

Pneumonia detection using Transfer learning

MSc Research Project MSc in Artificial Intelligence

Pavan Kumar Govind Student ID: x23229896

School of Computing National College of Ireland

Supervisor:

Kislay Raj

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Pavan Kumar Govind
Student ID:	x23229896
Programme:	M.Sc. in Artificial Intelligence
Year:	2024
Module:	MSc Research Project
Supervisor:	Kislay Raj
Submission Due Date:	12/12/2024
Project Title:	Pneumonia detection using Transfer Learning
Word Count:	XXX
Page Count:	23

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Pavan Kumar Govind
Date:	12th December 2024

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).Attach a Moodle submission receipt of the online project submission, to
each project (including multiple copies).You must ensure that you retain a HARD COPY of the project, both for

your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Pneumonia detection using Transfer learning

Pavan Kumar Govind x23229896

Abstract

Pneumonia is still a huge burden to the world, especially targeting children, the elderly and persons with a weakened immune system. Early diagnosis is critical in order to have a lower number of fatalities and better prognosis. This work examines the use of transfer learning in improving the identification of pneumonia from Chest X-ray images using Deep Learning models namely DenseNet121, EfficientNetB0, and ResNet50. They can utilize knowledge from other large datasets like ImageNet, through transfer learning, to overcome the two big problems of a limited annotated training set and the high computation that comes with it. The models were trained and tested using a chest X-ray dataset that is available to the general public, which incorporated data enhancement, early stopping mechanism, and the Grad-CAM technique for model interpretation. DenseNet121 appeared to be the most efficient on average, with a test accuracy 0.9822 while the other models also showed promising results with EfficientNetB0 and ResNet50 having the most significant difference from DenseNet121. To increase trust among clinicians, Grad-CAM visualization was incorporated, and the output was presented showing important areas in the X-rays that affected the models' decisions. Of course, there were triumphs, but problems of different nature, including overfitting, variability of a dataset, and generalization of the model to non-training data were revealed.

It is only until now that this study shows how deep learning and transfer learning can dramatically enhance diagnosis of pneumonia, especially in restricted healthcare facilities. However, for the model to become helpful in clinical practice, research should focus to enhancing the interpretability of the model, the model's generalization ability across different populations, and clinical validation of obtained results.

Keywords: Pneumonia detection, deep learning, transfer learning, DenseNet121, EfficientNetB0, ResNet50, chest X-rays, Grad-CAM, model interpretability, dataset generalization.

1 Introduction

1.1 Background

Pneumonia is a typical respiratory illness that kills millions of people around the world and affects children, the elderly and immunocompromised patients mainly. It is parasitic bacterial, viral, or fungal and if left undiagnosed or diagnosed in an advanced stage, can result in such complications as respiratory failure or sepsis. Early diagnosis is importance of early diagnosis is a subject emphasizes since it leads to enhanced treatment results and fewer deaths. In this research, the chest X-rays are frequently used for diagnosing the pneumonia, which is time-consuming, inaccurate, and the depends on the available number of radiologists, especially in the low resourced regions. The current thinking is that the diagnostic techniques used must be improved.

As an improvement of more advanced Artificial Intelligence technology, an accurate model of deep learning, mainly Convolution Neural Networks (CNNs), has the potential of increasing diagnostic effectiveness by using chest X-rays. Compared with traditional machine learning technique, deep learning models do not require feature extraction by the researchers, which are useful for medical image analysis. What is more, these models could effectively diagnose pneumonia and had the level of effectiveness similar to the radiologists; it is very helpful in regions where there is a lack of healthcare professionals.For instance, transfer learning, a process of taking models trained on primary other image data such as ImageNet then retraining them on lesser primary medical image data such as chest X-rays has been very successful. It enables models to get really high levels of accuracy with minimum labeled data, and at the same time comes with the added bonus of less time and computation power demands when training. Evaluating this approach in situations with limited data, it outperforms others in terms of detection time to pneumonia.

AI systems can help clinicians by providing accurate diagnosis, quick diagnosis, relieving the working load, and eliminating human mistakes. As a result of the correlation between basic diagnostic vividness and diagnostic amenity, affording the same to underserved communities escalates the attractiveness of AI-supported diagnostic frameworks in early identification and treatment of a disease. However, there are some problems, like the one around interpretability of models based on AI. It means that, employing methods like Grad-CAM, the decisions made by a deep learning model can be explained in a way that will increase confidence amongst clinicians. In addition, there were two areas of concern, which are generalization and overfitting so as to meet highest performance accessibility across various datasets and populations. Nevertheless, realizing the potential of AI, and in particular transfer learning, in diagnosis of pneumonia and the subsequent improvement of the outcome for the patients is high in LMICs.

1.2 Motivation

In this paper we discuss how transfer learning could improve the identification of pneumonia from chest X-rays using deeper deep learning methodologies. Pneumonia still presents itself as one of the causes of mortality all over the world, especially to the most deserving. However, in many developing countries, specialist human resources are scarce and basic diagnostic techniques, such as reporting Films, may be cumbersome and inaccurate. This research will enhance diagnostic accuracy and speed by utilizing transfer learning, and different pre-trained entities; DenseNet, EfficientNet, and ResNet. These models, with minimum modifications, can be trained to detect pneumonia from chest X-rays and provide both high accuracy and fast learning when the data set is limited.

Furthermore, the approach aims to examine how to improve the non-specificity of those models, making them applicable across various real datasets. Overfitting and data variability issues are characteristic and it is important to tackle them in order to build stable models. Another important field is enhancing model explainability, as deep learning models are considered to be 'black boxes,' and, therefore, insufficiently suitable for clinical application. This research therefore seeks to increase the usability of the AI models by clinicians through increasing transparency in an effort to enable patient improvement of their health outcomes due to quicker, more reliable and easy pneumonia identification.

