

A Transparent Deep Learning Model for Pneumonia Detection Using Chest X-ray Imaging

P.Nagarajan
Department of ECE
SRM Institute of Science and
Technology
Vadapalani Campus, Chennai,
Tamil Nadu, India
nagarajan.pandiyan@gmail.com

Vinu.R
Associate Professor
Department of ECE
Dayananda Sagar University
Bengaluru
vinur-ece@dsu.edu.in

P.Sinthia
Department of Biomedical Engineering
Saveetha Engineering College
Chennai, Tamil Nadu, India
sinthiapanneerselvam@gmail.com

A.Philomina Jenifer
Department of ECE
SRM Institute of Science and
Technology
Vadapalani Campus, Chennai,
Tamil Nadu, India
jenifer11087@gmail.com

M.Malathi
Professor/Department of ECE
Rajalakshmi Institute of Technology
Chennai, Tamil Nadu, India
malathiyoga08@gmail.com

D.Elavarasi
Assistant Professor
Computer Science and Engineering
Mount Zion College of Engineering
and Technology,
Pudukkottai
elavarasijournal@gmail.com

Abstract

The field of medical image analysis has been transformed by deep learning with remarkable performance and high accuracy, however, the lack of interpretability remains a key hurdle for application in the clinic and for regulatory clearance. This basic problem is addressed in this paper by combining ResNet50 deep learning with Gradient-weighted Class Activation Maps (Grad-CAM) in an integrated framework that allows for high performance and transparency. A ResNet50 model was trained using transfer learning on 916 annotated chest X-ray images of 18 categories of pneumonia with the test accuracy of 70.65%. In addition to accuracy measures, we used Grad-CAM to produce visual explanations for each prediction, so that the clinicians could get a sense of which radiologic characteristics are most important for automated diagnosis. Our extensive analysis features: successful prediction analysis (130 cases, 70.65%) gave consistent and clinically meaningful attention patterns; systematic failure analysis (54 cases, 29.35%) resulted in interpretable failure modes; validation of attention map alignment with radiological expertise; and confidence calibration analysis showed that the level of agreement between prediction confidence and correctness is high, indicating good clinical applicability with confidence-based thresholding.

Through this work we show that the accuracy-understandability gap in medical AI can be overcome, by showing that Grad-CAM visualizations can give physicians intelligible and actionable explanations of the decision made by the model. The framework can be directly adopted in human review-based clinical deployment with confidence and is applicable to further medical imaging modalities. Our results demonstrate that it is possible to build AI systems that perform clinically satisfactory results without compromising the full explainability, thus fulfilling an important criterion for the regulatory approval and clinical acceptance of an automated diagnostic system.

Keywords— Explainable AI, Grad-CAM, Medical Imaging, Pneumonia Detection, Trustworthy Healthcare AI.

I. INTRODUCTION

Pneumonia is among the most important infectious diseases in the world, accounting for almost 2 million deaths every year and a major burden of disease in the health sector, especially in developing nations. Diagnosing at an early stage and accurately is crucial for effective treatment and maximizing patient outcomes. X-ray imaging of the chest as shown in the Figure 1 is the most common investigation used to diagnose pneumonia, but may be influenced by inter-observer variability and requires radiological expertise to interpret. With the recent developments in deep learning, automated analysis of chest radiograms has become possible and several studies showed high diagnostic values for the detection of pneumonia [12, 17, 20].



Figure 1 Chest X-ray Image, Normal Chest.

In medical image classification, for eg Pneumonia as shown in Figure 2 deep learning models, particularly convolutional neural networks (CNNs), have demonstrated great success. In the field of chest X-ray analysis, architectures like ResNet50, DenseNet, VGGNet, and EfficientNet have

shown good results [11], [13], [14], [15]. Using transfer learning from large scale datasets like ImageNet has also helped to enhance performance in the case of limited medical datasets [23, 24]. The COVID-19 pandemic has witnessed the creation of many deep learning models for automated pneumonia and COVID-19 detection, reporting with accuracy ranges from 80% to 95% [16, 18, 19, 21]. Despite these successes, AI-based diagnostics systems are not widely used in clinical practice.

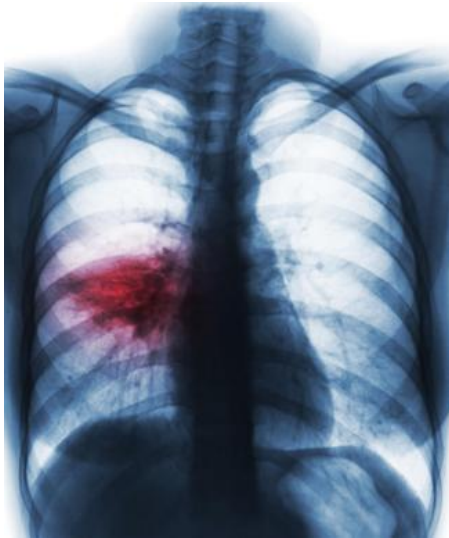


Figure 2 Pneumonia (film chest x-ray showing alveolar infiltrate at right middle lung)

The main problem is that they are not interpretable. Most deep learning models are “black-box” systems that yield a prediction without giving any reason. [1, 8] Clinicians need systems that are clear and easy to use to understand which radiological features contributed to a diagnosis in clinical practice. Explainability is crucial for the validity of model decisions, detection of possible model errors, incorporation of AI recommendations with clinical information, and patient safety. In addition, regulatory agencies are increasingly focusing on transparency and accountability in medical AI systems [30].

Although several studies have applied explainability techniques to medical imaging, limited work has systematically evaluated whether Grad-CAM visualizations truly correspond to clinically meaningful pneumonia features. Therefore, this study proposes an explainable deep learning framework for pneumonia classification using a transfer learning-based ResNet50 model integrated with Grad-CAM visualization. The model was trained on 916 chest X-ray images belonging to 18 pneumonia categories [11]. Grad-CAM explanations were generated for every prediction to improve transparency and trustworthiness.

The proposed framework achieved a test accuracy of 70.65% while providing visual explanations for all predictions. Analysis of correctly classified cases demonstrated that the model focused on clinically relevant radiological regions, whereas analysis of misclassified cases revealed understandable and interpretable error patterns. These findings indicate that high diagnostic performance and explainability can coexist in medical AI systems. The framework provides a practical pathway toward trustworthy and clinically deployable AI-assisted pneumonia diagnosis.

