

Deep Attention Network for Pneumonia Detection Using Chest X-Ray Images

Sukhendra Singh¹, Sur Singh Rawat², Manoj Gupta³, B. K. Tripathi⁴, Faisal Alanzi⁵,
Arnab Majumdar⁶, Pattaraporn Khuwuthyakorn⁷ and Orawit Thinnukool^{7,*}

¹Department of Information Technology, JSS Academy of Technical Education, Noida, India

²JSS Academy of Technical Education, Noida, India

³Department of Electronics and Communication Engineering, JECRC University Jaipur, Rajasthan, India

⁴Harcourt Butler Technological University Kanpur, India

⁵Department of Electrical Engineering, Prince Sattam Bin Abdulaziz University, College of Engineering,
Al Kharj, 16278, Saudi Arabia

⁶Faculty of Engineering, Imperial College London, London, SW7 2AZ, UK

⁷College of Arts, Media and Technology, Chiang Mai University, Chiang Mai, 50200, Thailand

*Corresponding Author: Orawit Thinnukool. Email: orawit.t@cmu.ac.th

Received: 15 May 2022; Accepted: 22 June 2022

Abstract: In computer vision, object recognition and image categorization have proven to be difficult challenges. They have, nevertheless, generated responses to a wide range of difficult issues from a variety of fields. Convolution Neural Networks (CNNs) have recently been identified as the most widely proposed deep learning (DL) algorithms in the literature. CNNs have unquestionably delivered cutting-edge achievements, particularly in the areas of image classification, speech recognition, and video processing. However, it has been noticed that the CNN-training assignment demands a large amount of data, which is in low supply, especially in the medical industry, and as a result, the training process takes longer. In this paper, we describe an attention-aware CNN architecture for classifying chest X-ray images to diagnose Pneumonia in order to address the aforementioned difficulties. Attention Modules provide attention-aware properties to the Attention Network. The attention-aware features of various modules alter as the layers become deeper. Using a bottom-up top-down feedforward structure, the feedforward and feedback attention processes are integrated into a single feedforward process inside each attention module. In the present work, a deep neural network (DNN) is combined with an attention mechanism to test the prediction of Pneumonia disease using chest X-ray pictures. To produce attention-aware features, the suggested network was built by merging channel and spatial attention modules in DNN architecture. With this network, we worked on a publicly available Kaggle chest X-ray dataset. Extensive testing was carried out to validate the suggested model. In the experimental results, we attained an accuracy of **95.47%** and an F- score of **0.92**, indicating that the suggested model outperformed against the baseline models.



This work is licensed under a Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Keywords: Attention network; image classification; object detection; residual networks; deep neural network

1 Introduction

Recent developments in artificial intelligent (AI) have improved the quality of services and people's experience in many diverse domains. Clinical Practice is one of them. Clinicians take assistance from these deep learning-based prediction tools which save their time, making their prediction much more reliable. Furthermore, these models can provide an accurate prediction of disease and interpret the severity of infection [1] in case of disease. Pneumonia is a disease that is caused by infections in the lungs and it is detected by examining chest X-ray radiographs. Pneumonia disease is one of the leading causes of death worldwide, especially among the younger and old age population. Pneumonia is an infection of the lungs caused by microbes, with inflammation as a result.

The inflammation causes the liquid to enter the lung tissue, making breathing more difficult. If in the acute case, it may cause death in the absence of proper timely medication. Pneumonia is detected by radiologists by examining a chest X-ray radiograph (CXR). If they are inconclusive about infection in CXR, then they recommend computer tomography (CT) scan of the lungs to examine finer details. Apart from examining the CXR, radiologists use computer-aided detection systems [2–4] which run on deep learning (DL) techniques which can recognize patterns in images as shown in Fig. 1. Such a system saves time and increases the efficiency of radiologists and it minimizes disagreement [5] among multiple radiologists in some cases. Fig. 2 shows samples of chest radiographs of a healthy person and the person suffering from Pneumonia.

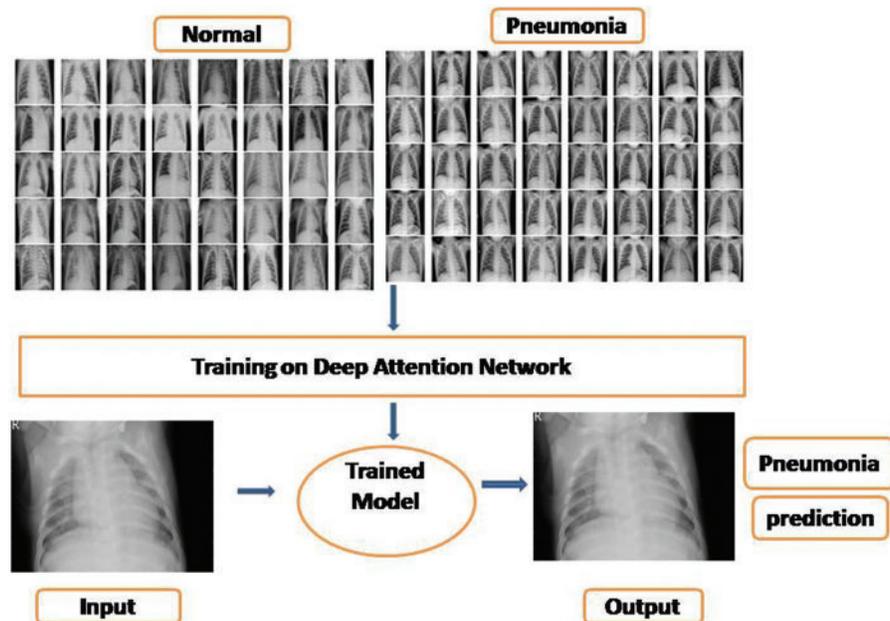


Figure 1: Proposed architecture' high level view

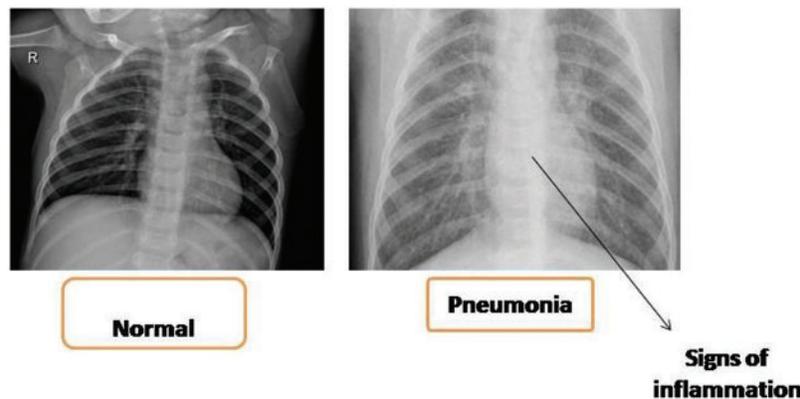


Figure 2: Sample of chest radiograph of normal and pneumonia patient

Previously researchers used hand-crafted feature [6] extraction to convert raw images into relevant features for classifying a medical image for easier classification. This necessitates domain knowledge to do feature engineering, which entails extracting relevant features, changing them, and deleting unnecessary ones. Decision trees, random forest [7], support vector machines [8] have delivered noticeable performance. Deep neural networks [9–11] have substituted feature engineering with feature learning where better features are learned, and it has demonstrated more promising results. Complex features are extracted from input images and further categorized in the architecture's final levels. Transfer learning has been applied with pre-trained classifiers [12,13] with better performance. Although the attention mechanism [14–17] was primarily used for language processing-related tasks like image captioning [18–22], language translation [23].