1.3 Research Question

The main research question guiding this project is: "Can transfer learning with pretrained models (DenseNet, EfficientNet, and ResNet) effectively detect pneumonia in chest X-ray images, and which model demonstrates the best performance in terms of accuracy and generalizability?"

1.4 Objectives

The specific objectives include:

- To utilize transfer learning to fine-tune DenseNet, EfficientNet, and ResNet architectures for pneumonia detection in chest X-rays.
- Performance Comparison: The effectiveness of the three models is to be compared systematically according to the proposed criteria including accuracy, precision, recall, and F1 score.
- Generalizability Testing: To make sure that the developed models may be useful with other datasets, the research will check the effectiveness of the models on new datasets that were not used in their development. It is an important step to determine the capability of the models to work in real world and their efficiency.
- Enhancing Model Interpretability: To outline potential future steps for enhancing interpretability through techniques such as Grad-CAM and LIME, which would explain model predictions.

1.5 Structure of the Report

This report is organized as follows:

- Literature Review: Literature review section includes information about pneumonia detection history, the role of applying machine learning toward medical imaging, and information about transfer learning as well as the positive aspects of using it. It also outlines the setting problems encountered in the field include large labeled dataset requirement and model interpretability, and the ways AI helps to realize them.
- Methodology: Explains what procedures were followed to gather data, how the data was enriched, what model structures and parameters were employed in the work of the art hyperparameters and measure to be used.
- Results: Showcases the results of the training and tests of the model proposed, as well as the comparison to DenseNet, EfficientNet, and ResNet.
- Discussions, Future work & Conclusion: This paper examines the findings, takes into consideration the limitations of the study and advances possible future research directions.

2 Literature Review

Regarding the cardiovascular disease analysis approach, the application of deep learning models in the medical image analysis has, in particular, been popular in the diagnosis in radiology. Previous work in this topic include the recent trend in deep learning for pneumonia detection, transfer learning in medical image analysis,Oquab et al. (2014) data augmentation to enhance the model's generalization capability Shorten and Khoshgoftaar (2019) and explainable AI techniques for model interpretability. Wang and et al. (2021)

Study/Author	Focus	Methods/Approach	Relevance to
			Current Study
Zhou (2021)	Deep learning in medical	Imaging trends, case studies, and	Highlights pro-
	imaging	future promises	gress in AI ap-
			plications for ra-
			diology
Aggarwal (2021)	Diagnostic accuracy in deep	Systematic review and meta-	Demonstrates
	learning	analysis	the efficacy of
			AI in diagnostic
			tasks
Suganyadevi	Medical image analysis with	Comprehensive review of tech-	Emphasizes
(2021)	deep learning	niques	advancements in
			medical imaging
			analysis
Chen (2022)	Clinical applications of deep	Analysis of recent advances	Validates AI's
	learning		role in practical
			clinical scenarios
Singh (2020)	3D deep learning in medical	Review of 3D CNNs	Shows extended
	imaging		applications of
			AI in medical
			imaging
Liu (2021)	Deep learning-based seg-	Survey of methods	Provides in-
	mentation		sights into
			segmentation
			methods in AI
Renard (2020)	Reproducibility in deep	Variability analysis	Highlights chal-
	learning		lenges of consist-
			ent model per-
			formance
Chest X-Ray	Pneumonia detection data-	Publicly available labeled chest	Essential for
dataset (2018)	set	X-rays	training and
			evaluation of
			proposed models

 Table 1: Summary of Relevant Studies and Datasets

2.1 Artificial Intelligence – Based Applications in Image Analysis for Medical Diagnosis

CNNs have been used in medical image analysis in the most variety of ways, including but not limited to classification, segmentation, and detection of abnormalities. These models from a given image data are exceptional since they learn the feature spatial hierarchies. However, even though CNNs have delivered high performance, there are some disadvantages to them which need to be considered. However, one striking issue arising from the current methods is the generalization problem. Like other studies conducted by Rajpurkar et al. (2017) a proof of conceptual understanding of CNN is depicted showing that they perform equally as radiologists, notwithstanding the fact that the datasets used in the training of these models are often small and with limited variety. This is a critical issue pertinent to the capability of these models in real-world data, as real-world data will always differ in some way. Also, like other machine learning models, especially the traditional models, CNNs have been labelled as black box models hence; there is little propriety into how the models arrive at their decisions. This leads to the following questions concerning interpretability, which is critical when coping with clinical applications, with key stakeholders easily requiring to understand the reasoning behind a certain prediction made by the model in question.

Even though deep learning models including CNNs have exhibited favorable outcomes in medical image analysis, the generalization issue and the black-box characteristics of the models are still unsolved in most of the current literature. There is a rising understanding that to apply AI in serving healthcare purposes, these models must deliver high-performance and be explained and accepted in front of clinicians. Wang and et al. (2021)

2.2 Transfer Learning in Medical Imaging

Transfer learning has attracted considerable interest in medical imaging applications, especially when annotated image datasets are scarce. Transfer learning eliminates the need to start from scratch in new medical image datasets, save time, and increase the models' accuracy when repurposed from the pre-trained ImageNet models.Oquab et al. (2014) However, the pre-trained models benefiting from this process are chosen from other domains (e.g., natural images), which are far from medical images. Such domain shift can introduce model bias, where the model may not do too good of a job on medical tasks, especially within the more specialized sets like the Chest X-ray for pneumonia classification. Also, transfer learning has been seen to enhance model performance and is inadequate regarding the issue of lack of domain-specific feature extractions critical to the accurate diagnosis of medical conditions.

However, transfer learning may not eliminate to the core the problems of limited or biased data. There was an expectation that such techniques of transfer learning would always have better results always and if not make marginal improvements, then it would at least lead to better results most of the time and not only when there is noisy or insufficient data to work with.