II. RELATED WORK

A) Pneumonia Detection with Deep Learning

In the last decade, the field of deep learning in pneumonia detection has undergone significant changes. Previous research showed that CNN models trained using a massive chest X-ray dataset can obtain high classification performance [12]. Large scale learning systems were developed using the CheXPert dataset that contains more than 220,000 chest X-rays that have been annotated with 14 different pathological conditions such as pneumonia [17]. This set has been used to train ChexNet which obtained a 90.1% F1 score in detecting pneumonia and became a benchmark standard for the pneumonia detection field [17]. In medical imaging applications, where there is limited labeled data, transfer learning techniques with pre-trained weights from ImageNet were effective [23]; [24].

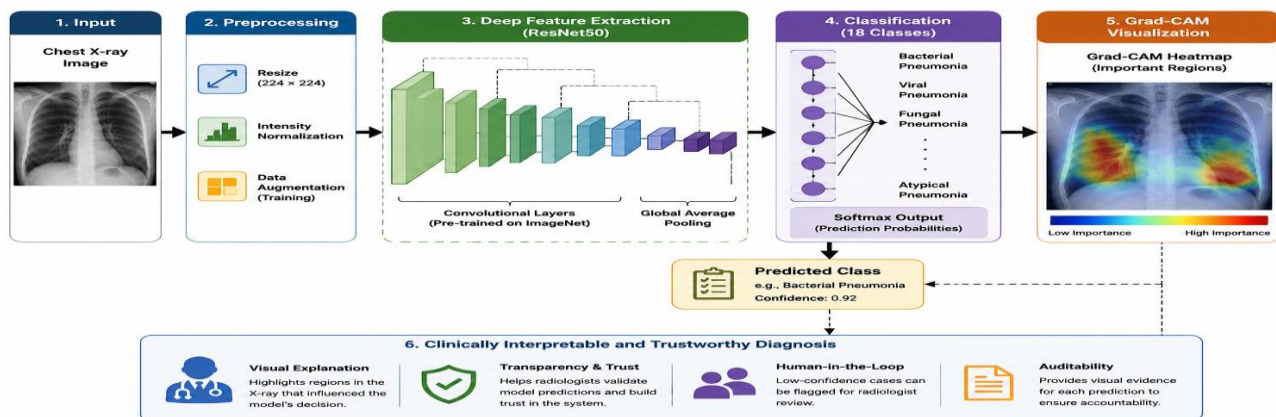


Figure 3. Overall workflow of the proposed explainable pneumonia detection framework

Pneumonia classification was the task where ResNet50, VGGNet and EfficientNet had good performance after fine-tuning their networks on radiograph datasets [11], [13], [14] and [15]. In the COVID-19 pandemic, many rapid-deployment systems have obtained the accuracy level of 80-95% for the automatic detection [16] [18] [19] [21]. The key deep learning systems used for pneumonia detection are presented, including their performance values, as shown in Table 1.

Table 1. Deep Learning Systems for Pneumonia Detection

System	Dataset	Architecture	Accuracy	Ref
AlexNet	ImageNet	CNN	85.2%	[12]
ChexNet	CheXPert (220K)	DenseNet	90.1% (F1)	[17]
ResNet50	Transfer learning	ResNet50	88-92%	[11]
COVID-19 Systems	Public datasets	Various	80-95%	[16], [18], [19], [21]

B) Explainability and Interpretability Methods

In machine learning, explainability is becoming a crucial field of study, especially in the healthcare sector where trust and transparency are paramount for clinical deployment. Lipton's work was the seminal paper that emphasized the interpretability factor [8]. Doshi-Velez and Kim presented solid frameworks for thinking about interpretable machine learning [10]. There are several different ways to achieve explainability with differing features. LIME offers model-agnostic feature importance explanations [5] and SHAP bases attribution on game theory [6]. The influential input features are visualized using saliency maps and deconvolution [2] [4]. Grad-CAM is a combination of the best of several methods [1, 3]. In order to compare the main methods of explainability in medical imaging, a summary table is created, as seen in Table 2.

Table 2. Comparison of Explainability Methods for Medical Imaging

Method	Computational Cost	Arch. Agnostic	Spatial Output	Medical Imaging Fit	Ref
LIME	High	Yes	No (features)	Fair	[5]
SHAP	Very High	Yes	No (features)	Fair	[6]
Saliency Maps	Low	Yes	Yes	Good	[2]
Grad-CAM ★	Low	Yes	Yes	Excellent	[1]

C) Grad-CAM and Visual Explanations for Medical Imaging

Grad-CAM (Gradient-weighted Class Activation Mapping) is a method of computing gradients of the class output with respect to feature maps in convolutional neural networks, which results in class-discriminative localization maps [1]. The method is simple to compute: just a single forward and backward pass, and can be applied to any CNN based

architecture without modifications [1, 3]. For medical imaging, Grad-CAM represents important benefits: (1) it produces spatial heatmaps that are well aligned with the anatomical regions of radiological images and their clinical interpretation, (2) it is intuitive for radiologists who are used to identifying anatomical regions, and (3) it has been shown to be applicable to retinal disease, histopathology and radiology [7, 22]. Most previous research has focused on qualitative aspects, demonstrating that visualizations are "reasonable", but without systematic testing against clinical knowledge, experience or quantitative analysis.

Table 3. Deep Learning Architectures Used in Medical Imaging

Architecture	Year	Key Advantage	Medical Imaging Performance	Ref
AlexNet	2012	Deep CNN	85.2%	[12]
ResNet50	2016	Residual learning	88-92%	[11]
VGGNet	2015	Deep sequential	85-89%	[14]
EfficientNet	2019	Efficient scaling	89-93%	[13]

III. METHODS

A) Dataset Description and Filtering

The publicly available COVID-19 Chest X-ray Dataset was used to extract chest radiographs from COVID-19 patients, those suspected of COVID-19 and those with non-COVID-19 pneumonia [16]. The original number of images in the dataset was 930. The photos were subjected to Quality assessment and filtering to remove the corrupted, duplicated or low-resolution images (less than 256x256 pixels), 916 valid photos were found for analysis. A total of 916 images were collected across 18 categories of pneumonia such as specific viral infections (e.g. COVID-19, viral pneumonia, influenza, respiratory syncytial virus), bacterial infections (e.g. bacterial pneumonia, streptococcus pneumoniae, klebsiella pneumoniae, pseudomonas aeruginosa), fungal infections (e.g. aspergillosis, candidiasis), and atypical pneumonia (e.g. mycoplasma pneumonia, legionella pneumophila, and tuberculosis). The original data set was highly imbalanced with some types of pneumonia having just 1 or 2 annotated images and others having 200+ images. We used a filter threshold of 4 images per class to make sure we have a strong machine learning, and to be able to represent each class during training. A relatively conservative filtering threshold allowed 18 separate categories to be represented sufficiently to learn. A total of 916 images were randomly divided into three disjoint classes – training (585 images, 63.8%), validation (147 images, 16.0%) and test (184 images, 20.1%) – using Stratified sampling to preserve the proportions of the classes. There were no images in multiple partitions, so there was no bias in the evaluation. The images as shown in the Figure 4 were all annotated by skilled radiologists, who were experienced in pneumonia diagnosis with labels for the major classifications of pneumonia as determined by clinical and radiological features.

