

Pneumonia Detection from Chest X-ray Images Using Transfer Learning by Fusing the Features of Pre-trained Xception and VGG16 Networks

A. S. M Abdus Shafi, MD. Mareful Hasan Maruf, and Sunanda Das

Department of Computer Science and Engineering

Khulna University of Engineering & Technology

Khulna-9203, Bangladesh

abdusshafi99@gmail.com, marufhasan212@gmail.com, and sunanda@cse.kuet.ac.bd

Abstract—Pneumonia is said to be the “Silent Killer” disease caused by the infection of viruses, bacteria, or fungi in the lung alveoli. It bears an extensive risk for people, especially children in some developing nations. The ecumenic way to detect pneumonia is from Chest X-ray data. But it has some complications to diagnose pneumonia if the lung has gone through some surgery, bleeding, the superabundance of fluids, or lung cancer. So, it is necessary to take the help of Computer-Aided Diagnosis (CAD) which can collaborate the doctors to detect pneumonia. Many deep learning methods are applicable to detect pneumonia. Our research introduces a new model generated from the fusion of two different transfer learning models, the Xception model and the VGG16 model. Our research includes image pre-processing using image normalization and augmentation. We took two different transfer learning models namely Xception, and VGG16 for the feature extraction, then added some layers, made a fusion, and lastly added some extra dense layers to develop the proposed model. We took 5216 images of two classes named ‘NORMAL’ and ‘PNEUMONIA’ images to train our model. We took 5216 images to train the model in ‘NORMAL’ and ‘PNEUMONIA’ form. The results were tested with 624 images belonging to two classes. The proposed model achieved accuracy, precision, recall, and f1-score of 91.67%, 92.30%, 89.92%, and 90.87% respectively. The extensive experimental analysis demonstrates the viability of the proposed approach for various test samples.

Keywords—Classification, Pneumonia Detection, Transfer Learning Model, Xception, VGG16.

I. INTRODUCTION

Pneumonia is a disease that has to be treated properly and still has to be prevented [1]. It kills over 1.4 million people annually and it is also responsible for nearly 18% of all children under the age of five worldwide. It ranks among the deadliest diseases for children below five years of age, according to the World Health Organization (WHO, Geneva, CH) [2]. Knowing how to deploy Computer-Aided Diagnosis (CAD) effectively should make it easier for patients and medical professionals to identify various sorts of anomalies in medical images, allowing for more precise disease analysis and diagnosis [3]–[5]. Pneumonia can be diagnosed using chest X-rays, CT scans of the lungs, chest ultrasonography, needle biopsy, and chest MRI [6]. But the greatest way to identify pneumonia is using a chest X-ray. CT imaging is not as popular as X-ray imaging since it frequently takes a lot longer and may

not be possible to find enough high-quality CT scanners in many developing areas. X-rays have become the most widely used and accessible diagnostic imaging method compared to other techniques, and they are essential for healthcare delivery and epidemiological research [7], [8]. Also, various countries around the world lack skilled medical professionals and radiologists, whose estimations of such diseases are so crucial for the patient [9], [10].

Deep learning has recently become popular in a variety of fields [11]–[13]. Convolutional Neural Networks (CNNs) are supposed to be the key emerging fields for Deep Neural Networks (DNNs) and lots of deep learning accomplishments are dependent on them. To identify an image, CNN techniques show more precise results than humans nowadays. Also, Deep learning has significance in the area of artificial intelligence that can assist the medical professional with the accuracy as well as further enhance the detection of childhood pneumonia. The algorithm’s input can be the raw data and it is possible to abstract the target qualities needed for an individual task from raw data by using feature extraction and stratified sampling. The last stage includes mapping the learned characteristics to the learned attributes, and no human work is required during the whole process. The availability of enormous datasets and the development of more powerful GPUs are two key reasons that are motivating researchers to improve learning rates. Technological improvements have enabled resource-intensive models to perform as good as human practitioners for various medical imaging tasks, including diagnosis of pneumonia, detection of diabetic retinopathy, classification of skin cancer, detection of arrhythmia, identification of bleeding. Therefore, researchers’ preferred technology for image classification and diagnosis from medical images is CNN or similar deep learning techniques.

In our research, we proposed a model that can automatically determine if a patient has pneumonia or not. It is mainly the fusion of two different transfer learning models with extra dense layers. The suggested approach helps to extract characteristics from an X-ray image that automatically defines the existence of pneumonia and indicate whether pneumonia is present or not.