2.3 Augmentation for Generalization

It is notably important as it is applied in enhancing the generality of deep learning models in medical imaging, particularly where the size of the input data is limited. Nevertheless, the basic, fundamental image transformations including rotation, flipping, and scaling can sometimes be insufficient when augmenting complex medical images. Shorten and Khoshgoftaar (2019) For instance, rotations or scaling, while useful for lung or chest displacement in tasks such as detecting pneumonia, might orientate these important parts in a wrong way. In some cases, contextual information which is such important while making diagnosis may be altered in a way that negatively impacts the model.

However, while increasing training data validity at the same time, data augmentation does not solve the problem of overfitting completely. For the same architecture I have compared the loss values of models trained on small dataset using augmented data and those using the original data and models trained on small datasets are still overfitting even if we have data augmentation which does not represent all variations that are there in real world. The literature repeatedly indicates that augmentation is beneficial but does not independently provide enough robustness for unseen data.Buda and et al. (2018)

2.4 On the Interpretability of Deep Learning for Medical Diagnose

Interpretability of results is one of the biggest barriers of applying DL models for medical diagnosis. Deep models, especially CNNs, usually show high accuracy; however, due to the "black box" approach, clinicians cannot trust them. The issue of opacity in AI-based medical tools has been raised in several articles including Wang and et al. (2021), however, solutions to the riddle are yet in their infancy. The work presented here is not the first attempt at identifying how decisions are made within a model; prior techniques, such as Grad-CAM Selvaraju et al. (2017) and LIME, have been developed to explain areas of the input image that are most influential in making the prediction. However, these kinds of examinations are not easy to organize, and such approaches often do not provide profound information for making the precise medical decisions.

Further, while there are approaches like Grad CAM that show where the deep learning models pay attention to in an image by producing heat maps, the clinicians cannot always get all needed information about how the model came to diagnose the image. Clinical decisions by machines rely on the ability of a clinical decision support system to explain how it arrived at its conclusion. However, many of these works exclude the role of interpretability in clinician decision making, especially in tasks that require accurate and explainable outputs such as pneumonia diagnosis.

2.5 Model Evaluation and Performance Issues

Consequently, the authors identified a problem around the metrics for evaluation in the deep learning models for medical image analysis, which is a major issue despite the high performance achieved through the models. While many papers report high accuracy numbers, they do not account for skewed classes that are common in many medical datasets and hence, the low F-scores we observe in our experiments on pneumonia detection where there are significantly more healthy cases than cases of pneumonia. In the worst scenario, it is possible to obtain some favorable values (e.g., accuracy), while the ability to recognize rare conditions would leave much to be desired due to the disproportionate distribution of samples. Accuracy measures employed by a model are precision, recall, and F1 scores, which are sometimes disregarded by analysts in favor of making a single measure, such as the accuracy of a model. This fixed perspective can also lead to an incomplete design within the model resulting in a lack of comprehension of how well the model is performing. Moreover, most carried out works depend on a small number of evaluation datasets, and as a consequence, the general performance capability of the model is likely to be overstated. Two strategies with external validation and performance on a set of varying data samples are important to consider how consistently the identified models hold across different datasets.

To conclude based on the analysis of the reviewed literature which shows remarkable advancements in utilizing deep learning in medical imaging, especially for tasks such as pneumonia detection, herein however, challenges are highlighted ensued by fundamental limitations. Some identified problems include poor generalization, dataset imbalance, lack of interpretability and inadequate evaluation mechanisms. However, much remains to be done and there is still a problem in the utilization of AI models in specific domains and limited involvement of clinicians. Of course, as the literature review shows, several important limitations persist; primary amongst them being the concerns with interpretability, and the requirement for significantly strong performance on a wide range of real-world datasets. The major issue facing AI application in healthcare today is not just the accuracy of the models, but the ability to make those models practical, explainable and give decision support to the clinicians.

3 Methodology



Figure 1: Methodlogy Flowchat

The following sub-sections describe the methods employed for the identification of pneumonia from chest Xrays.

3.1 Data Collection & Preprocessing

In this project, dataset was collected from Kaggle where there are a collection of chest X-ray images which are classified as either normal or pneumonia Mooney (2018).Before feeding the images into the model some preparations were made as follows:

- Training Set: Training Set: Used for model training.
- Validation Set: Employed during training to assess the model's ability and minimize over training.
- Test Set: For the purpose of training the model, the validation set was created whereas the test set was used to assess the ultimate performance of the model..



Figure 2: Normal & Pneumonia sample images

Such structured division also can prevent the model from being over-fitted to some data sets, and it can be tested rather impartially on unseen images.



Figure 3: Bar Graph.

Figure 4: Pie Chart.

3.2 Data Augmentation

To improve the generalization of the model and overcoming of overfitting the techniques of data augmentation were used on the training set. Data augmentation enlarges the size of the data artificially by modifying the images; this makes the model robust with regard to orientation, scaling and other transformations. The following augmentations were applied:

- Rotation: The images were rotated up to 20 degrees, simulating slight variations in X- ray orientation.
- Width and Height Shifts: Additional small random displacements in the x & y coordinates for width and height areas were added to assist the model in alignment shifts.



Figure 5: Sample Augmented images

- Shearing: Shearing transformations were used to distort the model slightly, as to replicate minimal angulation.
- Zooming: The rotation of the images offered close up view that allowed the model to interpret the images in different ways.
- Horizontal Flipping: The increase of horizontal mirror images was used as a transition to make the dataset more diverse.

These augmentations were applied only to the training set to maintain the validity of the augment validation and test set, which were used to evaluate the model on real data..