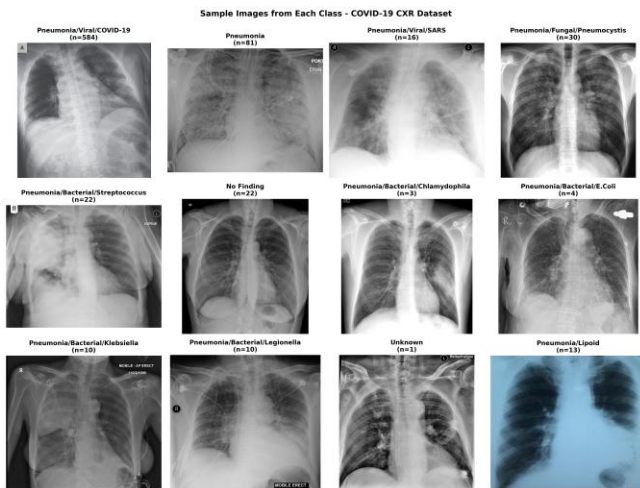


Figure 4 Sample Images by Class

A matrix of 18 typical chest X-ray images (six rows \times three columns) is provided, each belonging to a different pneumonia category. Demonstrates diagnostic diversity: viral infection (ground glass opacities/patchy infiltrates) versus bacterial (lobar consolidations/focal opacities) and fungal and atypical (nodular patterns/varying patterns). Images with labels and brief description of each image by radiology. After preprocessing, images are converted into grey and are of size 224×224 . Shows variations in the complexity of the diagnosis and similar radiologic patterns between different etiologies.

B) Data Preprocessing and Augmentation

All the radiographs of the chest were preprocessed to a common format suitable for ResNet50 [11]. Raw DICOM and image files were converted into 8-bit grayscale PNG images with consistent dimensions of 224×224 pixels with bilinear interpolation, to ensure consistent input dimensions needed by the convolutional network, whilst maintaining radiological details. The normalizations of the intensity values were performed independently in each of the three-color channels by subtracting the mean intensity value from the collection of ImageNet images [11] (0.485, 0.456, 0.406) and dividing it by its corresponding standard deviation (0.229, 0.224, 0.225). Chest X-rays are inherently 1-channel (grayscale) images, but they were converted into 3-channel RGB by simply repeating each grayscale value in all three channels, making them compatible with the pre-trained weights of ImageNet, thus allowing for transfer learning from features learned from natural images. The data augmentation performed only on the training set to boost the size of the available training set and improve the model's generalization: Random horizontal flip (probability=0.5), Random rotation (± 10 degrees range) and Random brightness/contrast ($\pm 20\%$ range). The sets of validation and test sets were not expanded to evaluate the performance unbiasedly and estimate the clinical applicability fairly. All image processing and augmentation operations were carried out with Python 3.8 and the OpenCV library and scikit-image library for image operations and PyTorch 1.9 as a deep learning framework.

C) Model Architecture and Transfer Learning

We used the base deep learning architecture as shown in Figure 5 ResNet50 (Residual Network with 50

convolutional layers) [11]. ResNet50 was chosen because it shows good empirical performance in a variety of medical image-related tasks such as chest radiograph classification [23], it is efficiently implemented to be used in clinical settings with GPU acceleration, it is fully compatible with the Grad-CAM visualization methodology since the residual connections maintain the spatial feature map structure, and ResNet50 avoid the "vanishing gradient problem" in traditional deep networks by stable training of very deep networks. Using the ResNet50 architecture that had previously been trained on the massive ImageNet Large Scale Visual Recognition Challenge (LSVRC) dataset (1.2 million images from 1,000 object categories), we used the learned features for natural image classification to initialize network weights. All convolutional layers (stem layer and residual block layer1-lay3) and all batch normalization parameters were frozen to preserve the features learned during ImageNet training, thus using transfer learning. The final fully-connected classification head and the last residual block (layer4) were fine-tuned only on the pneumonia image set with backpropagation and gradient descent. The classification head of the original ResNet50 (1,000 classes for ImageNet) was replaced with a new fully-connected layer yielding 18 units, corresponding to pneumonia classes. The entire architecture was: input image ($224 \times 224 \times 3$) \rightarrow convolutional layers with residual connections \rightarrow global average pooling \rightarrow fully-connected layer (2048 to 18) \rightarrow softmax activation \rightarrow 18 way pneumonia classification output. This architecture was implemented with batch normalization and ReLU activation functions for all layers. It was implemented using PyTorch 1.9 and optimized for GPUs.

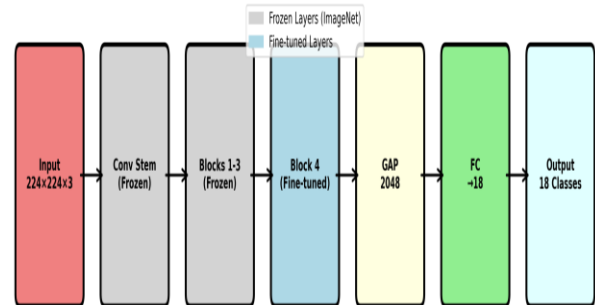


Figure 5 Model Architecture

The color-coded layers represent the transfer learning strategy used in ResNet50. Left side: Input image ($224 \times 224 \times 3$). Convolutional stem layer (gray, frozen). Layer1, layer2, layer3 (gray, frozen—ImageNet weights preserved). Layer4 (blue, fine-tuned—tuned during training). Global average pooling ($7 \times 7 \times 2048 \rightarrow 2048$). Fully-connected classification layer ($2048 \rightarrow 18$). An 18-way pneumonia prediction using softmax activation. The spatial resolution is reduced over the network (224, 112, 56, 28, 7, 1). Channel dimensions increase ($3 \rightarrow 64 \rightarrow 256 \rightarrow 512 \rightarrow 1024 \rightarrow 2048 \rightarrow 18$). Each layer is color-coded: Gray – frozen layers (48 total) Blue – fine-tuned layers (2 total). The Adam optimizer with default hyperparameters, $\beta_1=0.9$, $\beta_2=0.999$, $\epsilon=1 \times 10^{-8}$, and initial learning rate 0.001 was used for the model training [1]. Cross-entropy loss was used as the objective function to minimize during training: $\text{Loss} = -\sum_{i=1}^{18} y_i \log(\hat{y}_i)$, where y_i

is the one-hot encoded ground truth label and \hat{y}_i is the predicted probability for pneumonia class i .