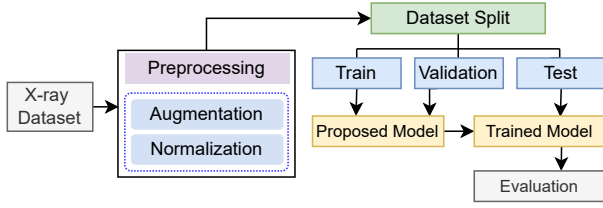


Fig. 1. Framework of the proposed method.

The following sections of this paper are organized as follows. Section II has a description of works related to our research. Section III contains an explanation of the system's whole procedure. Section IV describes the evaluation of the outcome that our system produced. In section V, we present our conclusion.

II. RELATED WORKS

As pneumonia detection is a significant issue for doctors in developing countries, several kinds of research have been performed to resolve the issue.

Jiang et al. [14] have applied a new model IVGG13 which can solve medical image recognition problems that arise when the VGG16 model is used. Their proposed model has better performance than the pre-trained VGG16 model. They applied augmentation for the pre-processing of the dataset and the training set was increased from 5216 images to 22,146 images and the test set was increased from 624 images to 1000 images. This IVGG13 model reduced the network depth more than the VGG16 model architecture and overfitting & under-fitting were avoided. They got higher accuracy for the 6000 data which was 89.1% and the highest recall for 10,000 data which was 99%.

Enes AYAN and Halil Murat ÜNVER [15] have distinguished the performance of two different CNN networks to diagnose pneumonia. They applied transfer learning including fine-tuning the parameters to train their model. They realized that VGG16 produces better performance than Xception in terms of accuracy (87%), specificity, pneumonia precision, and f1 score respectively. On the contrary, the Xception model performed better than the VGG16 model in terms of sensitivity, normal precision, and recall (94%). According to their results, the VGG16 model was better at identifying typical situations but less effective than the Xception model at identifying pneumonia cases.

For image classification of pediatric pneumonia, Liang et al. [16] opposed a persistent structure with an expanded convolution. By weighting parameters acquired from extensive datasets in the same field, a transfer learning technique was used to start the model. To categorize and diagnose pediatric pneumonia, they offered a deep learning method that combines expanded convolution with persistent thinking. According to the experimental findings of the test dataset, the method's recall rate on the problem of classifying children's pneumonia is 96.7%, and the f1-score is 92.7%.

Samir S. Yadav and Shivajirao M. Jadhav [17] investigated how to classify pneumonia using a dataset of chest X-rays and a CNN-based method. The methods they examined included

training a capsule network from scratch, transfer learning on two convolutional neural network models—the VGG16 model and the InceptionV3 model—and an LSVM classifier with local rotation variables which are also orientation-free variables. The authors thought about increasing the primary capsule layers number, increasing the capsule number in the primary capsule, averaging the results, and assembling multiple models. Then they added more Conv-Layer and evaluated the activation function.

Jain et al. [18] provided CNN models to identify viral or bacterial pneumonia with x-ray images. By adjusting different parameters, hyperparameters, and the number of convolutional layers, different convolutional neural networks were trained to distinguish x-ray images into two major categories, one is pneumonia and the other is non-pneumonia. The authors mentioned six models. The first two models, respectively, have two and three convolutional layers. The remaining models, namely VGG16, VGG19, ResNet50, and Inception-v3 are pre-trained. The validation accuracy for the first and second models is 85.26% and 92.31%, respectively.

Stephen et al. [19] proposed a customized CNN to distinguish and diagnose pneumonia from a given dataset of chest X-ray images. They constructed a CNN model from scratch for feature extraction. In the layer of feature extraction, every layer receives the output of the layer before it as input and uses that output as the input for the layer after it. Convolution layers with max-pooling and classification layers are all incorporated into the suggested architecture. Feature extractors consist of a max-pooling layer of stride 2×2 , Conv 3×3 , 32, Conv 3×3 , 64, Conv 3×3 , 128 and Conv 3×3 , 128, with a RELU activator in between them. They got 94% validation accuracy and the training loss was 0.1835.