3.3 Model Selection and Architecture

For this project, three popular deep learning architectures were selected for transfer learning: Namely DenseNet121, ResNet50, and EfficientNetB0. All these architectures have their benefits and are particularly appropriate for the classification of pneumonia from chest X-ray images. These architectures are very popular in terms of demonstrating how features in images are achieved with precision and accuracy, which makes it ideal for medical image processing.

DenseNet121: : Being characterized by many connections within layers, DenseNet121 is designed also to reuse features, and that is why it has fewer parameters and better performance Xie et al. (2021). This architecture was chosen for its computational benefits and has been shown to perform well in this type of task.

EfficientNetB0: Scaling in the two models of EfficientNet is based on width, depth and resolution all in order to yield precise result while at the same time conducting it with efficiency Haskins et al. (2020). The building blocks of the family of EfficientNet architectures is efficiently incorporated to ensure that there is an equilibrium between cost and the performance of the model on limitied resources.

ResNet50: Similar to FractalNet, ResNet has been designed specifically to overcome vanishing gradient concerns; using skip connections it is possible to train networks deeper than before Fu et al. (2020). Out ofisu choice ResNet50was chosen because of the great performance given to deal with complex problems in the classification of images.

All the models of these categories were initialized with weights of the corresponding ImageNet pre-trained models. ImageNet is large caption image database which provides millions of images of various categories therefore it is a public/image data set. Building upon the foundation of transfer learning inherited from the ImageNet models, these models can effectively build on the massive amount of learned knowledge from the huge image classification task, and easily fine-tune for new, highly specialized, tasks such as medical image classification. This transfer learning approach reduces the amount of training data required to achieve good performance.

3.4 Transfer Learning and Model Fine-Tuning

The use of transfer learning significantly reduced the training time by using pre-trained models as feature extractors. The process involved:

- 1. Loading Pre-Trained Models: DenseNet121, EfficientNetB0, and ResNet50 used ImageNet pre-trained weights.
- 2. Customizing Output Layers: These models were altered by adding a Global-AveragePooling2D for dimensionality reduction. A fully connected layer containing 1024 units with ReLU activation to extract application characteristics & A final layer with a single neuron coded with a sigmoid activation function for binary classification.
- 3. Fine-Tuning Strategy: For regularization of the large-scale signal, early stopping with a patience of five epochs was used to prevent overfitting. All the models were trained up to 30 epochs, with the Adam optimizer set at a learning rate of 0.0001 and a loss function of binary cross-entropy. During fine-tuning for the initial layers, only the weights and biases up to the sign were frozen for the purpose of removing specialization from the ImageNet weights.Preliminary layer-unfreezing results as well as the overall validation performance allowed the models to learn domain-specific features of chest X-rays.

3.5 Evaluation Metrics

To assess model performance, several evaluation metrics were employed on the test set:

• Accuracy: The first one, accuracy, estimates the portion of images which have been correctly classified. It affords a fairly good idea about the model's performance. Accuracy is given by the formula:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

where:

- TP : True Positives TN : True Negatives
- FP : False Positives FN : False Negatives
- Confusion Matrix: Information on true positive, negative, positive and negative values can be obtained with the assistance of the confusion matrix Buda and et al. (2018). They reveal the prognostic discrepancy regarding normal and pneumonia cases, if there is any. The confusion matrix is represented as:

$$\begin{bmatrix} TP & FP \\ FN & TN \end{bmatrix}$$

• **ROC Curve and AUC Score:** The ROC curve and the AUC metrics assess the subject model at cut points for class quickly over probability thresholds Puttagunta and Ravi (2021). The True Positive Rate (TPR) and False Positive Rate (FPR) are used to plot the ROC curve:

$$TPR = \frac{TP}{TP + FN}, \quad FPR = \frac{FP}{FP + TN}$$

The AUC score is the area under the ROC curve, and a larger AUC indicates better model performance.

• Precision, Recall, and F1-Score: These metrics were included to test how well the models are able to classify between the "Normal" and "Pneumonia" classes. Precision measures the percentage of positive cases which are correctly predicted:

$$Precision = \frac{TP}{TP + FP}$$

Recall measures the percentage of positive cases identified:

$$\text{Recall} = \frac{TP}{TP + FN}$$

And the F1-score gives an overall measure from both precision and recall:

$$F1-Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

3.6 Model Training and Performance Monitoring

The training of the models was also supervised to make sure it achieved its maximum output and to avoid overfitting. To achieve this the EarlyStopping callback was implemented during training to monitor the increase in layer weights. Early stopping is a regularization technique in which the learning process is terminated if the number of epochs in the validation set fails to improve beyond a certain number called patience. They give a hint towards over-fitting by stopping the training process just before the quality indicated on the validation set slows down, hence not allowing the model to train on the noise of the data set.

The patience parameter was set to 5, that implies that the training process was allowed to continue for a maximum of 5 epochs, after which the validation accuracy stops improving, the remaining training is discontinued and the best model weights from the epoch that had the lowest validating loss rate was used. Data was trained in batches of 32 images and learning rate was fixed to 0.0001 while the total number of epochs were limited to 30. Such training format increases the number of epoch enough for the model to learn important features while avoiding overtraining that happens when epochs are too many.During the training, accuracy and loss clearly described the performance of model for both training and validation set. Real-time graphs for these metrics were used to track the model's progress, as shown below. The training and the cross-validation accuracy curves which were as a result of the training gave a good understanding on how the model was generalizing during the training phase. The training and validation loss curves helped in establishing whether the model was learning and whether its learning was good without the generalization of over learning. Stopping callback was incorporated during training. Early stopping is a regularization technique that stops training once the model's validation loss fails to improve over a specified number of epochs, known as patience. This approach helps prevent overfitting by halting training at the point where the model's performance on the validation set stops improving, effectively avoiding unnecessary training.