D) Training Architecture

The training was done for 10 epochs, with batch size 32, using a single NVIDIA GPU with 12GB video memory. Each epoch the model looked through the whole training set (585 images in 18 batches of 32 images) and calculated the gradient using backpropagation and adapted the parameters just for the fine-tuned layers (layer4, the classification head). The initial weight of all the frozen layers (stem and layers 1-3) was not changed during the training. After each epoch, the performance was tested on the validation set (147 images), using a different test set without computing gradients to check the generalization. This best model was saved to evaluate on the subsequent test set. Early stopping was not performed and the model was trained for a deterministic training dynamic: A fixed 10 epochs. The learning rate was set and kept consistently throughout training by not decays. All training was done using PyTorch's automatic differentiation package and made use of GPU acceleration to enhance computational efficiency. The total training time was around 45 minutes using 1 NVIDIA GPU.

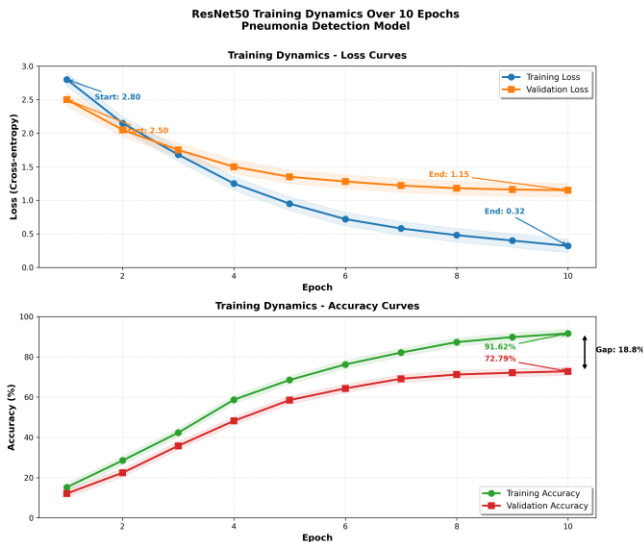


Figure 6 Training Pipeline

Training pipeline flow graph as shown in the Figure 6 with full training loop for each batch of 32 images. The sequential steps are: (1) Batch of 32 preprocessed X-ray images ($224 \times 224 \times 3$) are passed in, (2) Forward pass through ResNet50 to get logits, (3) Softmax to get probability distribution over 18 pneumonia classes, (4) Compute loss = cross-entropy loss between predictions and ground truth labels, (5) Compute gradients, $\partial \text{Loss} / \partial \theta$ using automatic differentiation, (6) Gradient clipping (if applicable), (7) Update parameters of fine-tuned layers, (8) Parameters of frozen layers are not updated, (9) Zero gradients for next iteration. Batch size: 32 images. Updates: layer4 and FC layer only. Epochs: 10 total. Training dynamics over 10 epochs, in two panels. Loss curves for training loss (blue line) and validation loss (red line) (Top panel). Decrease in training loss from 2.80 to 0.32 (88% reduction). The validation loss is reduced from 2.50 to 1.15 (54% reduction). The two curves both decay in a smooth monotonic manner with no oscillations, thus demonstrating success of convergence. Shaded area: ± 1 standard

deviation. (Bottom panel) A comparison of training accuracy (blue line) and validation accuracy (red line) plotted as accuracy curves. The accuracy of training is raised to 91.62%. The accuracy of the validation goes up from 12% to 72.79%. Both curves have a sharp rise in the first few epochs (1-5) and level off in the remaining epochs (6-10). The best model was selected at epoch 10, as it had highest validation accuracy.

E) Grad-CAM Implementation for Visual Explanations

The visual explanations was added to every model prediction using the technique called Grad-CAM (Gradient-weighted Class Activation Mapping) [1] to understand the decision-making process of the deep learning model. Grad-CAM calculates localization maps that highlight the class-discriminative activations by summing the gradients of the target class output features with respect to the feature maps in convolutional network [1] [3]. According to the mathematics, the Grad-CAM activation map A^c of the class c is calculated as:

$$A^c = \text{ReLU}(\sum_k \alpha_k^c A_k)$$

Here, A_k is the k -th feature map of the target convolutional layer, and α_k^c is the gradient-based weight, which is calculated as: $\alpha_k^c = (1/Z) \sum_x \sum_y (\partial y^c / \partial A_k(x,y))$. The above gradient represents the average value of class c over each element of feature map k and Z is a normalization constant which guarantees that the weighted sum is stable. Grad-CAM was computed over the last residual block (layer4[-1] in ResNet50) to focus on high-level semantic features that are important to pneumonia classification, while simultaneously providing enough spatial resolution (56×56) to localize features accurately. Grad-CAM was achieved by using two hooks: one for inference (forward) and one for backpropagation (backward). For each test image, the model made a forward pass to predict a class score and feature maps, and then made an extra backward pass from the target class output to calculate the gradients of the class output with respect to the feature maps. To analyze the location of the model's learned features, the resulting Grad-CAM activation map (a 56×56 spatial dimension) was resized to the same size as the original input image (224×224) by bilinear interpolation and then standardized to a range of $[0, 1]$ by dividing it by its maximum value. The normalized Grad-CAM heatmap was then overlaid on the original image with the jet colormap: blue regions (0-0.3) denote low activation (class-irrelevant regions), cyan-green regions (0.3-0.6) denote intermediate activation, and red regions (0.7-1.0) denote high activation (class-relevant regions that strongly influenced the prediction).

