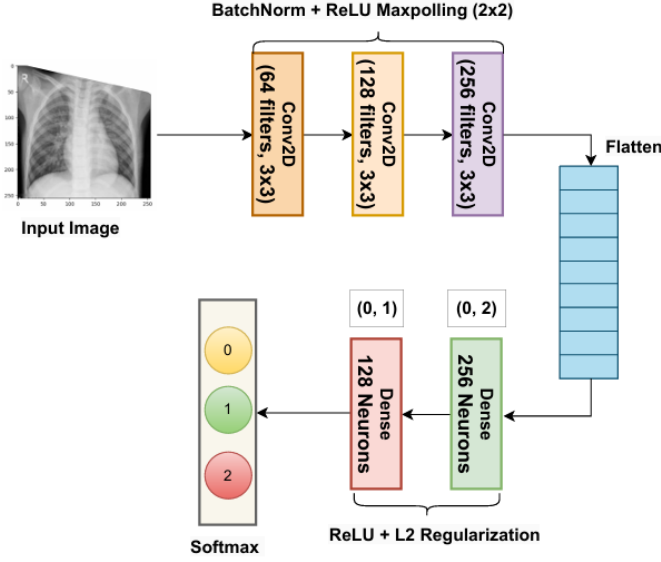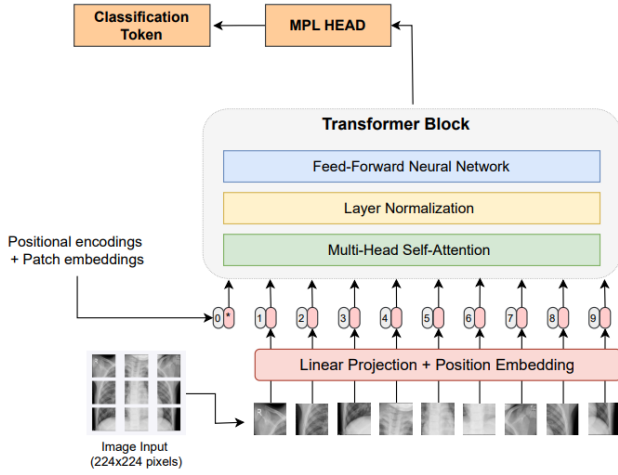
Fig. 3. Architecture diagram for Custom CNN.

| Feature | Custom CNN | TinyViT | Swin Transformer | VOLO | FocalNet |
|---|---|---|---|---|---|
| Input Image Size | 256x256x3 | 224x224x3 | 224x224x3 | 224x224x3 | 224x224x3 |
| Core Building Block | Conv2D Layers | Transformer Blocks | Swin Transformer Blocks | Vision Outlooker | Focal Modulation |
| Attention Mechanism | N/A | Self-Attention | Shifted Window Attention | Outlooker Attention | Focal Attention |
| Number of Layers | 3 Conv Blocks + Dense Layers | 12 Transformer Layers | 12 Transformer Layers | 19 Transformer Layers | 16 Transformer Layers |
| Nodes per Layer | Conv2D: 64, 128, 256. Dense: 256, 128, 3 | Embedding: 192 Attention. Heads: 3. | Embedding: 96 Attention. Heads: 3/6/12/24 | Embedding: 192 Attention Heads: 12 | Modulation: 192 Expansion Factor: 4 |
| Normalization | Batch Normalization | Layer Normalization | Layer Normalization | Layer Normalization | Layer Normalization |
| Dropout | 0.2, 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |
| Learning Rate | 1e-4 | 1e-4 | 1e-4 | 1e-4 | 1e-4 |
| Batch Size | 32 | 4 | 4 | 4 | 4 |
| Parameters (Approx) | 1M | 5M | 28M | 30M | 19M |

## V. RESULT AND ANALYSIS

**Model Performance:** For instance, we used Convolutional Neural Networks (CNNs) and advanced Vision Transformers (ViTs) to classify chest X-ray images into "COVID," "NORMAL," and "PNEUMONIA" categories. Sequentially, we got the CNN model results with a precision of 0.9829, a recall of 0.9829, and an f1-Score of 0.9829. Additionally, the TinyViT model achieved a precision of 0.9640, a recall of 0.9640, and an f1-Score of 0.9640. In addition, the Swin Transformer was followed by scores of 0.9722 (Precision), 0.9720 (Recall), and 0.9721 (F1-Score). Among all the models, the VOLO architecture emerged as the best-performing model with a precision of 0.9802, recall of 0.9800, and an f1-Score of 0.9800, showcasing its effectiveness in both the training and testing phases. Besides that, we obtained a precision of 0.9769, a recall of 0.9760, and an f1-Score of 0.9761 in the FocalNet model, confirming its robustness in complex image classification tasks. Across all evaluation metrics, our performance highlights consistent performance with reliability in real-world medical diagnostic applications. Figure 5 represents a visual representation of the Precision, Recall, and F1-score metrics across all evaluated models.