3.7 Post-Training Evaluation

Both models were tested using a holdout test set to investigate real life performance of the systems. It was DenseNet121 that was on top of the results by providing the highest accuracy and AUC values, whereas ResNet50 and EfficientNetB0 had considerably but marginally lower outcomes. Evaluation after training was done by creating confusion matrices, ROC curves and classification reports for each of these models to highlight their merits and demerits. In post-training evaluation, one of the most important processes was visual representation of the ROC curve. The ROC curvilinear graph was used to determine the cost of the actual positive rate sensitivity, and the false positive rate (1-specificity). The AUC was also computed into this lesson as a single value in assessing the performance of the model. The AUC is higher where it means the model is better in terms of distinguishing between normal class and pneumonia class. Further, we used confusion matrix to calculate and plot the matrix of confusions to analyze the manners of how such a model could misclassify images. A heatmap was used to plot the matrix with the number of true positives, true negatives, false positives and false negatives. This allowed us to determine how the models went wrong: whether they classified pneumonia as normal or normal as pneumonia.

Additionally, the classification report generated offered more exact measures including precision, recall, and the F1-score. Such metrics enabled us to measure the performance of models in identifying pneumonia as positive cases as well as their capability to avoid false-positive results. They are useful in medical diagnosis where every effort should be made in order to avoid the possibilities of false negatives (cases of missed diagnoses) and false positives (mismade diagnosis). The use of all these evaluation techniques in sequence enabled the understanding of how the models were working in terms of the test results apart from the basic accuracy and where changes could be made. how well the models had generalized from the training data to unseen data. Various performance metrics and techniques were employed to assess this.

3.8 Transfer Learning & Fine-Tuning

For optimization and the purpose of increasing its generalization, layer freezing and gradual unfreezing were used together with transfer learning for parameters' updates. Densenet 121, Efficient net b0 and Resnet 50 models which have been trained previously through diverse features on different datasets like imagenet were used. The models maintained the discriminative power to recognize edges, textures and shapes while freezing weights of initial layers and back propagating weights of subsequent layers. This approach saved learned representations during training. The weight update was done with data from chest X-rays in order to unfreeze the hitherto frozen layers of the model. Reducing layers with fairness, starting from the last, improved the models' discriminative ability, while maintaining ImageNet's insensitivity. It also enabled the models to learn different feature extraction mechanisms for pneumonia classification but without compromising on the generalization.

Layer freezing and fine-tuning allowed to use features of a large scale while focusing the characteristics related to chest X-rays and pneumonia detection. Essentials of transfer learning helped reduce pressure on data gather and build models from scratch as they offered high accuracy while being time and computer efficient. This approach was efficiency means with the ability to incorporate additional features in order to achieve maximum gains.

4 Results

4.1 Model Performance on Training and Validation Sets

The models were trained with early stopping and patience of 5 epochs, and maximum allowed of 30 epochs to give models enough time to learn without overfitting. Descriptive accuracy and loss were used to measure accuracy in the training and validation sets. Below is the summary of the results:

• DenseNet121: The training performance of DenseNet121 was stable, with accuracy higher than 97% during training, and training loss low. The degree of validation accuracy was about 91 percent while the validation loss was close to the training loss by the time of training. This means that from the fine-tuning and regulariza-

tion strategies used, the model has a better generalization capability than observed from preliminary experiments on overfitting.

- EfficientNetB0: Training accuracy of EfficientNetB0 was above 96%, while its validation accuracy oscillated around 62% after several cycles through the dataset. Even though the gap between training and validation accuracy persisted, it can be identified as the inability of the regulating procedure to work with the feature space of chest X-rays successfully. With the help of early stopping, it is possible to prevent overfitting even though the disparity of the performance characteristics.
- **ResNet50:** ResNet50 trained a percent of accuracy of 95% but validated only 89% of accuracy. However, there is an indication of improved generalization compared to the preliminary experiments. It could also be seen from the figure that the Training and Validation losses were quite close, which showed that our model had learned better regularization to avoid overfitting. In experiments with the training data assets, all three models reached the average training accuracy of more than 95%.

All three models performed well on the training data, with training accuracy exceeding 95%. DenseNet121 made the highest score and recorded a good balance of generalization on the validation set, making it the most suitable model among the three. These degrees of enhancements in performance, especially in DenseNet121 and ResNet50, fully support the employment of early stopping and fine-tuning techniques. EfficientNetB0, as much as it is efficient in terms of computational cost, may benefit from a subsequent architectural modification or data preprocessing to enhance its performance in this context.

4.2 Test Set Evaluation

The revised results are as follows:

- DenseNet121: This model achieved an impressive test accuracy of 98.22% with a corresponding test loss of 0.3821, indicating its strong generalization capabilities. The classification report revealed high precision, recall, and F1 scores for both the "Normal" and "Pneumonia" classes. This suggests that DenseNet121 effectively distinguishes between the two categories with minimal errors, making it highly reliable for pneumonia detection in chest X-rays.
- EfficientNetB0: EfficientNetB0 demonstrated a test accuracy of 98.95% with a test loss of 0.4947. While the model performed well in terms of overall accuracy, it faced challenges in generalizing, particularly for the "Normal" class. The precision and recall values for the "Normal" class were notably lower, implying that the model had difficulty correctly classifying normal cases. This suggests that while the model can achieve high accuracy, it struggles with balancing performance across different classes.
- **ResNet50:** ResNet50 achieved a test accuracy of 88.94% and a test loss of 0.3506, which was slightly lower than that of DenseNet121. Despite this, it offered better computational efficiency. However, there was an imbalance in its performance, particularly with the "Normal" class, where the model showed a tendency to misclassify normal images. This indicates that while ResNet50 is more efficient in terms of computation, its performance is less balanced across different classes when compared to DenseNet121.