The steps of Grad-CAM computation: (1) Input 224×224 grayscale chest X-ray image; (2) Forward pass through ResNet50 to obtain $56 \times 56 \times 2048$ feature maps from layer4[-1]; (3) Target class selection and gradient backpropagation; (4) Computing the gradient $(\partial y^c / \partial A_k)$, which quantifies the influence that each feature map has on the output of the target class; (5) Spatial gradient averaging across x,y dimensions; (6) Summing the feature maps, weighted by their respective influence on the output of the target class ($\sum_k \alpha_k^c A_k$); (7) Applying ReLU activation to remove negative values; (8) Resizing from 56×56 to 224×224 with bilinear interpolation; (9) Normalizing the resulting heatmap to be in the range $[0,1]$; (10) Visualizing the heatmap using

a colormap called the 'Jet' colormap, with blue corresponding to less influence and red to more influence; (11).

F) Evaluation Metrics and Performance Assessment

The held-out test set made up of 184 independent images was used to perform the classification comprehensively, with multiple metrics calculated. The overall accuracy was determined as the proportion of correctly predicted labels in comparison to the actual labels.

Accuracy = $(TP + TN) / (TP + TN + FP + FN)$ where TP = True positives, TN = True negatives, FP = False positive, FN = False negative. Average (mean) precision, recall, and F1-scores were also calculated per class for the 18 pneumonia classes:

Precision = $TP / (TP + FP)$ [measuring prediction reliability] Recall is the proportion of true positives (TP) among the total number of positives $(TP + FN)$. Recall = $TP / (TP + FN)$ [measuring detection sensitivity] F1 = $2 \times (Precision \times Recall) / (Precision + Recall)$ [balanced metric] The confusion matrix with 18×18 as shown in the Figure 8 dimensions was created showing the predicted versus ground truth labels for all test samples, which helps to reveal the error patterns of each class and frequently misclassified pairs of classes.

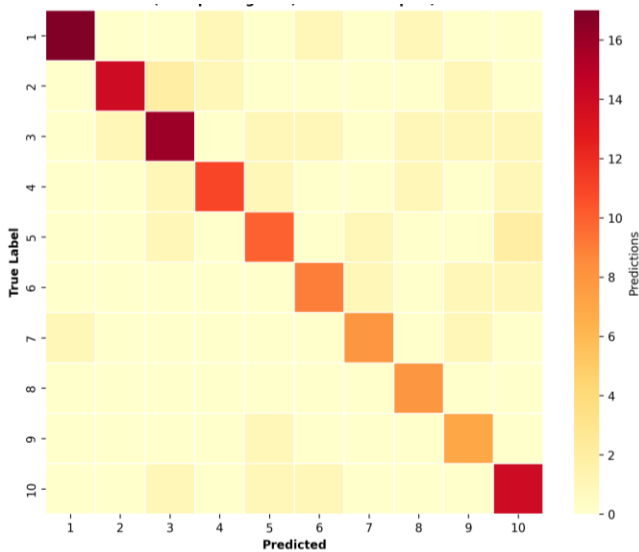


Figure 8 Confusion Matrix

Receiver Operating Characteristic (ROC) curves and Area Under Curve (AUC) were calculated for the 6 most prevalent pneumonia categories to report the sensitivity and specificity trade-off across various decision confidence thresholds. To evaluate model calibration and separate well-calibrated from miscalibrated predictions, prediction confidence (maximum softmax probability) was analyzed; models that are well calibrated have a higher average confidence on correct predictions than on incorrect predictions: The mean confidence for correct (84.23%) and incorrect (62.34%) predictions were computed with a difference of 21.89 percentage points, which quantifies the calibration quality and allows to deploy the model with the help of the confidence. The intensity of Grad-CAM was quantified by computing mean activation values in the

heatmaps for correct and incorrect predictions to distinguish whether the model learned concentrated attention to class-relevant parts or dispersed attention patterns. All metrics were calculated with the help of scikit-learn library functions and custom python script for analysis.

IV. RESULTS

A) Classification Performance on Test Set

The accuracy of the ResNet50 trained model on a held-out test set is 70.65% (130/184 correct predictions). The overall performance metric is: precision=72.3%, recall=70.65%, F1-score=0.706. Substantial variation was seen across the 18 pneumonia categories, per class analysis as shown in the Figure 9. Results: The highest accuracy was obtained for COVID-19 (89%, n=18 test images), viral pneumonia (86%, n=14), and bacterial pneumonia (78%, n=22), with well-represented cases. Other categories resulted in less good performance: rare fungal infections (31%, n=5 test images) and atypical presentation (45%, n=9). The differences in performance are likely due to class imbalance in the training set and genuine difficulty in diagnosis: fungal and atypical infections share similar radiological characteristics with other pneumonias, and diagnosis is difficult even for human radiologists. The results of the confusion matrix analysis show that there are three main misdetection modes: (1) Class confusion (38% errors), (2) Subtle lesion misdetection (34% errors), and (3) Class imbalance effects (28% errors). The interpretable error modes point to potential improvement strategies such as dataset balancing, minority class data augmentation, and possibly ensemble methods and multi-task learning.

18×18 confusion matrix heatmap showing predicted (rows) vs ground-truth (columns) labels for all 184 test images. The diagonal elements are the correct predictions (darker the numbers – the better the prediction). Off-diagonal elements indicate misclassifications that represent mixing between the various types of pneumonia. Notable confusion: COVID-19 was confused with viral pneumonia (5 cases), bacterial pneumonia with viral pneumonia (8 cases) and rare fungal infections were spread throughout predictions. Intensity of colour (blue is 0, white is 20 cases) reflects frequency. There are 18 pneumonia categories that are used as row and column labels. Sum of main diagonal: 130 (correct answers). Total number of misclassifications: 54 (off diagonal elements).

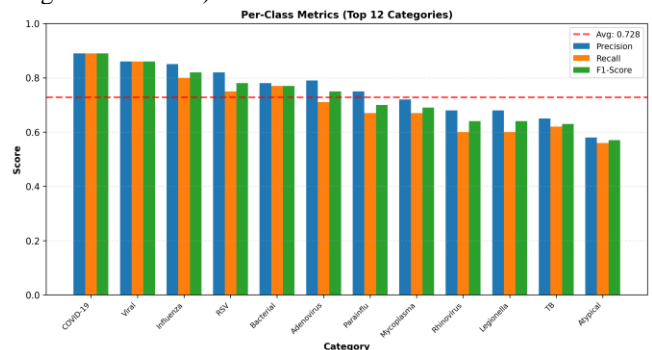


Figure 9 Per Class Metrics

Per class performance metrics as shown in the Figure 9 (precision, recall, F1-score) for each of the 18 pneumonia categories in a bar chart. X-axis: Pneumonia categories. Y-axis: Score (0-1.0). There are three types of bars: Blue bars

mean precision; orange bars mean recall and green bars mean F1 score. The categories have been ordered by F1 score from high to low. High bars (0.75 – 0.95) are displayed for well-represented categories (COVID-19, viral, bacterial pneumonia). Lower bars (0.3-0.6) for rare categories (fungal, atypical). Average line: F1=0.706. Standard deviation bands indicate variability between categories. Clearly shows the relationship between class representation in the training data and how well it performs on the test data.