The motivation behind the proposed work is inspired by the recent advances in deep learning-based segmentation models, a classical machine learning-based end-to-end model which was used to classify chest X-ray images using a hand-crafted set of features defining the local look. Also, the suggested design extracts attention-aware features, hence improving classification performance.

In this paper, we are proposing to apply the attention mechanism in baseline Convolution Neural Network (CNN) and Residual CNN for binary image classification tasks to classify CXR images into normal or Pneumonia. Normally, in CNN, though, each convolution layer processes the entire image regardless of what is foreground and what is background. Attention mechanism enables the model to learn only meaning features by only focusing on the main part of the image and ignoring the rest of the unnecessary part of the image so that it does not learn any feature from the least important portion of the image. The proposed model learns only attention-aware features which makes classification more accurate.

Attention is a mechanism that is inspired by the human brain's complex and powerful cognitive power to detect something. When the human eye observes some image, text, or object, it does not focus uniformly on the entire thing. It focuses mainly on the selected region, and the remaining is ignored while processing. With this same mechanism, only the most relevant features extracted from the selected region of the image have a role to play in classification while ignoring irrelevant parts of the image. In the proposed architecture, many attention modules yield attention-aware features. As the layers become deeper, different modules' attention-aware features modify themselves to provide greater performance.

The novelty of the proposed network is that it captures local, global, and spatial information from chest X-ray images in order to enhance diagnostic performance. By incorporating the channel and spatial attention modules, it has slightly increased the complexity of the architecture but the number of trainable parameters remain the same which does not impose any extra overhead on performance on the contrary it increases the diagnostic accuracy by an appreciable amount.

The following is the summary of contributions in this paper.

- i. The paper presents attention-aware CNNs to classify CXR images for the detection of Pneumonia. This attention-aware CNN architecture leads to higher accuracy in comparison to architecture without an attention mechanism. We have applied attention mechanisms in baseline CNN and Residual CNN architecture and compared the performances.
- ii. We assessed the performance of pre-trained classifiers on the same dataset compared to our findings.

The rest of the paper is structured as follows. Section 2 summarizes the findings of recent research studies in related areas. Section 3 discusses the material and methods in which we have described our chosen dataset along with its characteristics. Section 4 is about working principle and proposed architecture and the Section 5 provides details of experimentation and results from the analysis. Lastly, the Section 6 discusses the conclusion and its future scope.

2 Related Work

The motivation behind mimicking human attention was first seen in the domain of natural language processing, computer vision [24] in a view of reducing the computation complexity while processing an image. Moreover, to improve the performance, a model was introduced that would mainly focus on the specific region of interest in an image instead of the entire scene. The attention mechanism was further refined to machine translation model to address the issue involved. In the recent, it has been employed in a large number of DL models across various domains and tasks like text classification, image captioning, sentiment analysis and speech recognition [25,26].

Two types of attention mechanisms have arisen in the literature, both inspired by the human visual system. The first is a top-down technique, in which an iterative process selects the appropriate region from a pool of records about the scene. The bottom-up technique, on the other hand, identifies the most relevant and conspicuous locations along the visual path. The top-down strategy is iterative, and it is slower than the single-pass bottom-up approach. Furthermore, the bottom-up strategy selects the most relevant regions from the input data progressively, but the errors produced by these sequential processes increase with the depth of the process.

The attention mechanism has become a hot research area due to several reasons. Firstly, the attention mechanism involved in any model is impressive performance against state-of-the-art approaches. Secondly, the attention model can be jointly trained with a base recurrent neural network [27] or the CNNs by making use of the backpropagation approach. In addition, the induction of the transformer model was very much adopted in the tasks like image processing, video processing and recommendation system which has increased the performance of the attention model and the parallelized issue involved in the recurrent neural networks can be circumvented.

It is well known that the neural network involved in classification models the data as the numeric vector which comprises the low-level features, in which all the features are assigned the same weights irrespective of their capabilities. This issue has been addressed in the attention model [28] wherein the variable was assigned to different features according to their importance or in other words. The

attention model computes the weight distribution based on the input features and assigns higher values to the features with higher rank, which means that the attention model computes the weight distribution based on the input features and assigns higher values to the features with a higher rank.

To elaborate further, the attention layers in the attention mechanism are the alignment layer, attention weight, and the context vector. The working of the attention layer is to compute the alignment score between the encoded vector $h = \{h_1, h_2, \dots, h_n\}$ and a vector v . The softmax is applied to calculate the probability distribution α_i by normalizing over all the n elements of h where $i = 1, 2, \dots, n$ as given shown in the Eqs. (1) and (2).

$$\alpha_i = \frac{\exp(h_i v)}{\sum_{j=1}^n \exp(h_j v)} \quad (1)$$

$$O = \sum_{i=1}^n \alpha_i h_i \quad (2)$$

From the above equations, the larger value of α_i means that h_i contributes important information to vector v . Also the output O of the attention mechanism is a weighted sum of all elements in the encoded vector h_i .

In machine vision challenges, we now have a plethora of attention and DL algorithms, such as (1) attention-based CNN, (2) CNN transformer pipelines, and (3) Hybrid transformer [29,30]. The primary idea of the attention-based CNN was to discover the most important components of the feature maps in CNN so that redundancy could be reduced for machine vision applications. DL [31–33] has used the attention mechanism in a variety of machine vision applications, including object detection, picture captioning, and action recognition. The CNN is utilized to provide the features map to a transformer and act as a teacher to the transformer in the CNN transformer pipelines.

The attention mechanism has evolved as a key mechanism for improving the precision and efficiency of the system by focusing on the most informative input while ignoring clutter and noise.

2.1 Attention Mechanism

Attention is a useful technique for improving the performance of the Encoder-Decoder architecture on neural network-based machine translation applications. It's a connection between the encoder and the decoder that allows the decoder to read data from every encoder's secret state. The model can selectively focus on valuable bits of the input sequence using this framework, and thus learn the relationship between them [34]. This makes it easier for the model to deal with extended sentences. A neural network is a simplified model of the brain. In deep neural networks (DNNs), the attention mechanism is an attempt to emulate the similar behavior of selectively concentrating on a few relevant things while ignoring others. Attention is used by neural networks to improve input data [35]. Other sections of the data are diminished as a result of the effect, with the idea being that the network should focus more on that data. Gradient descent [36] is used to learn which parts of the data are more significant than others, which is dependent on the context. In the 1990 s, attention-like processes such as multiplicative modules, sigma pi units, and hyper networks were introduced.

An attention function is a mapping of a query and a set of key-value pairs to output, where the query, keys, values, and output are all vectors. The output is a weighted sum of the values, with the weights assigned to each value set by the query's key compatibility function.

2.2 Attention Mechanism for Image Classification

As the human eye rapidly scans the entire image, image attention is the act of locating a target region that requires attention. This target region is given additional weight (distribution) in order to acquire the most exact data about the target while suppressing other irrelevant data. Soft attention is the most popular method since it is a fully differentiable procedure that allows CNN models to train from start to finish. Most soft attention models build an attention template, which they then utilize to locate unique components, in order to align the weights of discrete segments in a sequence or an image.

The hard attention mechanism, unlike soft attention, is a random, non-differentiable process that determines the significance of particular regions one at a time rather than identifying the image's key features as a whole. The weight of the arithmetic mean of attention can be calculated from an image's attention spectrum when using attention learning for picture categorization. In the same way that standard natural language processing may be used to gather image-based attention, the technique can be used to gather image-based attention.