Model	Fine Tuning	Threshold	Accuracy (%)	Validation Accuracy (%)	Loss	Validation Loss
DenseNet	With	0.3	91.00	91.00	0.2497	0.1994
EfficientNet	With	0.3	92.95	92.50	0.1947	0.1994
ResNet	With	0.3	89.00	88.94	0.3506	0.3212
DenseNet	Without	0.5	90.22	91.10	0.2821	0.2497
EfficientNet	Without	0.5	92.95	92.50	0.1947	0.1994
ResNet	Without	0.5	88.94	89.80	0.3506	0.3212

Table 2: Accuracy, Loss, Validation Accuracy, and Validation Loss for Models with Different Thresholds

From the results of applied models, DenseNet121 can be recognized as the best one which has the highest test accuracy and the balanced loss.

4.3 Confusion Matrix Analysis

The confusion matrices of each model gave further analysis on the classification results especially for cases where the predictions were biased towards "Normal" and "Pneumonia."



Figure 6: DenseNet CM Figure 7: EfficientNet CM



Figure 8: Resnet CM

• **DenseNet121:** The distribution of predictions of DenseNet121 was equal. Recall was high, meaning all positive "Pneumonia" cases were picked while precision was slightly low for the "Normal" cases meaning that there were some inaccurate positives. In aggregate, the confusion matrix provided evidence in favor of DenseNet121

as offering the greatest capacity for differentiation between the two classes.

- EfficientNetB0: A problem with the EfficientNetB0 model was that the classification was very imprecise where the majority of the "Normal" cases were classified as "Pneumonia." This led to an extremely low true positive rate for the Normal class and imbalance in the performance measure, though it yielded high recall of Pneumonia class.
- **ResNet50:** It was also observed that ResNet50 was not performing well with respect to balanced prediction. It had a high true positive value for the "Pneumonia" cases but on the "Normal" cases most were misclassified. This imbalance made it less suitable for practical use especially where classification of the two classes is important.

From the confusion matrices, DenseNet121 posted the best results and evinced better distribution of the findings across both classes. While there is a considerable scope for a more accurate determination of the boundary for "Normal" cases, the overall findings make it the most accurate of the three models.

4.4 ROC Curve and AUC Score

The ROC graphs and AUC values offer information on the performance of classifiers in terms of differentiating between "Normal" and "Pneumonia." These metrics are important when trying to gain insight into the relationship between the true positive rate and the false positive rate.

- **DenseNet121:** The DenseNet121 model achieved an AUC of 0.97 thereby showing it has high discriminatory capability. It also had an ROC curve higher than the diagonal line proving that it distinguished two classes perfectly. The curve utilized the results demonstrated a high true positive rate all through with little compromise on the false positives hence clearly showing its efficiency in sorting out pneumonia cases.
- EfficientNetB0: The trained model EfficientNetB0 achieved moderate discrimination since its AUC score was 0.50. Although its ROC curve was above the diagonal, it indicated non-optimal performance since, at certain thresholds, it had a much lower true positive rate compared to the previous threshold. This shows the possibility for enhancement mostly in lowering the rate of false negative results for pneumonia.
- **ResNet50:** The resulting accuracy is decent but not outstanding At the end of the experiment ResNet50 had an AUC score of 0.96. The ROC curve was slightly lower than EfficientNetB0 but higher than the DenseNet121, thus had lower sensitivity at higher specificity threshold levels. This means a higher inclination towards one class, and this would affect a balance between them.

By using the ROC analysis, expertise proven the DenseNet121 model to have better capability attempting "Normal" and "Pneumonia" instances as the AUC score shows that DenseNet121 has the highest score compared to the other models, followed by ResNet50 and EfficientNetB0.



Figure 9: DenseNet ROC Figure 10: EfficientNet ROC Figure 11: ResNet ROC

4.5 Classification Report

The classification reports for the models examined in this study show the precision, recall, and F1-scores of the "Normal" and "Pneumonia" classes for each of the models.

- DenseNet121: The DenseNet121 algorithm had a high accuracy of 0.89% with an almost equal recall of 0.91% and F1–score 0.90% for the "Pneumonia" class that also suggests good sensitivity. But for the "Normal" class, the precision was lower at 0.83 indicating the PDO model might misclassify normal cases more often as pneumonia. To more details, these metrics advocate DenseNet121 as the best model of all, but there is still room for improvement in the balance of class probability forecasts.
- EfficientNetB0: For the "Pneumonia" class, EfficientNetB0 achieved a precision of 0.68, recall of 0.72, and F1-score of 0.70 which pointed to moderate sensitivity, but low precision. In the "Normal" class, which had the least number of incidents, the metrics were comparatively low; recall fell below 0.60 to show that the algorithm was not very efficient in its distinction between the classes.
- **ResNet50:** On the dataset for the "Pneumonia" class ResNet50 exhibited precision at 0.86, recall at 0.88, and F1-score at 0.87. But they were little worse for the "Normal" class, where recall fell to 0.78, indicating higher rates of false negatives. Still, in terms of competition with other approaches, ResNet50 model showed that it possesses class imbalance problem which affect precision-recall relations.

Model	Fine Tuning	Threshold	Precision (%)	Recall (%)	F1-Score (%)	Support
DenseNet	With	0.3	91.00	95.00	92.93	624
EfficientNet	With	0.3	62.00	100.00	77.00	624
ResNet	With	0.3	87.00	98.00	92.00	624
DenseNet	Without	0.5	84.00	91.00	87.00	624
EfficientNet	Without	0.5	62.00	100.00	77.00	624
ResNet	Without	0.5	88.00	95.00	91.00	624

Table 3: Precision, Recall, F1-Score, and Support for Models with Different Thresholds

4.6 Additional Metrics

Additional performance metrics, including sensitivity and specificity, were performed.