Han et al. [20] proposed a framework for detecting pneumonia in chest X-rays that leverages radiomics features and contrastive learning. The authors extracted the features from a pre-trained ResNetAttention Model and fine-tuned the model before testing. Chhikara et al. [21] also implemented feature extraction from a transfer learning model named InceptionV3 and then fine-tuned the model as well. Images are preprocessed by utilizing filtering, equalization, gamma correction, and compression techniques before being sent to the model.

III. METHODOLOGY

The detection system of pneumonia took an image dataset from chest X-rays as input and predicts the possibility of being affected by pneumonia as the output. We have provided a flow diagram of our proposed model in Fig1.

A. Dataset

The dataset that was predominantly used in this research was made available in 2018 by Kermuny et al. [22]. The Guangzhou Women and Children's Medical Center selected 5862 x-ray images including both anterolateral and postero-lateral images that were specifically chosen from retrospective young patients under the age of 5. This dataset was split into three different sets (i.e., training, testing, and validation sets) and each set contains two groups ('NORMAL' images shown

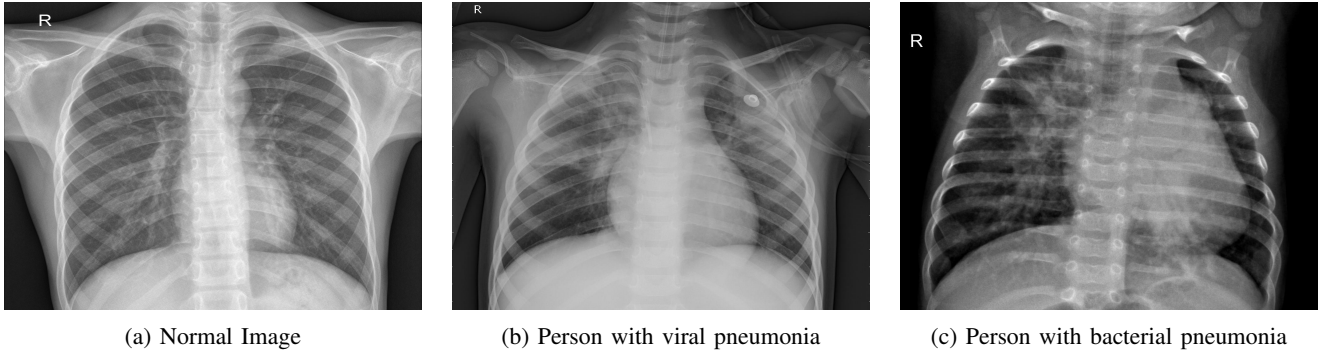


Fig. 2. Visualization of the dataset where (a) represents normal class, and (b), (c) both represent pneumonia class category.

in Fig. 2a and images with ‘PNEUMONIA’). It has two types of pneumonia diseases, one is viral pneumonia shown in Fig. 2b and the other is bacterial pneumonia shown in Fig. 2c. All chest x-ray imaging is a regular component of the patient’s clinical treatment. To increase the validation accuracy, 547 pneumonia-affected images along with 357 normal images were provided to the validation set whereas 3,339 pneumonia-affected images and 995 normal were provided to the training set. Also, 390 pneumonia-affected images and 234 normal images are provided to the testing set. Considering everything, the images are sometimes indistinct because of outside factors including the scanning place, habituation, and the history of a patient with other illnesses. In the qualitative research, we will improve the image to make it easier to interpret.

B. Image Preprocessing

The dataset we took has unbalanced positive and negative samples, which caused significantly fewer normal images than pneumonia images not only in training datasets but also in testing datasets. It can cause poor prediction accuracy. So, before training and testing the data, some image denoising was performed. Some data augmentation methods are implemented to increase image quality and size artificially. This preprocessing of data helps to prevent the problem of over-fitting and increases the stimulus generalization ability of the proposed model. Firstly, we rescaled the image and performed feature-wise normalization. The 20-degree rotation range can be regarded as the angle at which images are rotated randomly around the center of the images both clockwise and anti-clockwise which is shown in Fig. 3b and 3c. A 0.2% movement of the image was represented by the height displacement in the vertical direction in Fig. 3e. Also, the height displacement is same as the width displacement horizontally, except that the angle changed to 0.2% vertically in Fig. 3f. The size of the image will be arbitrarily increased or decreased by 0.45 times in the process described above to improve the image. Lastly, the image must be flipped horizontally as shown in Fig. 3d. Furthermore, we normalized each pixel value between 0 to 1 of the image by dividing each value by 255.