Because it effectively pulls characteristics from the data, the DNN is suitable for pixel-wise categorization of images. The attention mechanism, which mimics how people perceive images, is important for quickly and precisely acquiring crucial information. With CNNs [37], all convolution layers process all of the image's features and details. To average the features and details of the entire image, use numerous convolution layers and a global average pooling [38] in the final layer. The last affine fully connected layer in this network estimates the image category. The more the background and other non-required information affect the classification results, the smaller the area of the full image taken. A huge amount of data is used to eliminate errors in findings, and the neural network learns not to output background features, thus global average pooling is unaffected.

One method is to branch the output of one layer in the CNNs, resulting in a single image produced by two or more convolutional layers. Here, we set sigmoid, which serves as the convolution output's activation function, to produce a value between 0.0 and 1.0 for each pixel. Sigmoid is a function that saves input values between 0 and 1. The output of the convolution function is multiplied by the original output. Two additional layers estimate the original output's focal areas. The values approaching 0 aren't worth focusing on. We're interested in areas where the values are multiplied by values close to one because the output is close to the original number.

Configuring processing in this way results in most of the sigmoid values that are near 0 not being used at all in the downstream recognition process, as those values would be completely disregarded. Configuring a neural network to estimate the area of focus using the result of which is the most used approach of using attention for image classification [39].

3 Materials and Methods

3.1 Dataset

The dataset is divided into three folders (train, test, and val), each of which contains a subdirectory for each image type (Pneumonia/Normal). There are 5,856 JPEG X-Ray images in total, divided into two groups (Pneumonia/Normal). From retrospective cohorts, anterior-posterior chest X-ray images from children aged one to five years old at Guangzhou Women and Children's Medical Center in Guangzhou were chosen. Chest X-ray imaging was performed as part of the patients' routine medical care. Quality control was performed on all chest radiographs, with any scans that were poor quality or unreadable being discarded. Two professionals then evaluated the diagnostic for the photographs before they were accepted for use in the AI system. The dataset for each class is shown in [Tab. 1](#).

Table 1: Class wise distribution of dataset

| Class | No. of images |
|---------------|---------------|
| Pneumonia (P) | 4273 |
| Normal (N) | 1583 |

We have divided the dataset 75% into training set and 80% of the remaining for test set and 20% for validation set as shown in the [Tab. 2](#).

Table 2: Splitting of dataset into train, test and validation

| | No. of images | No. of images from P class | No. of images from N class |
|-----------------|---------------|----------------------------|----------------------------|
| Training data | 4392 | 3205 | 1187 |
| Validation data | 292 | 226 | 66 |
| Test data | 1172 | 842 | 330 |

3.2 Data Augmentation

We must purposefully extend our dataset to avoid the risk of overfitting. We have the ability to expand the amount of your current dataset. Small adjustments to the training data are made to reproduce the variations. Methods for changing the array representation while keeping the label the same while changing the training data are known as data augmentation strategies. Grayscale, horizontal and vertical flips, random crops, colour hiccups, translations, rotations, and a variety of additional enhancements are common. By making just a few of these changes to our training data, we can quickly double or quadruple the amount of training examples and construct a very robust model.

We used data augmentation transformation techniques such as 30 degree sequence rotation, zooming within a range of 0.2, width shift by 0.1, height shift by 0.1, and horizontal flipping.

3.3 Implementation Environment and Tools Used

Python 3.6 was used here in this work for the implementation. We have utilized Google Colaboratory (<https://colab.research.google.com/>), which uses the Jupyter notebook environment, because we needed faster GPUs for our tests. Colaboratory provided a Tesla K80 GPU to speed up the processing. In addition, the Google Colaboratory environment provides 12 GB of GDDR5 VRAM and 13 GB of RAM for free. We utilised the Keras and Tensorflow libraries in versions 2.2.2 and 1.10.0, respectively, to create DL models. The parameter settings utilized during the experiment are listed in [Tab. 3](#).

Table 3: Parameter settings used in entire experimentation

| Parameter | Value |
|---------------|----------------------|
| Image size | 150 × 150 |
| Loss function | Binary cross entropy |

(Continued)

Table 3: Continued

| Parameter | Value |
|-------------------------|--|
| Dropout | 0.2 |
| Stride | 1 |
| Batch size | 16 |
| Optimizer | Adam |
| Epochs | 50 |
| Learning rate reduction | ReduceLROnPlateau(monitor = 'val_accuracy', patience = 2, verbose = 1, factor = 0.8) |

4 Proposed Architecture

4.1 Working Principle

Image attention is the act of locating a target place that requires attention when the human eye rapidly sweeps the global image. This target region is given additional weight (distribution) in order to acquire the most exact data about the target while suppressing other irrelevant data. The most prevalent strategy is soft attention, which is a fully differentiable procedure that allows CNN models to train from beginning to end. In order to align the weights of discrete segments in a sequence or an image, the most soft attention models create an attention template, which they subsequently use to discover unique components. Unlike soft attention [36–39], the hard attention mechanism is a non-differentiable, random process that decides the significance of individual regions one at a time rather than identifying the image's main elements as a whole.

The weight of the arithmetic mean of attention can be determined from an image's attention spectrum when using attention learning for image categorization. The attention model calculates and selects the most relevant region of the feature spectrum for the final classification job by learning the attention of the CNN feature spectrum, and supplies the maximum attention input (weight distribution). When the attention weight is added to the last layer of CNN features, the original attributes are muted to varying degrees. The original feature spectrum is blended with the weighted feature spectrum to solve this suppression. The completely connected layer receives the fusion spectrum after that.

In the second fully connected layer, the attention feature spectrum is coupled with the fully connected feature spectrum in the channel direction before being delivered to the classification layer for classification, which has been adjusted by the global average pooling dimension. The Softmax layer's output error is back propagated during the back propagation process, and the parameters are modified using a random gradient descent approach to lower the final loss function value and bring the network closer to convergence.

The output of the CNN's final convolutional layer is acquired using the soft attention technique. This data is entered into the attention model, which then generates the proper attention spectrum. The following network uses the output attention feature spectrum as an input, and the attention spectrum is utilized to weight the original feature spectrum.

Here I is the input image in Fig. 3. Each of the n parameters in the attention model ($a_1, a_2, \dots, a_i, \dots, a_n$) describes a different element of the image. The importance of each a_i in relation to the input I determines the return value of the model's attention spectrum (more precisely, the weight

values of the n parameters). You can locate the area that requires the most attention by filtering the input image through this output. Our proposed model's layer-by-layer structure is shown in Fig. 4.

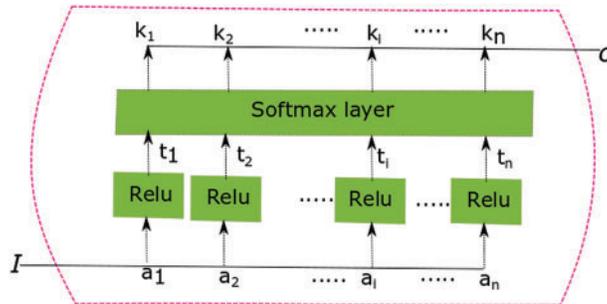


Figure 3: Block diagram of attention mechanism

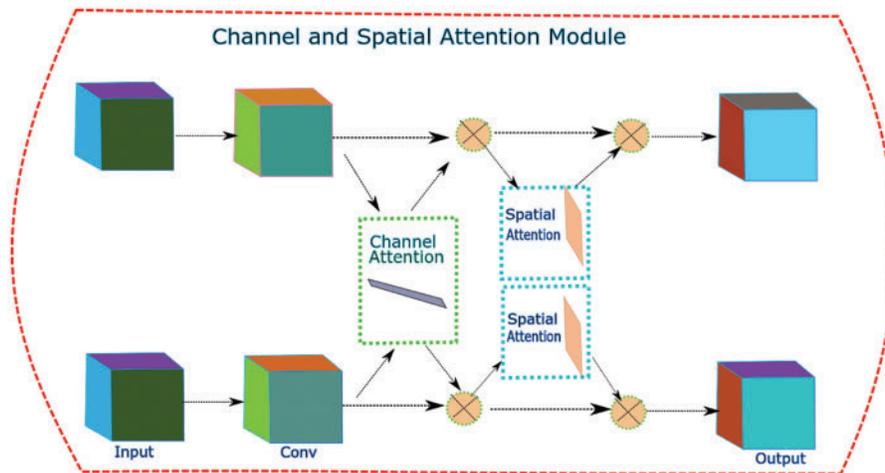


Figure 4: System design of proposed architecture

4.2 Layer Wise Organization of the Proposed Model

Convolutional Block: It extracts feature map in the form of high-level features and then these are passed to channel attention and spatial attention blocks to enhance the feature maps.