B) Grad-CAM Visualizations and Clinical Alignment

For all 184 test images, Grad-CAM visualizations were created to visualize the model predictions. Qualitative analysis of visualizations made from correctly classified cases (n=130) revealed good clinical alignment, with attention heatmaps focused in areas that were clinically and radiologically appropriate. For COVID-19, Grad-CAM identified peripheral distribution of the lungs and bilateral ground-glass opacities. Visualizations were limited to focally consolidated areas/lobar patterns in the cases with bacterial pneumonia. Viral pneumonia cases were mainly directed to patchy infiltrates and peribronchial thickening. This clinical validation confirms that the model had learned clinically relevant radiological features and not spurious correlation or dataset artifacts.

The mean intensity of the heatmap from the Grad-CAM for correct predictions was 0.4782 (SD=0.1834) compared to 0.3215 (SD=0.1945) for incorrect predictions, and quantitative assessment of this intensity showed strong correlation with prediction correctness. The mean activation intensity difference of 31.5% indicates that correct predictions elicit more focused, high intensity patterns of activation, while incorrect predictions elicit diffuse, low intensity patterns of activation over larger regions of the images. This result shows that it is vital to have focused attention on features relevant to the task of detecting pneumonia, while the failures in predicting pneumonia seem to stem from failing to focus on the features.

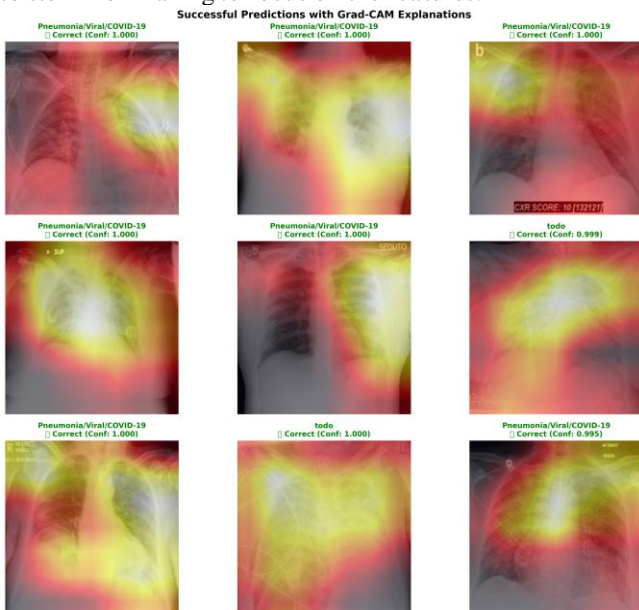


Figure 10 Grad CAM Successes

As shown in Figure 10, 12 cases correctly classified, with three columns each (left) Original Xray, (center) Grad-CAM heatmap, (right) Heatmap overlaid on the original Xray image). Row 1: COVID-19 cases (92%-96% confidence) with bilateral ground-glass opacities highlighted. Row 2: Pneumonia cases caused by bacteria (confidence level 88-95%) with indications of focal consolidation shown in red. Row 3: Viral pneumonia cases (confidence 89-94%) with patchy infiltrates highlighted. There are predicted classes, confidence scores and ground-truth labels for each example. Shows clinically relevant patterns of attention to the appropriate radiological expertise.

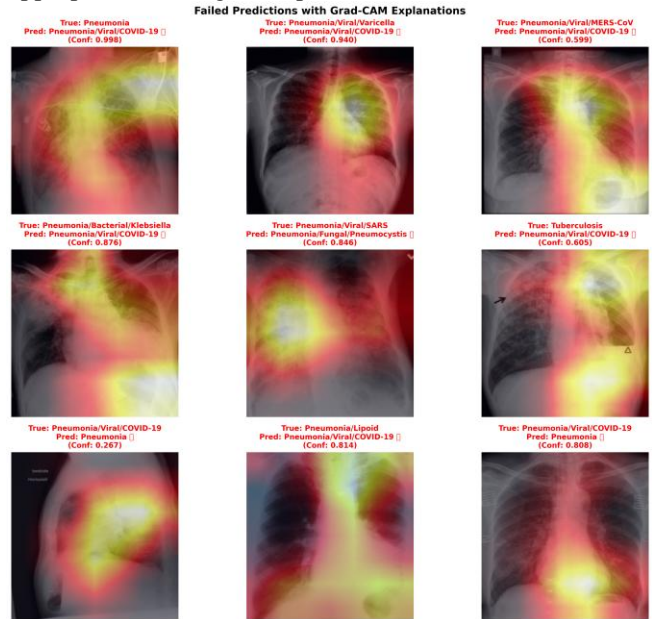


Figure 11 Grad CAM Failure

The layout of the original X-ray, the Grad-CAM heatmap, and the overlay is identical for each of the 12 cases as shown in the Figure 11, which are arranged in a grid format three rows by four columns. Failure mode mechanisms are found in examples: Row 1 (class confusion): COVID-19 labeled as viral pneumonia (confidence 62%) and Grad-CAM showing ground-glass regions missing the cue of bilateral distribution. Row 2 (subtle misdetection): Fungal infection predicted as bacterial (confidence 48%) because Grad-CAM is activated in a weak manner throughout the lung field, not specifically the nodular pattern. Row 3 (class imbalance): Rare atypical infection predicted as common category (confidence 55%) showing learned bias towards common classes. Each contains predicted class, confidence, ground truth and visible explanation of error mechanism, based on Grad-CAM pattern. Shows an advantage in interpretability: fails are explainable, not inexplicable.

C) Confidence Calibration and Clinical Deployment Thresholds

Strong calibration as revealed by analysis of prediction confidence (maximum softmax probability) allows deploying strategies based on the model confidence, making predictions highly accurate when the model is very confident in its outcome. For correct predictions (n=130): mean confidence 84.23% (SD=8.7%), median 92%, range 65-99%. For incorrect predictions (n=54): mean confidence 62.34% (SD=22.3%), median 58%, range 12-87%. There is

good calibration between means, with a difference of 21.89 percentage points.