C. Proposed Model

The first stage includes mapping the learned characteristics to the learned attributes, and no human work is required during

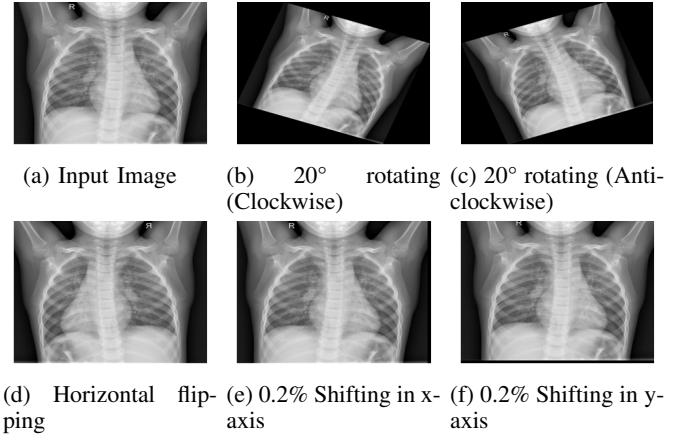


Fig. 3. Different types of augmentations used in the proposed system.

the whole process. We have applied Transfer Learning in our research. Transfer learning means transferring the parameters of a trained model to help a new model train. Feature extraction and classification of the image can perform better if large-scale well-annotated datasets are available. Therefore, transfer learning makes the best use of the ability to extract the feature on a large number of image datasets like ImageNet which can accurately represent the class boundaries. We have applied many transfer learning methods, such as ResNet-50, InceptionV3, VGG16, Xception, etc. But none of this could perform better individually. But when we made a fused model by merging two different models (The Xception model and the VGG16 model), we got a better result.

1) **Xception**: Xception is a CNN architecture that is completely based on depth-wise separable convolution layers. It initially applies the filters to every data separately before compressing the input data with 1×1 convolution all at once. It is a well-built version of Inception architecture which stands for “Xtreme Inception”. A ReLU activation follows the activation of two convolutional layer blocks in the entry flow. The number of filters, filter size (also known as kernel size), and strides are all specific according to the model. Several Distinguishable convolutional layers are also included here. Some max-pooling layers are also included. The strides are also stated when there is more than one. There are also