Batch Normalization: Batch normalization [40] is a network layer that allows each layer to learn in a more autonomous manner. It's used to normalize the output of the previous layers. When batch normalization is employed, learning becomes more efficient, and it can also be used to prevent model overfitting.

Channel attention: A Channel of Attention A CNN module that focuses on channel-based attention is Module. Using the inter-channel relationship of features, we create a channel attention map. Channel attention focuses on 'what' is significant given an input image since each channel of a feature map is viewed as a feature detector. It compresses the spatial dimension of the input feature map to compute channel attention effectively. We initially perform average-pooling and max-pooling processes to aggregate spatial information from a feature map, resulting in two separate spatial context

descriptors, F_{Avg}^c , F_{max}^c , which stand for average-pooled features and max-pooled features, respectively [41–43].

Spatial Attention: Based on the inter-spatial relationship of features, it generates a spatial attention map. In contrast to channel attention, which is concerned with the location of a channel, spatial attention, which is complementary to channel attention, is concerned with the location of an informational component. To compute spatial attention, we first use average pooling and max pooling processes along the channel axis and concatenate them to build an efficient feature descriptor. To produce a spatial attention map that encodes where to highlight or suppress, we apply a convolution layer to the concatenated feature descriptor [44–48].

4.3 Resnet-50 with Attention

The 50-layer ResNet architecture was merged with attention layers in between. This proposed model was previously compared to the standard ResNet-50 model. We started with a convolutional layer and added channel and spatial attention. After that, we used a convolutional block, channel attention, and spatial attention, followed by two identity blocks. A convolutional block, channel attention, spatial attention, and five identity blocks follow this model, which is followed by a convolutional block, channel attention, spatial attention, and five identity blocks. Finally, there are two identification blocks, a convolutional block, a channel attention block, and a spatial attention block. Fig. 5 shows the architecture of channel attention and spatial attention, which are situated between ResNet layers [48–51].

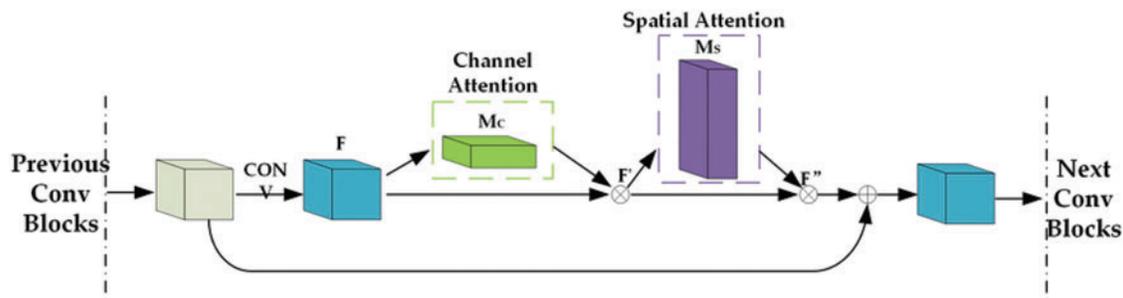


Figure 5: Integrating attention with a CNN network

5 Experimental Results and Its Analysis

We have experimented with an attention mechanism with a baseline, CNN, and Resnet50 on the chosen dataset and observed the improvement in the performance parameters and we have compared the same with a pre-trained model on the same dataset. To compare performance, we worked on the same dataset, keeping all hyper-parameter values the same on the same number of epochs. Tab. 4 shows the confusion matrix, Tab. 5 shows the performance of a pre-trained model using transfer learning on the same dataset. From Tab. 6 below, we observed improvement in all performance parameters, when we compared the performance of baseline CNN with and without the attention mechanism and there we got a significant rise of 12% in accuracy. Attention mechanism did not increase the learnable parameters. Figs. 6 and 7 show confusion matrix of baseline CNN with attention network and the receiver-operating curve (ROC) curve. Fig. 8 shows plot for accuracy, loss, ROC, precision and recall (PRC) for CNN with attention network.

Table 4: Confusion matrix

| Prediction/Actual | Positive | Negative |
|-------------------|----------|----------|
| Positive | TP | FP |
| Negative | FN | TN |

Table 5: Results of pre-trained model on the dataset

| Pre-trained model | Accuracy | Precision | Recall | F1-measure | Parameters |
|-------------------|----------|-----------|--------|------------|------------|
| VGG16 | 92.14% | 0.89 | 0.91 | 0.93 | 14,714,688 |
| VGG19 | 89.90% | 0.85 | 0.88 | 0.87 | 20,024,384 |
| ResNet50 | 84.24% | 0.84 | 0.82 | 0.83 | 23,587,712 |
| InceptionV3 | 89.42% | 0.86 | 0.90 | 0.88 | 21,802,784 |
| Xception | 86.64% | 0.83 | 0.88 | 0.85 | 22,910,480 |
| InceptionResNetV2 | 86.17% | 0.89 | 0.89 | 0.89 | 58,331,648 |
| NasNetLarge | 88.14% | 0.86 | 0.84 | 0.85 | 84,916,818 |

Table 6: Comparison of results of attention mechanism on a baseline CNN with and without attention

| | Accuracy | Precision | Recall | F1-Score | AUC | Parameters |
|-------------------------------|----------|-----------|--------|----------|------|------------|
| Baseline CNN | 82.50% | 0.8156 | 0.8562 | 0.81 | 0.83 | 459,969 |
| Baseline CNN (with attention) | 94.28% | 0.8488 | 0.9696 | 0.9052 | 0.95 | 459,969 |