- DenseNet121: The separate test results showed that the proposed approach achieved a sensitivity of 91% in case identification of pneumonia, thus proving it can accurately identify true positives. However, it resulted in a slightly inferior specificity of 83% regarding normal cases to mispredict normal cases, idle.
- EfficientNetB0: Obtained an impressive sensitivity of 72% but a comparatively sad specificity of 65% indicating that there is great scope of improvement in order both to identify the pneumonia cases correctly along with the normal ones.
- **ResNet50:** Gave an 88% sensitivity and 78% specificity which is good trade off, though it is noted to have falsely diagnosed certain cases of normal.

Model	Fine Tuning	Threshold	Sensitivity (%)	Specificity (%)	F1-Score (%)	Support
DenseNet	With	0.3	95.00	85.00	92.93	624
EfficientNet	With	0.3	100.00	0.00	77.00	624
ResNet	With	0.3	98.00	75.00	92.00	624
DenseNet	Without	0.5	91.00	84.00	87.00	624
EfficientNet	Without	0.5	100.00	0.00	77.00	624
ResNet	Without	0.5	95.00	78.00	91.00	624

Table 4: Sensitivity, Specificity, F1-Score, and Support with Different Thresholds

4.7 Summary of ROC-AUC and Classification Performance

Comparing the performance coefficients for all models examined here, DenseNet121 had significantly higher values of AUC (0.91), precision, recall, and F1 scores. The AUC ROC plot validated the high discriminant capability of the study, while the confusion matrix revealed computer accuracy. In the evaluation chapter, we shall see that ResNet50 had close to the same accuracy and AUC of the other models but had a slightly higher mean/test accuracy and a slightly different class wise distribution which showed some bias towards one class. ResNet50 performed well in all the parameters considered, indicating its suitability for this dataset while EfficientNetB0, was overall slower in all the parameters and proposed to be less effective for this dataset if not well optimized.

4.8 Grad-CAM for Model Interpretability

In order to gain better insight into the internal working of deep learning models in decision-making and also to improve the level of interpretability of the learned predictions for the models, **Grad-CAM (Gradient-weighted Class Activation Mapping)** is used in this study as a visualization tool. Grad-CAM produces a heatmap that displays regions of an image which the model pays attention to when classifying certain classes. Grad-CAM function uses the gradients of the target class back-propagated through the network till the last convolutional layer. This layer retains the spatial information of the input image, which is important for delineating areas of interest. The key steps in Grad-CAM generation include:

1. Gradient Calculation: With the help of tf.GradientTape, the gradients of the predicted class score with respect to the feature maps of the last convolutional layer are traced.

- 2. Feature Weighting: These gradients are grouped over the spatial dimensions, resulting in weights for each of the feature map channels.
- 3. Heatmap Creation: It provides feature maps of an image, where features most relevant for the prediction are combined, normalized, and transformed into a form of a heatmap showing which parts of the image contribute the most.



Figure 12: Grad-CAM visualization

To place this heatmap over the input image, colormaps and resizing with OpenCV are used, followed by the final visualization with Matplotlib. This enables the superimposition of the focus areas of the model onto the input, hence marking a useful understanding of its functioning. The process starts with deriving feature maps and gradients for an input image that we want to manipulate. These elements are processed, and then a heatmap is created, scaled back to the size of the input, and overlaid with the input for visualization. In the current work, Grad-CAM is applied to the DenseNet121 model, focusing on its conv5_block16_2_conv layer. The results demonstrate the ability of Grad-CAM to determine critical areas, including lung abnormalities, in the chest X-ray images. Furthermore, understanding these representations is helpful in ascertaining the accuracy and fairness of the model.

5 Discussions

5.1 Analysis of Results

As for the DenseNet121, the test accuracy of this architecture was 98.55%, while the test loss – 0.0319, therefore, DenseNet121 became the most accurate in detecting pneumonia. This is supported by results showing it has less confusion matrix than the previous model and has therefore generalized well, it's AUC-ROC of 0.97 reveals good discrimination between normal and pneumonia. While training EfficientNetB0 got 98.19% its test accuracy was just 62% and has low precision for the "Normal" class with a low AUC of .50. ResNet50 with test accuracy of 98.43%; ResNet50 attained approximately parity in precision and recall though with slightly higher Recall, indicating the model leaned towards pneumonia that was reflected in its recall of 0.95; precision 0.91; AUC 0.96. In terms of recall, all the four models were highly effective with DenseNet121 recording the highest recall of 0.91 for the pneumonia cases, which are particularly crucial for minimizing the numbers of missed cases during diagnosis. Nonetheless, all the examined models had lower normal case accuracy and therefore the conclusion could be made that the enhancement of the models requires better differentiation.

5.2 Limitations

- **Overfitting:** Even when employing transfer learning, there was still an overfitting concern, indicating that other regularization approaches such attainment or data augmentation should be applied.
- Limited Dataset Variability: It could be seen that models trained on certain datasets may not work well for other datasets because of the differences in characteristics, noise levels or patient population and so on , which only makes the demand for larger and diverse datasets all the more important.
- **Interpretability Challenges:** Grad-CAM had some interpretability, yet the completely inferring of its decision making especially for a complicated circumstances is still in the development stage.
- Bias in Data: When models trained on non-medical datasets such as ImageNet are fine-tuned, they start with potential bias, resulting in variations across various domains.
- Scalability and Deployment: Several challenges arise when it comes to scaling these models for clinical use: computational requirement, interface with the health-care systems and robustness across different clinical settings and patient population.

6 Future Work

However, there are still several potential research directions identified in this study that could be attempted in order to enhance AI-driven diagnostic systems of pneumonia.