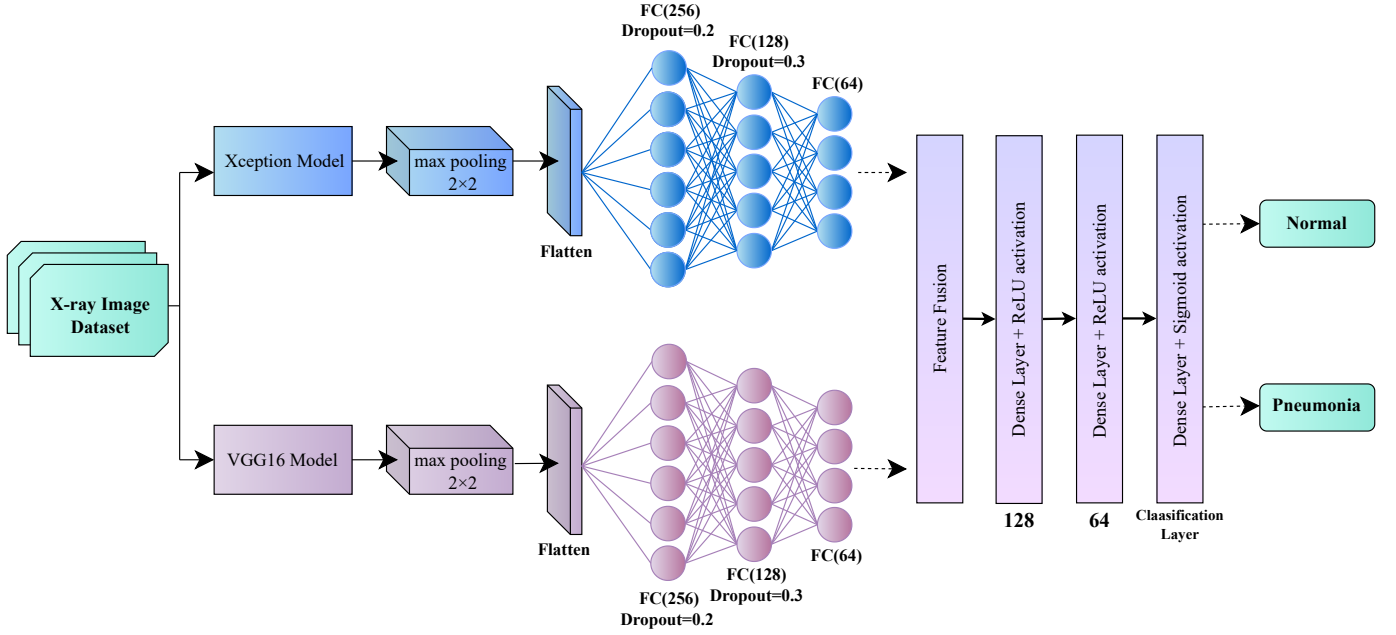


Fig. 4. The Proposed Model Architecture.

skip connections, where the two tensors are combined using 'ADD'. Additionally, it displays the input tensor's structure for each flow. For instance, if an image of size $299 \times 299 \times 3$ was used, and then it would obtain an image of size $19 \times 19 \times 728$ after the entry flow. The image size, types of various layers, the filter number, filter shape, type of pooling, amount of iterations, and option to add a fully-connected layer at the end are all included in the Middle flow and Exit flow, respectively. Additionally, batch normalization is used after each Convolutional and Separable Convolutional layer.

2) **VGG16**: VGG is the elaborated form of the Visual Geometry Group from Oxford. VGG16 has 16 layers-deep CNN. This network receives images of dimensions (224, 224, 3) as input. 64 channels with a 3×3 kernel size having the same padding are in the first and second layers. There are two layers each having 128 convolution layers and a max-pooling layer of stride 2×2 and filter size 3×3 . A max-pooling layer of stride 2×2 that is similar to the preceding one makes up the following layer. After that, 256 filters of size 3×3 are scattered over two convolution layers. Then there are two sets of three convolution layers. A max-pooling layer comes after that. Each layer is padded in the same way having 512 filters of size 3×3 . Then the block of two convolution layers receives the image. In order to avoid the image's spatial characteristics, 1-pixel padding is applied after every layer of CNN.

3) **Model Architecture**: We passed the chest X-ray image dataset to two individual pre-trained models where a global max-pooling of stride 2×2 was applied to the extracted features of each pre-trained model separately. We extracted the features of the Xception model upto a depth-wise separable convolution layer namely 'block14_sepconv2' with a size of $3 \times 3 \times 2048$. We also extracted the features of the VGG16

model upto the 'block5_conv3' layer which has a size of $6 \times 6 \times 512$. Then we flattened the inputs and added three fully connected dense layers consisting of 256 neurons having 0.2 dropouts, 128 neurons having 0.3 dropouts, and 64 neurons respectively. The layers have a 'ReLU' activation function. Then, we fused the outputs of two different models. After that, we also added two dense layers having 128 and 64 neurons. Finally, we added a dense layer for the classification by applying a 'sigmoid' activation function to get the predicted output. Fig. 4 shows the proposed model which explains how our model detects pneumonia from chest X-ray image datasets.

D. Performance Metrics

The performance of our proposed model was scrutinized by the following performance metrics: firstly precision rate, then recall, f1-score, and at last, accuracy. Then the Receiver Operating Characteristic (ROC) and precision-recall curves are calculated from these results. We used a confusion matrix to calculate each of these parameters. The estimated values are used to determine all evaluation metrics, where TP, TN, FP, and FN stand for corresponding true positive, true negative, false positive, and false negative consequences.

1) **Precision**: It stands for the ratio of observations predicted accurately as positive to the total number of observations predicted as positive. Intuitively, it refers to the capacity to avoid classifying a non-positive sample as positive. Its optimal value is 1 and the poorest value is 0. The formula for precision is:

$$P = \frac{TP}{TP + FP} \quad (1)$$

2) **Recall**: Recall measures about how many positive class predictions were produced using all of the positive results in the dataset. It has the formula given below:

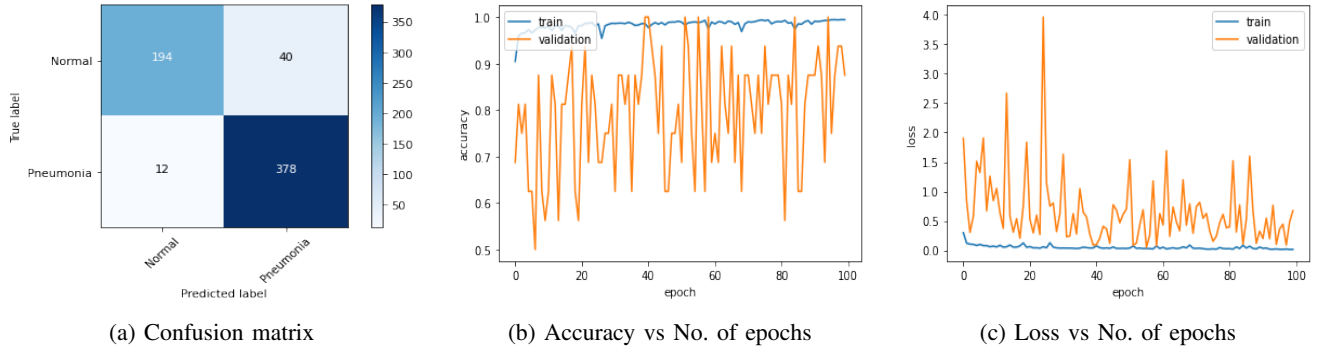


Fig. 5. Different evaluation curves of the proposed model where (a) represents Confusion matrix, (b) represents Accuracy versus epoch graph, and (c) denotes the Loss versus epoch graph.

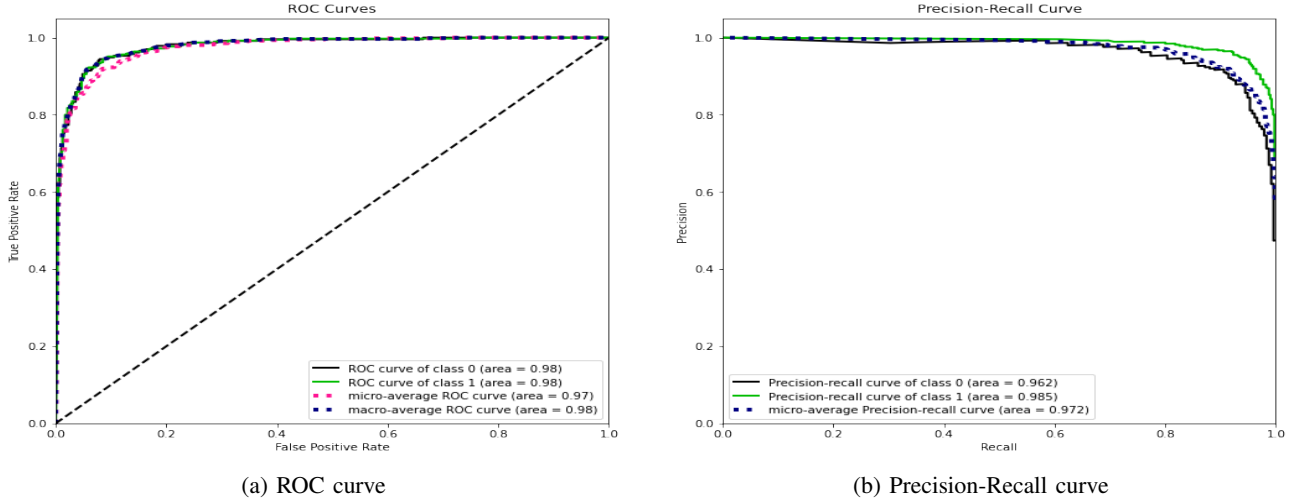


Fig. 6. Receiver operating Characteristics curve and Precision-Recall curve of the proposed model.

$$R = \frac{TP}{TP + FN} \quad (2)$$

3) **F1-score**: F1 score describes a harmonic mean of accuracy and recall. It can be 1 at best and 0 at least. F1 score gets an equal percentage of contribution from precision and recall. Let the precision be P and the recall be R, then its formula is:

$$F1 - score = \frac{2 * P * R}{P + R} \quad (3)$$

4) **Accuracy**: It refers to the proportion of accurate predictions over all of the predictions. Its formula is:

$$A = \frac{TP + TN}{TP + FN + FP + TN} \quad (4)$$

IV. RESULT ANALYSIS

Initially, there was no flattening or extra added dense layer before the fusion of two individual models. But the model achieved the accuracy of 87% for 65 epochs and 89% for 100 epochs. We did not expect that comparatively lower accuracy. So, we applied max-pooling to both models individually and then flattened them. After the flattening, we added some extra dense layers and fusioned them. After that, we trained the