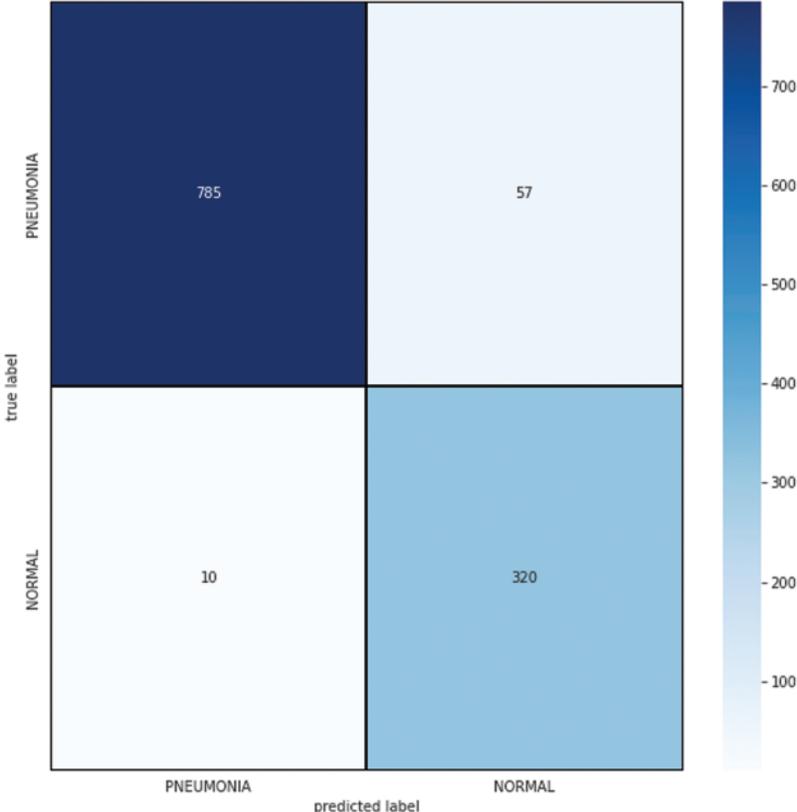


Figure 6: Confusion matrix

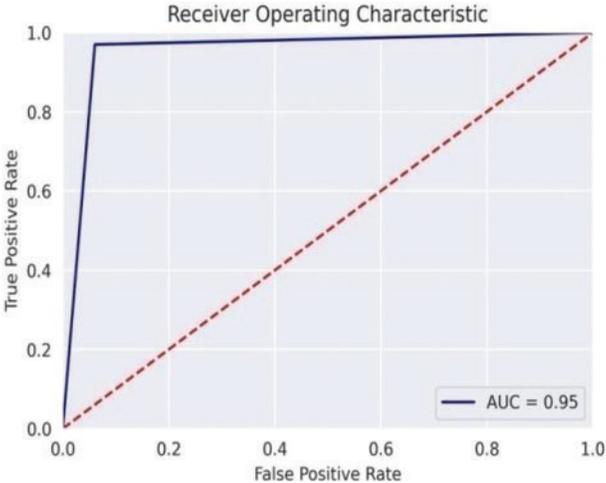


Figure 7: ROC curve

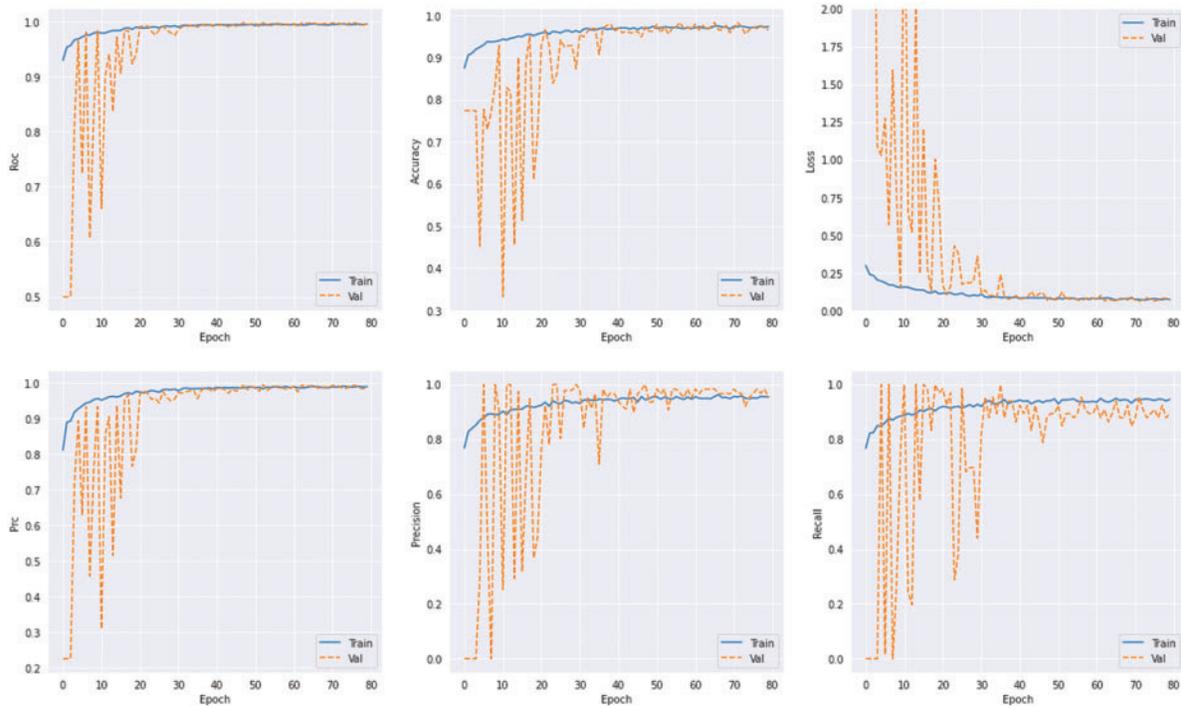


Figure 8: Plot for accuracy, ROC, Loss, Prc, precision and recall

Evaluation Indicators

To compare the performance of various DL architectures, most common evaluation indicators are classification accuracy, sensitivity, specificity, Area under Curve (AUC) and confusion matrix.

Classification Accuracy: This indicator tells us about correct predictions made by our trained model on unseen new images from the test set. The mathematical formula of classification accuracy is as given in Eq. (3)

$$\begin{aligned} \text{Classification Accuracy} &= \frac{\text{No. of correct prediction on test set}}{\text{size of test set}} \\ &= \frac{(\text{TP} + \text{TN})}{(\text{TP} + \text{TN}) + \text{FP} + \text{FN}} \end{aligned} \tag{3}$$

True Positive (TP): Sample that is correctly classified as positive same as ground truth.

False Positive (FP): Samples that is wrongly classified as positive as opposed to ground truth.

True negative (TN): Sample that is correctly classified as negative same as ground truth.

False negative (FN): Samples that is wrongly classified as negative as opposed to ground truth.

Sensitivity: It is the proportion of positive samples that were classified as positive. Higher sensitivity indicates that the system can accurately predict the presence of disease and false negative cases will be very less. It is mathematically defined as in Eq. (4)

$$\text{sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \tag{4}$$

Specificity: It shows the proportion of negative samples that were classified as negative. Higher specificity means that the system is making correct predictions about healthy people. It is defined as in Eq. (5).

$$\text{specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (5)$$

Precision and Recall are defined as in Eq. (6).

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (6)$$

F-Score is harmonic mean Precision (P) and Recall (R). It is defined as in Eq. (7)

$$\text{F - score} = \frac{2 * P * R}{(P + R)} \quad (7)$$

The Receiver Operator Characteristic (ROC) curve is graph between sensitivity and (1-specificity). It is used to analyses the trade-off between sensitivity and specificity.

Area under Curve (AUC): Higher AUC means that the model can better distinguish between positive and negative classes.

Confusion Matrix: Gives the predictions made by the model in terms of TP, TN, FP and FN as shown in the Tab. 5.

Performance Evaluation: Finally the trained model from the proposed architecture is tested against test-set and various performance metrics like accuracy, sensitivity, specificity, F-score, AUC are computed. The confusion matrix was also evaluated for the model.

We experimented with our proposed model on Kaggle dataset for three split ratios, 20%, 30% and 40%. Although this dataset was imbalanced, we overcame this problem by under sampling [41–43].

Tab. 7 shows the performance improvement of attention mechanism on Resnet50 and here an improvement of 10% accuracy was obtained if attention mechanism was applied with Resnet 50. To find out appropriate initial value of learning rate, we have performed this experiment both in baseline and in Resnet50 by initializing learning rate with 0.01, 0.001 and 0.0001 and keeping batch size fixed as 16. Tab. 8 shows test accuracy and test loss for learning rates 0.01, 0.001 and 0.0001.