**Performance Visualization:** In Figure 6, the graph shows us the training and validation accuracy and loss curves over 20 epochs for the custom CNN model. This graph shows us the confluence of the model achieving high training accuracy with minimal overfitting, where the close alignment of training and validation curves indicates that the model generalizes well
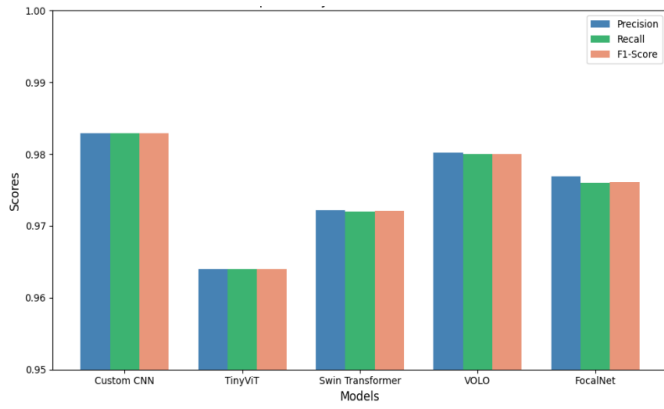


Fig. 4. Architecture Diagram for Transformer model.

model optimization. Figure 4 shows the architecture diagram for the Transformer model.

**Comparison of Models Architecture:** The custom CNN model is an interpretable architecture with completely attached layers. Additionally, we operated the swim transformer to counteract local and global feature learning with ˜28M parameters. Besides that, FocalNet, with 19 M parameters, has been used to achieve a balanced focus with increased computational efficiency. As a result, each model illustrates a trade-off between intricacy, efficiency, and quality extraction abilities. Table 4 shows the architectures executed in this project with the Convolutional Neural Network (CNN), the Vision Transformer (ViT), and the FocalNet Model highlighting their critical features and individual characteristics.

Fig. 5. Comparison graph among the model.



Fig. 7. ViT model (VOLO-D1) training and validation accuracy and loss.

to unseen data. Furthermore, minor fluctuations in validation accuracy and loss are observed around epoch 12, likely due to slight overfitting or noise. *The* figures 6 (a) and (b) show the training accuracy, validation accuracy, training loss, and validation loss graphs for custom CNN. In Figure 7, the graph shows us the training and validation accuracy and loss curves over 50 epochs for the ViT model (VOLO-D1). In Figure, the graph shows us the training and validation accuracy and loss curves over 50 epochs for the ViT model (VOLO-D1). Consequently, in the accuracy curve of the graph, the training accuracy improves consistently, stabilizing around 98%. In contrast, the validation accuracy shows slight changes in the early epochs but eventually aligns closely with the training accuracy, reaching approximately 96%. The loss curve of graph B demonstrates a constant decrease in training loss, indicating effective learning. In contrast, the validation loss follows a similar trend, stabilizing around 0.1 after initial instabilities with good generalization and minimal overfitting. Hence, it shows substantial accuracy, effectiveness, and low loss on the graph. Figures 7 (a) and (b) show the training and validation loss and accuracy for the ViT model (VOLO-D1).
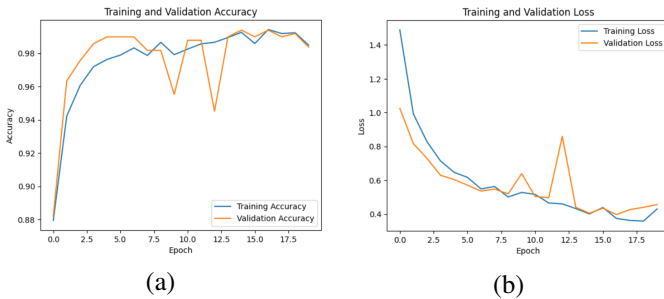


Fig. 6. Custom CNN training and validation accuracy and loss.