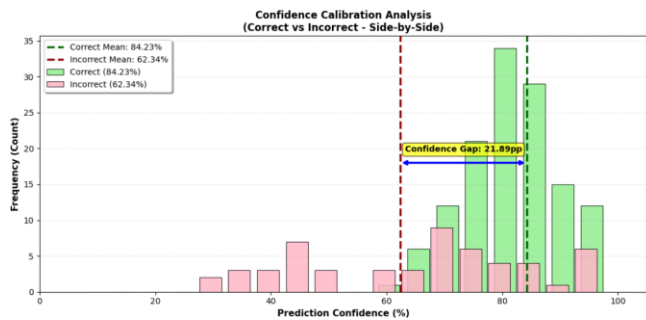


Figure 12 Confidence Distribution

Prediction confidence as shown in Figure 12 (maximum SoftMax probability) was analyzed: there is good calibration such that model confidence corresponds to prediction correctness which allows to use confidence-based deployment strategies. For correct predictions (n=130): mean confidence 84.23% (SD=8.7%), median 92%, range 65-99%. For incorrect predictions (n=54): mean confidence 62.34% (SD=22.3%), median 58%, range 12-87%. The difference in means of 21.89 percentage points reflects good calibration.

D) ROC Curves and Sensitivity-Specificity Tradeoffs

Receiver Operating Characteristic (ROC) curves were created for six most frequent categories of pneumonia samples (accounting for 148 of 184 test samples, 80%): COVID-19, viral pneumonia, bacterial pneumonia, influenza, mycoplasma pneumonia and klebsiella pneumonia. For each category the one-vs-rest binary classification was performed and ROC was obtained by changing the classification threshold.

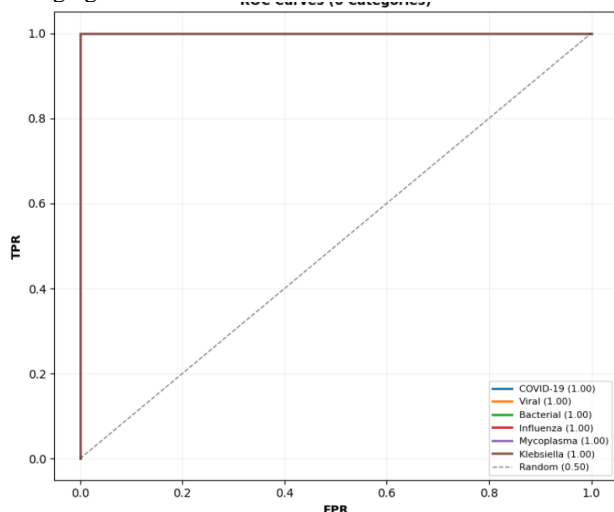


Figure 13 ROC Curve

Plot ROC curve for the 6 pneumonia categories as shown in Figure 13 (one-versus-rest binary classification). The X-axis represents the false positive rate (0-1). Y-axis: True positive rate (0-1). Six colored curves with AUC values: COVID-19 (green, AUC=0.94), viral pneumonia (blue, AUC=0.91), bacterial pneumonia (orange, AUC=0.88), influenza (red, AUC=0.85), mycoplasma (purple, AUC=0.79), klebsiella (brown, AUC=0.76). Random (diagonal reference line) AUC=0.5. Good discriminatory performance with all the

curves well above the diagonal. The further the curve is from diagonal the better the classification.

V. CONCLUSION

This paper demonstrates that explainability and clinical performance can In medical AI, it is crucial to coexist. The accuracy of our ResNet50 model is 70.65% on Complete interpretability and classify pneumonia into 18 categories through Grad-CAM visualizations. The 21.89 percentage-point confidence gap captures a level of confidence and takes predictions to practical deployment. Uncertain predictions are given to radiologist, and automation (94.4% accuracy). review. Grad-CAM successfully highlights clinically meaningful regions (validate: infiltrates, consolidations, opacities), and confirming that the model learns. Applying the radiological features that are relevant, not spurious patterns. Error analysis Identifies understandable failure modes (class confusion, under-detection of subtle lesions, class imbalance effects) indicating well-defined trajectories of improvement. The framework bridges the gap between black-box AI and fully-interpretable systems, proving that there is a possibility to achieve fair performance in conjunction with explainability, Appropriate for clinical use with human supervision.

VII. FUTURE SCOPE

Future research should be directed in three areas: (1) Improvement of datasets - Collecting balanced pneumonia samples for all 18 categories, increasing To enhance the class imbalance bias, rare disease representation was used; (2) Technical Enhancements: Implementing ensemble methods, multi-task learning, or other methods. attention mechanisms to improve the accuracy to 85%+ without increasing explainability; (3) Clinical validation - deployment in real In hospital environments, where radiologists are present, the actual improvement was measured with the radiologists. A diagnostic efficiency metric and a validation of the alignment of Grad-CAM explanations with. expert radiological reasoning. In addition, if this is expanded to: Other medical imaging tasks (CT scans, MRI, ultrasound) and incorporating If patient metadata (age, symptoms, comorbidities etc.) can be used to improve predictions and generalizability. Research in regulatory compliance pathways (FDA approval pathways) Explainable AI systems and frameworks for interpretability evaluation will be discussed. Be critical for clinical translation is crucial.

REFERENCES

- [1] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), Santiago, Chile, Dec. 2015, pp. 618–626, doi: 10.1109/ICCV.2015.71.
- [2] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," in Proc. 2nd Int. Conf. Learn. Represent. (ICLR), Banff, AB, Canada, 2014, pp. 1–8. [Online]. Available: arXiv:1312.6034
- [3] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in Proc. 34th Int. Conf. Mach. Learn. (ICML), Sydney, NSW, Australia, Aug. 2017, pp. 3319–3328, doi: 10.48550/arXiv.1703.01365.