Table 7: Comparison of results of attention mechanism on a residual Network with and without Attention

| | Accuracy | Precision | Recall | F1-Score | AUC | Parameters |
|-------------------------|----------|-----------|--------|----------|------|------------|
| Resnet50 | 84.53% | 0.8556 | 0.9190 | 0.9090 | 0.91 | 28,333,313 |
| Resnet50 with attention | 95.73% | 0.8825 | 0.9787 | 0.9281 | 0.96 | 28,333,313 |

Fig. 9 is showing training accuracy and loss variation in 50 epochs wrt learning rates as 0.01, 0.001, and 0.0001. From here, we observe that with an initial value of learning rate as 0.001, the model is offering better accuracy. Fig. 10 presents the Validation Accuracy, Loss variation in Learning rate in Resnet 50 with Attention. Fig. 11 presents the training accuracy and test loss respectively for batch sizes 8,16,32 and 64 in Resnet50 with Attention. Fig. 12 present the validation accuracy and test loss respectively for batch sizes 8,16,32 and 64 in Resnet50 with Attention.

Table 8: Impact on accuracy, loss in baseline CNN and resnet50 with attention mechanism

| Learning rate | Batch size | Epochs | Attention with baseline | | Attention with Resnet | |
|---------------|------------|-----------|-------------------------|---------------|-----------------------|--------------|
| | | | Test accuracy | Test loss | Test Accuracy | Test loss |
| 0.01 | 16 | 50 | 93.94% | 0.1946 | 94.710% | 0.1218 |
| 0.001 | 16 | 50 | 95.22% | 0.1206 | 95.307% | .1557 |
| 0.0001 | 16 | 50 | 94.11% | 0.1451 | 94.80% | .1550 |

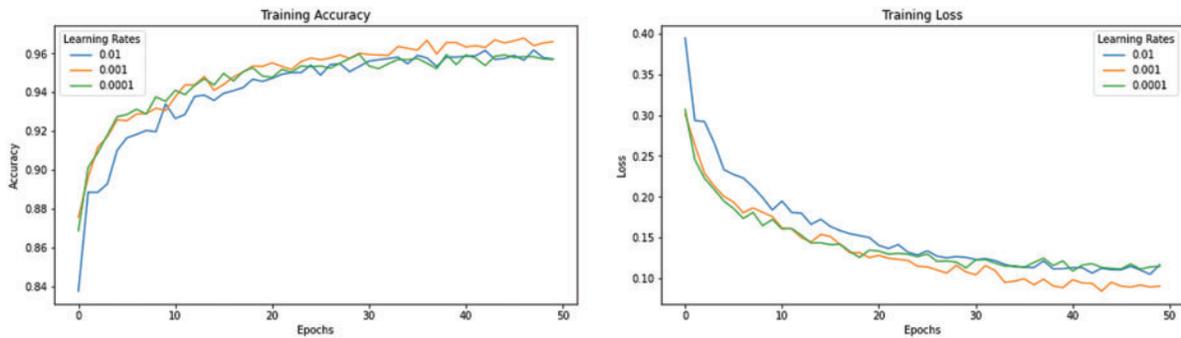


Figure 9: Training accuracy, loss variation in learning rate in Resnet50 with attention

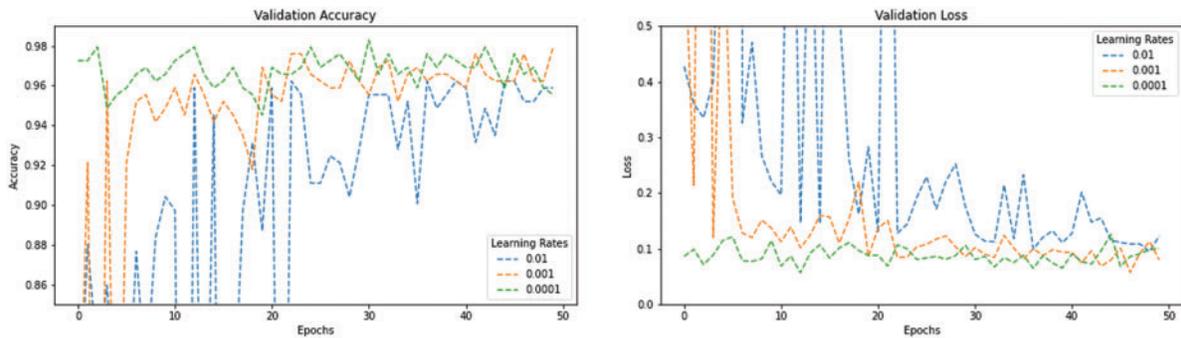


Figure 10: Validation accuracy, loss variation in learning rate in Resnet50 with attention

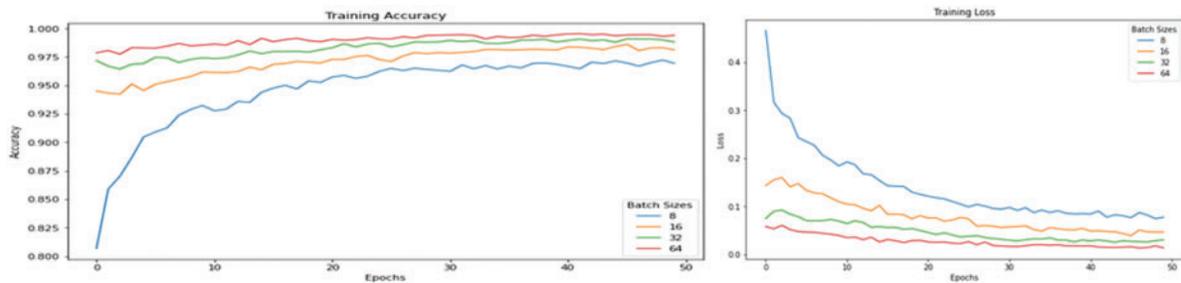


Figure 11: The training accuracy and test loss respectively for batch sizes 8,16,32 and 64 in resnet50 with attention

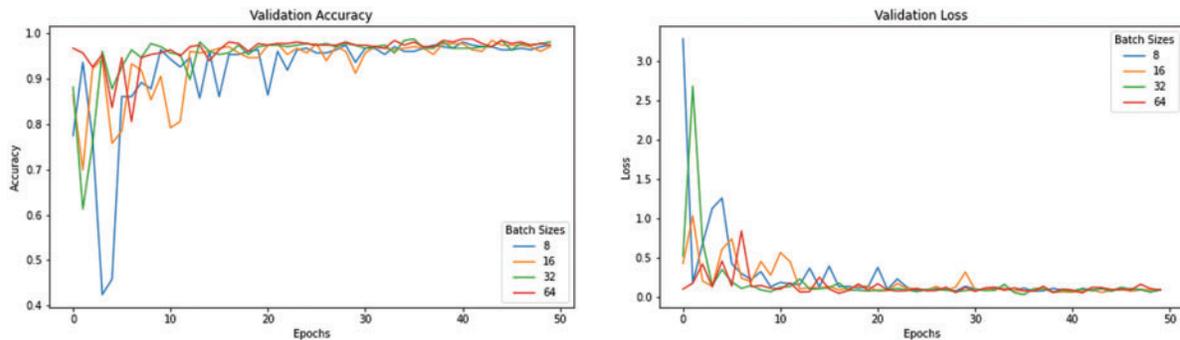


Figure 12: The validation accuracy and test loss respectively for batch sizes 8, 16, 32 and 64 in Resnet50 with attention

The better test accuracy was observed with learning rate 0.001 so we initialized learning rate with 0.001 and repeated the experiment by varying batch size as 8, 16, 32 and 64. From [Tab. 9](#), we found that with initial learning rate 0.001 and batch size 64, attention mechanism with Residual network (Resnet50) gave a test accuracy of 95.47%.