**Classification Result:** For the PNEUMONIA, 175 instances were correctly classified as PNEUMONIA, three instances were misclassified as "NORMAL," and 2 instances were misclassified as "COVID." However, for the "NORMAL" class, 177 instances were correctly classified as "NORMAL," two were misclassified as "PNEUMONIA," and one instance was misclassified as "COVID." Moreover, for the "COVID"
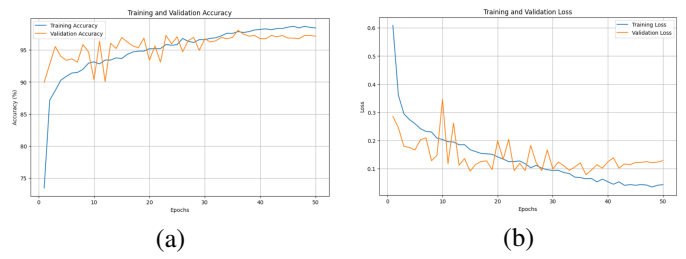
class, 182 instances were correctly classified as "COVID," two instances were misclassified as "NORMAL," and 0 instances were misclassified as "PNEUMONIA." Overall, the Custom CNN model performs very well in our system, with minimal misclassifications and a small number of "NORMAL" misclassified as "PNEUMONIA." The confusion matrix for Custom CNN is illustrated in Figure 8. Like the custom CNN,
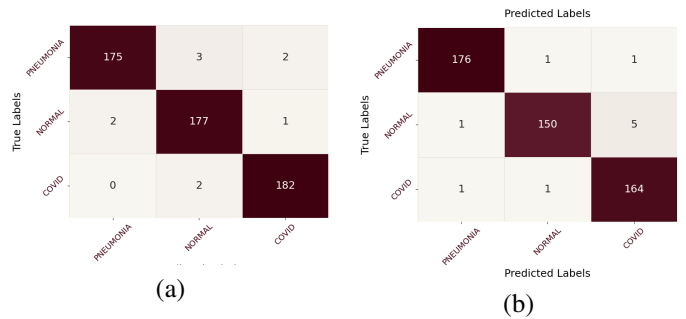


Fig. 8. Confusion matrix for the custom CNN and VOLO-d1 model.

the confusion matrix of the VOLO-D1 Model performed well. However, for the "PNEUMONIA" class, 176 instances were correctly classified as "PNEUMONIA," one was misclassified as "NORMAL," and one was misclassified as "COVID." Furthermore, for the NORMAL class, 150 instances were correctly classified as NORMAL, one instance was misclassified as PNEUMONIA, and five instances were misclassified as "COVID." Besides that, in the "COVID" class, 164 instances were correctly classified as "COVID," one misclassified as "NORMAL," and one misclassified as "PNEUMONIA." The confusion matrix for the VOLO-D1 Model is illustrated in Figure 9.

Significantly, almost all of the images were correctly classified as "PNEUMONIA", "COVID" and "NORMAL" as shown in Figure 10. Notably, we used a custom CNN model to classify the diseases.

**Accuracy Comparison Table:** Our system has five different deep learning architectures: Sequential CNN, TinyViT, Swin Transformer, VOLO, and FocalNet.However, we trained, validated, and tested our dataset categorized into "COVID", "NORMAL," and "PNEUMONIA" classes.

In the following table, we have demonstrated the strong interpretation of the custom CNN model with a training accuracy of 98.46%, validation accuracy of 98.37%, and test Accuracy of 98.29%, which offers to extract meaningful image features. However, TinyViT achieved a training accuracy of
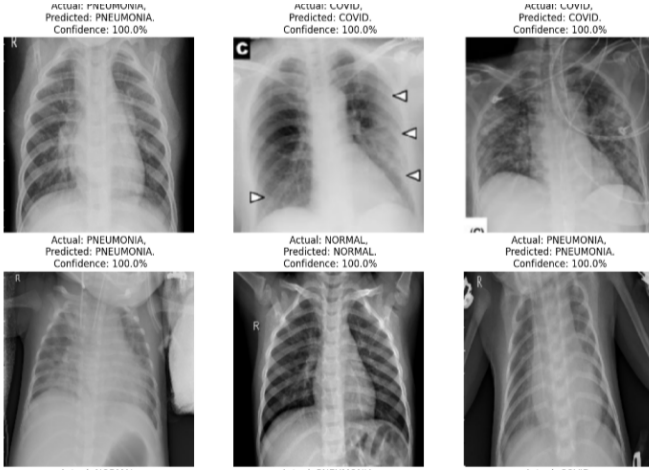
Fig. 9. Screenshot of classification of the "PNEUMONIA", " COVID" and "NORMAL" images.

98.72%, validation accuracy of 96.40%, and test accuracy of 96.40%. Again, Swin Transformer followed closely with 98.35%, 97.20%, and 97.20%, respectively. Consequently, the VOLO model excelled with a Train Accuracy of 98.43%, validation accuracy of 98.00%, and test accuracy of 98.00%, outperforming all other models across all datasets. Moreover, FocalNet achieved an excellent result with train accuracy of 99.32%, validation accuracy of 97.60%, and test Accuracy of 97.60% for medical image classification tasks. In conformity with our result, we can say that CNN remains a steadfast architecture, whereas transformer-based models like VOLO and FocalNet offer exceptional accuracy and inference capabilities. Table 5 compares the results of various models executed in our system.