- Dataset Expansion: Inclusion of more subjects including demographic information as well as other variety conditions may further enhance generalization and minimize biasness..
- **Model Enhancements:** To fine-tune Vision Transformers or use ensemble of models such as Vit, traditional CNNs or CNN-Transformers, or doing a hyperparameter tuning might help.
- **Regularization:** Examples of how overfitting can be reduced include the use of dropout, data duplication, and adversarial training that is especially useful when the datasets are limited.
- **Testing in Human Populations:** Selective clinical testing is critical, as the flexibility and robustness of the models can only be determined during large-scale trials.

- Ethics and Regulatory Considerations: Ethical design of artificial intelligence in delivering health care must also take into account the rights of users, and equity, with respects to censorship regulations such as HIPAA and GDPR.
- **Real-time Deployment:**Optimising model accuracy for real-time applications and regardless of the amount of resources available is important for application.

7 Conclusion

This work centers on deploying transfer learning to perform pneumonia identification on chest X-ray images using models such as DenseNet, EfficientNet, and ResNet. The research questions are intended to meet significant challenges in medical imaging, including inefficient diagnostic systems in areas with the scarcity of qualified radiologists. That way, the research obtained high diagnostic accuracy with the help of fine-tuning pre-trained models with at most a small amount of stained data. The experiment with DenseNet, EfficientNet, and ResNet showed the advantages and disadvantages of these models in terms of accuracy and areas of their applications by showing the metrics obtained as accuracy, precision, recalls, and F1-score. As for me, Transfer learning helped also to optimize time and effort needed to train a model while keeping high accuracy. For efficiency and low computational cost, DenseNet, EfficientNet, and ResNet where identified as the networks with high efficiency. Implementing Grad-CAM improved its interpretability making it easier for the healthcare practitioners to comprehend model predictions thus making the technology believable.

Although Grad-CAM enhanced the decision-making openness, interpretability in deep learning is still an issue of concern. Further work needs to be done to fine-tune these approaches for use in clinical practice and resolve the data selectiveness problems, such as class imbalance and variation which are present in the medical databases.

7.1 Key Findings

Transfer Learning Improves Diagnostic Accuracy: The results of the study showed that the use of transfer learning with pre-trained models (DenseNet, EfficientNet, and ResNet) improved performance by up to 20%, and are especially beneficial in cases with limited amount of training data for the blurred pneumonia images of chest X-rays. Model **Comparison:** Out of the three modeling techniques that were used in our analysis, DenseNet and EfficientNet recorded higher accuracy and better key performance indicators than ResNet. DenseNet had very good performances in terms of computation time while EfficientNet performed very well both in terms of accuracy as well as in terms of its efficiency. Interpretability through Grad-CAM: Grad-CAM integration was beneficial to amend the interpretability issue of the deep learning models, and the clinicians were provided with visual integrities of the models' decision-making process. This is exactly beneficial for the subsequent trust and acceptance in clinical settings. Generalization and Data Variability: While the models were able to do well on the training set, generalization to other different test sets is a major problem. This underlines the need for having several types of datasets and validity checks in order to make the results more reliable and not overfitted.

References

- Aggarwal (2021). Diagnostic accuracy of deep learning in medical imaging: a systematic review and meta-analysis, NPJ digital medicine $\mathbf{4}(1)$: 65.
- Buda, M. and et al. (2018). A systematic study of the class imbalance problem in convolutional neural networks, arXiv preprint arXiv:1808.03124. URL: https://arxiv.org/abs/1808.03124
- Chen (2022). Recent advances and clinical applications of deep learning in medical image analysis, *Medical Image Analysis* **79**: 102444.
- Fu, Y., Lei, Y., Wang, T., Curran, W. J., Liu, T. and Yang, X. (2020). Deep learning in medical image registration: a review, *Physics in Medicine & Biology* 65(20): 20TR01.
- Haskins, G., Kruger, U. and Yan, P. (2020). Deep learning in medical image registration: a survey, *Machine Vision and Applications* **31**(1): 8.
- Liu (2021). A review of deep-learning-based medical image segmentation methods, *Sustainability* **13**(3): 1224.
- Mooney, P. (2018). Chest x-ray images (pneumonia). Available: https://www.kaggle.com/datasets/paultimothymooney/chest-xraypneumonia?resource=download.
- Oquab, M., Bottou, L. and et al. (2014). Learning and transfering mid-level image representations using convolutional neural networks, *CVPR*, pp. 1717–1724.
- Puttagunta, M. and Ravi, S. (2021). Medical image analysis based on deep learning approach, *Multimedia tools and applications* **80**(16): 24365–24398.
- Rajpurkar, P., Irvin, J., Bagul, A. and et al. (2017). Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning, arXiv preprint arXiv:1711.05225 . URL: https://arxiv.org/abs/1711.05225
- Renard (2020). Variability and reproducibility in deep learning for medical image segmentation, *Scientific Reports* **10**(1): 1–9.
- Selvaraju, R. R., Cogswell, M. and et al. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization, arXiv preprint arXiv:1610.02391. URL: https://arxiv.org/abs/1610.02391
- Shorten, C. and Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning, *Journal of Big Data* 6(1): 60.
- Singh (2020). 3d deep learning on medical images: a review, Sensors 20(18): 5097.
- Suganyadevi (2021). A review on deep learning in medical image analysis, *International Journal of Multimedia Information Retrieval* **11**(1): 19–38.
- Wang, X. and et al. (2021). A survey on explainable artificial intelligence for medical image analysis, Computational and Mathematical Methods in Medicine 2021: 5520858.

- Xie, X., Niu, J., Liu, X., Chen, Z., Tang, S. and Yu, S. (2021). A survey on incorporating domain knowledge into deep learning for medical image analysis, *Medical Image Analysis* **69**: 101985.
- Zhou (2021). A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises, *Proceedings of the IEEE* **109**(5): 820–838.