Table 9: Test Accuracy w.r.t batch size in baseline and Resnet50 with Attention

| Batch size | Learning rate | Epochs | Baseline | | Resnet50 | |
|------------|---------------|-----------|---------------|---------------|---------------|---------------|
| | | | Test accuracy | Test loss | Test accuracy | Test loss |
| 8 | 0.001 | 50 | 96.16% | 0.095 | 94.53% | 0.1545 |
| 16 | 0.001 | 50 | 95.98% | 0.110 | 94.62% | 0.1559 |
| 32 | 0.001 | 50 | 95.81% | 0.120 | 95.22% | 0.1964 |
| 64 | 0.001 | 50 | 95.13% | 0.1160 | 95.47% | 0.2409 |

Discussion: In our experiment, we observed the classification performance on the dataset using pre-trained classifiers with transfer learning. We also experimented with two CNN architectures, the first one was a baseline CNN and the other one was a Resnet50. Then we integrated the attention module in this architecture and observed the improvement in performance metrics and based on [Tabs. 5–9](#), we see a rise of 12% accuracy in a baseline CNN when the attention mechanism is applied and an 11% rise in accuracy when attention is applied with Resnet50. We also further experimented with varying learning rates and batch sizes where we got the highest accuracy of 95.47% with the learning rate being initialized with 0.001 and batch size of 64.

6 Conclusion

In this paper, we illustrate how attention mechanism can be applied for image classification. Then, we summarize recent breakthrough by other researchers using attention mechanism in the field of computer vision and (NLP). Moreover, we applied attention mechanism in baseline CNN and in Residual network (Resnet50) for classification of chest X-ray images for detection of Pneumonia. The results indicated that we observed significant improvement in classification accuracy and in other performance parameters even for a baseline CNN. In addition, we compared the attention network's performance with performance of retrained model and found that CNN network with attention

mechanism offered better results when compared with retrained transfer learning based architecture. In the chosen dataset, we faced the problem of data imbalance and limited size and we overcome it with data augmentation. Furthermore, we experimented with varying learning rate and batch size and got the **95.47%** test accuracy with initial learning rate being 0.001 and batch size 64. Other performance parameters were improved when attention mechanism was integrated with CNN. Quaternion CNN an extension of real valued CNN performs better for color image classification. Attention mechanism with Quaternion CNN can be future extension of the proposed work.

Acknowledgement: The authors would like to thanks the editors of CMC and anonymous reviewers for their time and reviewing this manuscript and we are grateful to the ERAWAN project for high-performance computers.

Funding Statement: This research work was supported by Chiang Mai University.

Conflicts of Interest: The authors declare that they have no conflicts of interest to report regarding the present study.

References

- [1] J. P. Cohen, L. Dao, P. Morrison, K. Roth, Y. Bengio *et al.*, “Predicting COVID-19 pneumonia severity on chest X-ray with deep learning,” *Cureus*, vol. 12, no. 7, pp. 478–499, 2020.
- [2] J. Shiraishi, Q. Li, D. Appelbaum and K. Doi, “Computer-aided diagnosis and artificial intelligence in clinical imaging,” in *Seminars in Nuclear Medicine*, vol. 41, no. 6, pp. 449–462, 2011.
- [3] M. Chowdhury, T. Rahman, A. Khandakar, R. Mazhar and M. Abdul *et al.*, “Can AI help in screening viral and COVID-19 pneumonia?,” *IEEE Access*, vol. 8, no. 2, pp. 132665–132679, 2020.
- [4] M. Gromet, “Comparison of computer-aided detection to double reading of screening mammograms: Review of mammograms,” *American Journal of Roentgenology*, vol. 190, no. 4, pp. 854–869, 2008.
- [5] E. Krupinski, “Computer-aided detection in clinical environment: Benefits and challenges for radiologists,” *Radiology*, vol. 231, no. 1, pp. 7–19, 2004.
- [6] A. Olatunbosun and S. Viriri, “Deep learning approach for facial age classification: A survey of the state-of-the-art,” *Artificial Intelligence Review*, vol. 54, no. 1, pp. 179–213, 2021.
- [7] J. Z. Wang and X. Xue, “Multi-class support vector machine,” *Support Vector Machines Applications*, vol. 134, pp. 23–48. Springer, Cham, 2014.
- [8] N. Wu, J. Phang, J. Park, Y. Shen, Z. Huang *et al.*, “Deep neural networks improve radiologists’ performance in breast cancer screening,” *IEEE Transactions on Medical Imaging*, vol. 39, no. 4, pp. 1184–1194, 2019.
- [9] R. Rout, P. Parida, Y. Alotaibi, S. Alghamdi and O. I. Khalaf, “Skin lesion extraction using multiscale morphological local variance reconstruction based watershed transform and fast fuzzy C-means clustering,” *Symmetry*, vol. 13, no. 11, pp. 2085–2105, 2021.
- [10] X. Yu, H. Chen, M. Liang, Q. Xu and L. He, “A transfer learning-based novel fusion convolutional neural network for breast cancer histology classification,” *Multimedia Tools and Applications*, vol. 81, pp. 11949–11963, 2022.
- [11] G. Suryanarayana, K. Chandran, O. I. Khalaf, Y. Alotaibi and A. Alsufyani, “Accurate magnetic resonance image super-resolution using deep networks and Gaussian filtering in the stationary wavelet domain,” *IEEE Access*, vol. 9, no. 2, pp. 71406–71417, 2021.
- [12] S. Singh and B. K. Tripathi, “Pneumonia classification using quaternion deep learning,” *Multimedia Tools and Applications*, vol. 81, no. 2, pp. 1743–1764, 2022.
- [13] G. Li, F. Liu, A. Sharma, O. I. Khalaf and Y. Alotaibi, “Research on the natural language recognition method based on cluster analysis using neural network,” *Mathematical Problems in Engineering*, vol. 2021, pp. 567–579, 2021.