TABLE III
TRAIN, TEST, AND VALIDATION ACCURACY.

| Model | Train Accuracy (%) | Validation Accuracy (%) | Test Accuracy (%) |
|---|---|---|---|
| Sequential CNN | 98.46 | 98.37 | 98.29 |
| TinyViT | 98.72 | 96.40 | 96.40 |
| Swin Transformer | 98.35 | 97.20 | 97.20 |
| VOLO | 98.43 | 98.00 | 98.00 |
| FocalNet | 99.32 | 97.60 | 97.60 |

## VI. CONCLUSION

In this paper, we offer an approach for classifying chest X-ray images into three categories: "COVID," "PNEUMONIA," and "NORMAL" where it can bridge the crucial gap in medical diagnostics. Consequently, our system used custom CNN and transformer models such as the TinyViT, Swin Transformer, VOLO, and FocalNet for automated disease detection, ensuring significant accuracy and constant outcomes. Among the transformer models, VOLO emerged as the most promising, achieving the highest accuracy of 98%. In contrast, the sequential CNN model attained an accuracy of 98.29%, which can control the critical challenges in medical diagnostics. Additionally, FocalNet achieved 98%, outperforming other models by the state of art images. We aim to expand the application by merging multiple models to enhance accuracy on other medical imaging tasks.

REFERENCES

[1] World Health Organization, "Public health emergency of international concern," 2025. Accessed: 2025-01-13.
[2] W. contributors, "Covid-19 pandemic," 2025. Accessed: 2025-01-13.
[3] World Health Organization, "Impact of covid-19 on people's livelihoods, their health, and our food systems," 2020. Accessed: 2025-01-13.
[4] M. N. Absur, K. F. A. Nasif, S. Saha, and S. N. Nova, "Revolutionizing image recognition: Next-generation cnn architectures for handwritten digits and objects," in *2024 IEEE Symposium on Wireless Technology & Applications (ISWTA)*, pp. 173–178, IEEE, 2024.
[5] S. Ashwini, J. Arunkumar, R. T. Prabu, N. H. Singh, and N. P. Singh, "Diagnosis and multi-classification of lung diseases in cxr images using optimized deep convolutional neural network," *Soft Computing*, vol. 28, no. 7, pp. 6219–6233, 2024.
[6] M. Mousavi and S. Hosseini, "A deep convolutional neural network approach using medical image classification," *BMC Medical Informatics and Decision Making*, vol. 24, no. 1, p. 239, 2024.
[7] I. U. W. Mulyono, E. H. Rachmawanto, C. A. Sari, and M. K. Sarker, "A high accuracy of deep learning based cnn architecture: classic, vggnet, and restnet50 for covid-19 image classification," *TELKOMNIKA (Telecommunication Computing Electronics and Control)*, vol. 22, no. 5, pp. 1187–1195, 2024.
[8] P. D. Rinanda, D. N. Aini, T. A. Pertiwi, S. Suryani, and A. J. Prakash, "Implementation of convolutional neural network (cnn) for image classification of leaf disease in mango plants using deep learning approach," *Public Research Journal of Engineering, Data Technology and Computer Science*, vol. 1, no. 2, pp. 56–61, 2024.
[9] S. Srinivasan, D. Francis, S. K. Mathivanan, H. Rajadurai, B. D. Shivahare, and M. A. Shah, "A hybrid deep cnn model for brain tumor image multi-classification," *BMC Medical Imaging*, vol. 24, no. 1, p. 21, 2024.
[10] A. Sriwastawa and J. A. Arul Jothi, "Vision transformer and its variants for image classification in digital breast cancer histopathology: A comparative study," *Multimedia Tools and Applications*, vol. 83, no. 13, pp. 39731–39753, 2024.
[11] A. A. M. Omer, "Image classification based on vision transformer," *Journal of Computer and Communications*, vol. 12, no. 4, pp. 49–59, 2024.
[12] M. L. ABIMOULOUD, K. BENSID, M. Elleuch, M. B. Ammar, and M. KHERALLAH, "Vision transformer based convolutional neural network for breast cancer histopathological images classification," *Multimedia Tools and Applications*, pp. 1–36, 2024.
[13] V. A. Saputra, M. S. Devi, A. Kurniawan, *et al.*, "Comparative analysis of convolutional neural networks and vision transformers for dermatological image classification," *Procedia Computer Science*, vol. 245, pp. 879–888, 2024.
[14] S. Kumar, "Covid19-pneumonia-normal chest x-ray images." Mendeley Data, V1, 2022.