upgraded version of our model with 100 epochs and gained an accuracy of 91.67%. From Table II it is seen that our proposed model has achieved a higher rate of accuracy and a higher rate of precision in comparison to the other research works upon the same datasets. It took an average of 77 seconds to train each epoch of the model. In our testing dataset, there are 234 images that are normal and 390 images that are pneumonia-affected. Fig. 5a represents that the confusion matrix had 194 true positive samples and 40 false positive samples. It means the model can accurately predict 194 images of NORMAL class from 234 images. On the other hand, it showed 12 false negative samples and 378 true negative samples. It means the model can accurately predict 378 images of PNEUMONIA class from 390 chest x-ray images. From that result, it was calculated that our proposed model has produced a 92.30% precision rate with 89.92% recall and 90.87% f1-score as well. Table I shows the overall results of our proposed model. The training accuracy of the model were increasing and training loss were decreasing with respect to the number of epochs as shown in Fig. 5b and Fig. 5c. In the receiver operating characteristics (ROC) curve, as shown in Fig. 6a, the false positive rate was plotted on the x-axis and the true positive

TABLE I. RESULT ANALYSIS OF OUR PROPOSED MODEL

Class	Precision (%)	Recall (%)	F1-Score (%)	Support
Normal	94.17	82.91	88.18	234
Pneumonia	90.43	96.92	93.56	390
Average	92.30	89.92	90.87	—

TABLE II. COMPARISON AMONG OUR MODEL AND OTHER EXISTING MODELS

Research work	Precision (%)	Recall (%)	F1-Score (%)	Accuracy (%)
Ayan et al. [15]	91.00	89.00	90.00	87.00
Liang et al. [16]	—	96.7	92.00	90.5
Han et al. [20]	—	—	92.7	88.6
Chhikara et al. [21]	90.7	95.70	93.10	90.01
Proposed	92.30	89.92	90.87	91.67

rate was plotted on the y-axis. Here, both classes have achieved 98% area in the ROC curve and we got the macro average of 98%. Also, in the precision-recall curve, as shown in Fig. 6b, the NORMAL class has achieved 96.2% and the PNEUMONIA class has gained 98.5% area. Our model has produced a better performance as the values for both classes are higher on the ROC curve and Precision-Recall curve.

V. CONCLUSION

CNNs have made remarkable development in a wide range of fields including medical research, and chest radiography has become one of their most popular applications. In our research, we suggested a new model for detecting chest diseases like pneumonia. We fused two different transfer learning models, the Xception model, and the VGG16 model. The computational time of training our proposed model is comparatively faster and the accuracy is better as well. Because, if we trained the dataset with the Xception model only, it would get an accuracy of 82% and if we also train the dataset with the VGG16 model only, it would get an accuracy of 87%. But, after the fusion of these two models and fine-tuning some parameters with some extra layers, we got an accuracy of 91.67%. Though this model seems to have a complex structure, it consumes less computational time and produces higher accuracy which are the important factors for a model's performance.