- [14] A. Seemendra, R. Singh and S. Singh, "Breast cancer classification using transfer learning," *Lecture Notes in Electrical Engineering*, vol. 694, no. 3, pp. 689–699, 2020.
- [15] B. K. Tripathi, "On the complex domain deep machine learning for face recognition," *Applied Intelligence*, vol. 47, no. 2, pp. 382–396, 2017.
- [16] J. Cao, Q. Chen, J. Guo and R. Shi, "Attention-guided context feature pyramid network for object detection," arXiv Preprint arXiv, vol. 2, no. 1, pp. 258–271, 2020.
- [17] G. Yang, Y. He, Y. Yang and B. Xu, "Fine-grained image classification for crop disease based on attention mechanism," *Frontiers in Plant Science*, vol. 56, no. 3, pp. 2077–2096, 2020.
- [18] K. Cho, A. Courville and Y. Bengio, "Describing multimedia content using attention-based encoder-decoder networks," *IEEE Transactions on Multimedia*, vol. 17, no. 11, pp. 1875–1886, 2015.
- [19] S. S. Rawat, S. Alghamdi, G., Kumar, Y. Alotaibi and O. I. Khalaf, "Infrared small target detection based on partial sum minimization and total variation," *Mathematics*, vol. 10, no. 4, pp. 671–687, 2022.
- [20] H. Gill, O. Khalaf, Y. Alotaibi and S. Alghamdi, "Fruit image classification using deep learning," *CMC-Computers, Materials & Continua*, vol. 7, no. 3, pp. 5135–5150, 2022.
- [21] G. Brauwers and F. Frasincar, "A general survey on attention mechanisms in deep learning," *Ieee Transactions on Knowledge and Data Engineering*, vol. 56, no. 3, pp. 1–19, 2021.
- [22] Y. Alotaibi and A. Subahi, "New goal-oriented requirements extraction framework for e-health services: A case study of diagnostic testing during the COVID-19 outbreak," *Business Process Management Journal*, vol. 28, no. 1, pp. 273–292, 2022.
- [23] A. Galassi Andrea, M. Lippi and P. Torrioni, "Attention in natural language processing." *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 10, pp. 4291–4308, 2020.
- [24] S. Kumar, P. Kumar, M. Gupta and A. K. Nagawat, "Performance comparison of median and wiener filter in image de-noising," *International Journal of Computer Applications*, vol. 12, no.4, pp. 27–31, 2010.
- [25] S. Singh, G. Rathna and V. Singhal, "Indian sign language recognition on PYNQ board," *Recent Advances in Computer Science and Communications*, vol. 15, no. 1, pp. 98–104, 2022.
- [26] S. Rawat, S. K. Verma and Y. Kumar, "Reweighted infrared patch image model for small target detection based on non-convex \mathcal{L}_p -norm minimization and TV regularization," *IET Image Processing*, vol. 14, no. 9, pp. 1937–1947, 2020.
- [27] M. Valle and R. Lobo, "Hypercomplex-valued recurrent correlation neural networks," *Neurocomputing*, vol. 432, pp. 111–123, 2021.
- [28] Z. Niu, G. Zhong and H. Yu, "A review on the attention mechanism of deep learning," *Neurocomputing*, vol. 452, pp. 48–62, 2021.
- [29] M. Gupta and A. K. Courville and A. K. Nagawat, "Design and implementation of high performance advanced extensible interface (AXI) based DDR3 memory controller," in *IEEE International Conference on Communication and Signal Processing (ICCSP)*, India, vol. 17, no. 11, pp. 1175–1179, 2016.
- [30] M. Gupta, H. Taneja and L. Chand, "Performance enhancement and analysis of filters in ultrasound image denoising," *Procedia Computer Science*, vol. 132, no. 1, pp. 643–652, 2018.
- [31] M. Hardt, B. Recht and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," in *33rd Int. Conf. on Machine Learning, ICML*, New York, USA, pp. 375–389, 2016.
- [32] W. Pinaya, S. Vieira, R. G. Dias and A. Mechelli, "Convolutional neural networks", *Machine Learning*, vol. 67, no vol. 5, pp. 173–191, 2020.
- [33] P. Rodríguez, D. Velázquez, G. Cucurull, J. M. Gonfau, F. X. Roca *et al.*, "Pay attention to the activations: A modular attention mechanism for fine-grained image recognition," *IEEE Transactions on Multimedia*, vol. 22, no. 2, pp. 502–514, 2020.
- [34] S. Rawat, S. K. Verma and Y. Kumar, "Infrared small target detection based on non-convex Triple tensor factorization," *IET Image Processing*, vol. 15, no. 2, pp. 890–910, 2021.
- [35] B. Zoph, E. D. Cubuk, G. Ghiasi, T. Y. Lin, J. Shlens *et al.*, "Learning data augmentation strategies for object detection," in *European Conf. on Computer Vision*, Glasgow, U.K, pp. 566–583, 2020.
- [36] Y. Alotaibi, "A new meta-heuristics data clustering algorithm based on tabu search and adaptive search memory," *Symmetry*, vol. 14, no. 3, pp. 623, 2022.

- [37] S. Rajendran, O. I. Khalaf, Y. Alotaibi and S. Alghamdi, "MapReduce-Based big data classification model using feature subset selection and hyperparameter tuned deep belief network," *Scientific Reports*, vol. 11, no. 1, pp. 1–10, 2021.
- [38] M. Mohd, A. Sah, and M. Abulaish. "Deepsbd: A deep neural network model with attention mechanism for social bot detection," *Ieee Transactions on Information Forensics and Security*, vol. 16, no. 2, pp. 4211–4223, 2021.
- [39] A. Alsufyani, Y. Alotaibi, A. Almagrabi, S. Alghamdi and N. Alsufyani, "Optimized intelligent data management framework for a cyber-physical system for computational applications," *Complex & Intelligent Systems*, vol. 46, no. 2, pp. 1–13, 2021.
- [40] M. Umar, Z. Sabir, M. A. Z. Raja, M. Shoaib, M. Gupta *et al.*, "A stochastic intelligent computing with neuro-evolution heuristics for nonlinear Sitr system of novel COVID-19 dynamics," *Symmetry*, vol. 12, no. 10, pp. 1628–1646, 2020.
- [41] Q. Xu, Y. Xiao, D. Wang and B. Luo. "CSA-MSO3DCNN: Multiscale octave 3D CNN with channel and spatial attention for hyper spectral image classification." *Remote Sensing*, vol. 12, no. 1, pp. 188–203, 2020.
- [42] K. He, X. Zhang, S. Ren and J. Sun. "Deep residual learning for image recognition," *Pattern Recognition*, vol. 68, no. 3, pp. 770–778, 2016.
- [43] M. Gupta, S. Upadhyay and A. K. Nagawat, "Camera calibration technique using tsai's algorithm," *International Journal of Enterprise Computing and Business Systems*, vol. 34, no. 2, pp. 190–213, 2011.
- [44] M. Umar, Z. Sabir, M. A. Z. Raja, M. Gupta, D. N. Le *et al.*, "Computational intelligent paradigms to solve the nonlinear SIR system for spreading infection and treatment using levenberg–Marquardt backpropagation," *Symmetry*, vol. 13, no. 4, pp. 618–635, 2021.
- [45] K. Nisar, Z. Sabir, M. A. Z. Raja, A. A. A. Ibrahim, J. J. P. C. Rodrigues *et al.*, "Evolutionary integrated heuristic with gudermannian neural networks for second kind of lane–emden nonlinear singular model," *Applied Sciences*, vol. 11, no. 11, pp. 4725–4747, 2021.
- [46] M. Gupta, J. Lechner and B. Agarwal, "Performance analysis of kalman filter in computed tomography thorax for image denoising," *Recent Advances in Computer Science and Communications*, vol. 13, no. 6, pp. 1199–1212, 2020.
- [47] M. Gupta, H. Taneja, L. Chand Vishnu Goyal, "Enhancement and analysis in MRI image denoising for different filtering techniques," *Journal of Statistics and Management Systems*, vol. 21, no. 4, pp. 561–568, 2018.
- [48] M. R. Haque, S. C. Tan, Z. Yusoff, K. Nisar, C. K. Lee *et al.*, "Automated controller placement for software-defined networks to resist DDoS attacks," *Computers, Materials & Continua*, vol. 68, no. 3, pp. 3147–3165, 2021.
- [49] M. Gupta, L. Chand and M. Pareek, "Power preservation in OFDM using selected mapping (SLM)," *Journal of Statistics and Management Systems*, vol. 22, no. 4, pp. 763–771, 2019.
- [50] K. Nisar, Z. Sabir, M. A. Z. Raja, A. A. A. Ibrahim, J. J. P. C. Rodrigues *et al.*, "Artificial neural networks to solve the singular model with neumann–robin, dirichlet and neumann boundary conditions," *Sensors*, vol. 21, no. 19, pp. 6498–6521, 2021.
- [51] A. Kumar, S. Chakravarty, M. Gupta, I. Baig and M. A. Albreem, "Implementation of mathematical morphology technique in binary and grayscale image," *Advance Concepts of Image Processing and Pattern Recognition*, Springer, vol. 1, no. 1, pp. 203–212, 2022.