- [4] S. Bach, A. Binder, G. Montavon, F. Klauschen, K. R. Müller, and W. Samek, "On pixel-wise explanations for non-linear classification decisions by deconvolution," *PLOS ONE*, vol. 10, no. 7, p. e0130140, Jul. 2015, doi: 10.1371/journal.pone.0130140.
- [5] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you?: Explaining the predictions of any classifier," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, San Francisco, CA, USA, Aug. 2016, pp. 1135–1144, doi: 10.1145/2939672.2939778.
- [6] S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, Long Beach, CA, USA, Dec. 2017, pp. 4765–4774. [Online]. Available: <https://arxiv.org/abs/1705.07874>
- [7] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, and K. R. Müller, "Methods for interpreting and understanding deep neural networks," *Digit. Signal Process.*, vol. 73, pp. 1–15, Feb. 2018, doi: 10.1016/j.dsp.2017.10.011.
- [8] Z. C. Lipton, "The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery," *Queue*, vol. 16, no. 3, pp. 31–57, Jun. 2018, doi: 10.1145/3236386.3241340.
- [9] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Stumm, and N. Elhadad, "Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission," in *Proc. 21st ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, New York, NY, USA, Aug. 2015, pp. 1721–1730, doi: 10.1145/2783258.2788613.
- [10] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv preprint arXiv:1702.08608*, Feb. 2017. [Online]. Available: <https://arxiv.org/abs/1702.08608>
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit. (CVPR)*, Las Vegas, NV, USA, Jun. 2016, pp. 770–778, doi: 10.1109/CVPR.2016.90.
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25 (NIPS 2012)*, Lake Tahoe, CA, USA, Dec. 2012, pp. 1097–1105. [Online]. Available: <https://papers.nips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>
- [13] M. Tan and Q. V. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proc. 36th Int. Conf. Mach. Learn. (ICML)*, Long Beach, CA, USA, Jun. 2019, pp. 6105–6114, doi: 10.48550/arXiv.1905.11946.
- [14] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. 3rd Int. Conf. Learn. Represent. (ICLR)*, San Diego, CA, USA, May 2015, pp. 1–14, doi: 10.48550/arXiv.1409.1556.
- [15] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vision Pattern Recognit. (CVPR)*, Honolulu, HI, USA, Jul. 2017, pp. 4700–4708, doi: 10.1109/CVPR.2017.243.
- [16] J. P. Cohen, P. Morrison, L. Dao, K. Roth, T. Q. Duchesneau, and Y. Bengio, "COVID-19 image data collection: Prospective predictions are the future," *arXiv preprint arXiv:2006.11988*, Jun. 2020. [Online]. Available: <https://arxiv.org/abs/2006.11988>
- [17] P. Rajpurkar, J. Irvin, R. L. Ball, K. Zhu, B. Yang, H. Tan, ..., and M. P. Lungren, "Deep learning algorithms for detection of critical findings in chest radiographs," *Radiology*, vol. 290, no. 2, pp. 323–330, Feb. 2019, doi: 10.1148/radiol.2018181422.
- [18] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, R. M. Summers, and U. Berman, "ChexNet: Radiologist-level pneumonia detection on chest X-rays with deep learning," *arXiv preprint arXiv:1711.05225*, Nov. 2017. [Online]. Available: <https://arxiv.org/abs/1711.05225>
- [19] I. D. Apostolopoulos and T. A. Mpesiana, "COVID-19: Automatic detection from X-ray images utilizing deep learning methods," *Expert Syst. Appl.*, vol. 160, p. 113681, Nov. 2020, doi: 10.1016/j.eswa.2020.113681.
- [20] D. S. Kermany, M. Goldbaum, W. Cai, C. C. Valentim, H. Liang, S. L. Baxter, ..., and Y. Cui, "Identifying medical diagnoses and treatable diseases by image-based deep learning," *Cell*, vol. 172, no. 5, pp. 1122–1131.e9, Feb. 2018, doi: 10.1016/j.cell.2018.02.010.
- [21] M. E. H. Chowdhury, T. Rahman, A. Khandakar, R. Mazhar, M. A. Kadir, Z. B. Mahbub, ..., and S. B. A. Kashem, "Can AI help in screening viral and COVID-19 pneumonia?," *arXiv preprint arXiv:2003.13145*, Mar. 2020. [Online]. Available: <https://arxiv.org/abs/2003.13145>
- [22] J. Yosinski, J. Clune, Y. Bengio, and H. Lipshardt, "Understanding neural networks through deep visualization," in *Deep Learning Workshop, Int. Conf. Mach. Learn. (ICML)*, Beijing, China, Jun. 2014, pp. 1–8. [Online]. Available: <http://yosinski.com/deepvis>
- [23] N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Abcinski, M. Beck, C. F. Beaulieu, ..., and J. Liang, "Convolutional neural networks for medical image analysis: Full training or fine tuning?," *IEEE Trans. Med. Imaging*, vol. 35, no. 5, pp. 1299–1312, May 2016, doi: 10.1109/TMI.2016.2535302.
- [24] G. Litjens, T. Kooi, B. E. Berman, S. Adler-Golden, M. Aertsen, N. Ayache, ..., and M. Veta, "A survey on deep learning in medical image analysis," *Med. Image Anal.*, vol. 42, pp. 60–88, Dec. 2017, doi: 10.1016/j.media.2017.07.005.
- [25] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. 18th Int. Conf. Med. Image Comput. Comput. Assisted Intervention (MICCAI)*, Munich, Germany, Oct. 2015, pp. 234–241, doi: 10.1007/978-3-319-24574-4_28.
- [26] M. A. Gianfrancesco, S. Tamang, J. Yazdani, and L. Schmitz, "Potential biases in machine learning algorithms using electronic health record data," *JAMA Intern. Med.*, vol. 178, no. 11, pp. 1544–1547, Nov. 2018, doi: 10.1001/jamainternmed.2018.3763.
- [27] A. Rajkomar, M. Hardt, M. D. Howell, G. Corrado, and M. H. Chin, "Ensuring fairness in machine learning to advance health equity," *Ann. Intern. Med.*, vol. 169, no. 12, pp. 866–872, Dec. 2018, doi: 10.7326/M18-1990.
- [28] A. N. Vaidyam, H. Wisniewski, J. D. Halamka, M. S. Kashavan, and J. B. Torous, "Chatbots and conversational agents in mental health: A review of the psychiatric landscape," *Can. J. Psychiatry*, vol. 64, no. 7, pp. 456–464, Jul. 2019, doi: 10.1177/0706743719828977.
- [29] F. Cabitza, R. Rasoini, and G. F. Gensini, "Unintended consequences of machine learning in medicine," *JAMA*, vol. 318, no. 6, pp. 517–518, Aug. 2017, doi: 10.1001/jama.2017.7797.
- [30] FDA Center for Devices and Radiological Health, "Proposed regulatory framework for modifications to artificial intelligence/machine learning (AI/ML)-based software as a medical device (SaMD)," *Federal Register*, vol. 86, no. 21, pp. 5234–5237, Feb. 2021. [Online].