REFERENCES

- [1] M. F. Hashmi, S. Katiyar, A. G. Keskar, N. D. Bokde, and Z. W. Geem, "Efficient pneumonia detection in chest xray images using deep transfer learning," *Diagnostics*, vol. 10, no. 6, p. 417, 2020.
- [2] D. H. Johnson, S.; Wells. Viral pneumonia: Symptoms, risk factors, and more. Available: <https://www.healthline.com/health/viral-pneumonia>. [Accessed on December 31, 2019].
- [3] J. E. Luján-García, C. Yáñez-Márquez, Y. Villuendas-Rey, and O. Camacho-Nieto, "A transfer learning method for pneumonia classification and visualization," *Applied Sciences*, vol. 10, no. 8, p. 2908, 2020.
- [4] M. Khan, M. Sharif, T. Akram, and R. Damaševicius, "Maskeli unas, r," *Skin lesion segmentation and multiclass classification using deep learning features and improved moth flame optimization. Diagnostics*, vol. 11, p. 811, 2021.
- [5] M. I. Sharif, M. A. Khan, M. Alhussein, K. Aurangzeb, and M. Raza, "A decision support system for multimodal brain tumor classification using deep learning," *Complex & Intelligent Systems*, pp. 1–14, 2021.
- [6] Pneumonia. Available:<https://www.radiologyinfo.org/en/info.cfm?pg=pneumonia>. [Accessed on December 31, 2019].
- [7] T. Cherian, E. K. Mulholland, J. B. Carlin, H. Ostensen, R. Amin, M. d. Campo, D. Greenberg, R. Lagos, M. Lucero, S. A. Madhi et al., "Standardized interpretation of paediatric chest radiographs for the diagnosis of pneumonia in epidemiological studies," *Bulletin of the World Health Organization*, vol. 83, pp. 353–359, 2005.
- [8] T. Franquet, "Imaging of pneumonia: trends and algorithms," *European Respiratory Journal*, vol. 18, no. 1, pp. 196–208, 2001.
- [9] A. M. Tahir, M. E. Chowdhury, A. Khandakar, S. Al-Hamouz, M. Abdalla, S. Awadallah, M. B. I. Reaz, and N. Al-Emadi, "A systematic approach to the design and characterization of a smart insole for detecting vertical ground reaction force (vgrf) in gait analysis," *Sensors*, vol. 20, no. 4, p. 957, 2020.
- [10] M. E. Chowdhury, K. Alzoubi, A. Khandakar, R. Khallifa, R. Abouhasera, S. Koubaa, R. Ahmed, and A. Hasan, "Wearable real-time heart attack detection and warning system to reduce road accidents," *Sensors*, vol. 19, no. 12, p. 2780, 2019.
- [11] M. N. T. Akhand, S. Das, and M. Hasan, "Traffic density estimation using transfer learning with pre-trained inceptionresnetv2 network," in *Machine Intelligence and Data Science Applications*. Singapore: Springer Nature Singapore, 2022, pp. 363–375.
- [12] S. Das, A. A. Fime, N. Siddique, and M. Hashem, "Estimation of road boundary for intelligent vehicles based on deeplabv3+ architecture," *IEEE Access*, vol. 9, pp. 121 060–121 075, 2021.
- [13] S. Das, M. S. Imtiaz, N. H. Neom, N. Siddique, and H. Wang, "A hybrid approach for bangla sign language recognition using deep transfer learning model with random forest classifier," *Expert Systems with Applications*, p. 118914, 2022.
- [14] Z.-P. Jiang, Y.-Y. Liu, Z.-E. Shao, and K.-W. Huang, "An improved vgg16 model for pneumonia image classification," *Applied Sciences*, vol. 11, no. 23, p. 11185, 2021.
- [15] E. Ayan and H. M. Ünver, "Diagnosis of pneumonia from chest x-ray images using deep learning," in *2019 Scientific Meeting on Electrical-Electronics & Biomedical Engineering and Computer Science (EBBT)*. Ieee, 2019, pp. 1–5.
- [16] G. Liang and L. Zheng, "A transfer learning method with deep residual network for pediatric pneumonia diagnosis," *Computer methods and programs in biomedicine*, vol. 187, p. 104964, 2020.
- [17] S. S. Yadav and S. M. Jadhav, "Deep convolutional neural network based medical image classification for disease diagnosis," *Journal of Big Data*, vol. 6, no. 1, pp. 1–18, 2019.
- [18] R. Jain, P. Nagrath, G. Kataria, V. S. Kaushik, and D. J. Hemanth, "Pneumonia detection in chest x-ray images using convolutional neural networks and transfer learning," *Measurement*, vol. 165, p. 108046, 2020.
- [19] D. Varshni, K. Thakral, L. Agarwal, R. Nijhawan, and A. Mittal, "Pneumonia detection using cnn based feature extraction," in *2019 IEEE international conference on electrical, computer and communication technologies (ICECCT)*. IEEE, 2019, pp. 1–7.
- [20] Y. Han, C. Chen, A. Tewfik, Y. Ding, and Y. Peng, "Pneumonia detection on chest x-ray using radiomic features and contrastive learning," in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2021, pp. 247–251.
- [21] P. Chhikara, P. Singh, P. Gupta, and T. Bhatia, "Deep convolutional neural network with transfer learning for detecting pneumonia on chest x-rays," in *Advances in bioinformatics, multimedia, and electronics circuits and signals*. Springer, 2020, pp. 155–168.
- [22] D. Kermany, K. Zhang, M. Goldbaum et al., "Labeled optical coherence tomography (oct) and chest x-ray images for classification," *Mendeley data*, vol. 2, no. 2, 2018.