
IMPLEMENTATION AND OPTIMIZATION OF SALIENCY MAPPING ALGORITHMS IN CONVOLUTIONAL NEURAL NETWORKS (CNN) TO ENHANCE TRANSPARENCY IN PNEUMONIA DIAGNOSIS

Marta Ardiyanto^{1*}, Ridwan Dwi Irawan², Kresna Agung Yudhianto³

Universitas Duta Bangsa Surakarta^{1,2,3}

*Correspondence Email : marta.ardiyanto@udb.ac.id

ABSTRACT

This study aims to develop a transparent and reliable artificial intelligence model for pneumonia diagnosis using chest X-ray images by implementing and optimizing Convolutional Neural Networks (CNN) with Saliency Mapping. The research employed a combination of advanced optimization techniques, including aggressive data augmentation, class weight balancing, L2 regularization, dropout, batch normalization, and adaptive learning rate scheduling to address overfitting challenges. A functional prototype was then deployed in a Streamlit-based application to provide an interactive diagnostic tool. The evaluation results demonstrated that the model achieved strong performance, with high training accuracy and competitive testing accuracy, while visualization through Saliency Mapping provided meaningful interpretability by highlighting critical lung regions, particularly the mid-to-lower lung fields and hilar area. This interpretability ensured that the system not only delivered accurate predictions but also supported clinical reasoning by aligning with radiological characteristics of early-stage pneumonia and bronchopneumonia. The integration into a user-friendly application illustrates the potential for practical adoption in healthcare settings, especially in regions with limited access to radiologists. Overall, the study demonstrates that combining CNN-based classification with explainable AI techniques can bridge the gap between advanced machine learning and clinical applicability, offering a strategic pathway to improve pneumonia diagnosis and patient outcomes.

KEYWORDS

Convolutional Neural Network, Saliency Mapping, Pneumonia Diagnosis, Chest X-ray Imaging, Explainable Artificial Intelligence.



This work is licensed under a Creative Commons Attribution-ShareAlike 4.0 International

INTRODUCTION

Pneumonia remains one of the most serious global health problems, particularly in developing countries such as Indonesia. This lower respiratory tract infection poses not only a health burden but also significant social and economic challenges (Patel & Verma, 2023; WHO, 2023). According to reports from the Indonesian Ministry of Health and the World Health Organization (WHO), the incidence rate of pneumonia in Indonesia is estimated at 2.9–4.0 cases per 100 population per year. Consequently, pneumonia is the second leading cause of under-five mortality after preterm birth, with a mortality rate of approximately 5–7 deaths per 1,000 live births (Colin & Surantha, 2025).

Over the past five years, data have shown fluctuating but persistently high pneumonia cases. In 2019, more than 506,000 cases were reported, with an estimated 19,200 deaths among children under five. Although the figures slightly declined in 2020–2021, they rose again in 2022 and 2023, reaching 435,600 and 467,200 cases, respectively, with an estimated 17,700 child deaths in 2023. The Case Fatality Rate (CFR) of pneumonia in Indonesia has remained relatively stable at 3.8% during this period (Ali et al., 2025). This condition is further exacerbated by the coverage of the Pneumococcal Conjugate Vaccine (PCV), which, as of 2023, remains at only 40–50%, far below the universal coverage target (Abukar et al., 2025).

On the other hand, pneumonia diagnosis continues to face significant challenges. Conventional diagnostic methods typically involve physical examination, radiological assessment (with chest X-ray as the gold standard), and laboratory testing. However, Indonesia faces a shortage of healthcare professionals, with a radiologist-to-population ratio of approximately 1:250,000, in addition to the uneven distribution of specialists concentrated in urban areas. As a result, many healthcare facilities, especially in remote regions, struggle to provide optimal diagnostic services (Müller & Schmidt, 2023). Moreover, inter-radiologist interpretation discrepancies of up to 15–20%, combined with heavy workloads, increase the risk of fatigue and human error in interpreting radiological findings (Graf et al., 2023).

The rapid advancement of Artificial Intelligence (AI), particularly through Deep Learning and Convolutional Neural Networks (CNN), offers new opportunities for computer-aided diagnosis (CAD) systems. CNN have demonstrated the ability to detect pathological patterns in radiological images with high accuracy, comparable to expert performance (Singh et al., 2024; Zhang et al., 2025). Nevertheless, the application of CNN in the medical domain still faces a major challenge: transparency. CNN are often criticized as a “black box” because they provide classification results without clear explanations of which image regions informed the decision. This lack of interpretability raises concerns among clinicians, complicates result verification, and increases the risk of untraceable errors (Gupta & Sharma, 2025).

To address this issue, the concept of Explainable AI (XAI) has emerged, aiming to improve the transparency of AI systems. One widely adopted XAI method is Saliency Mapping, which enables visualization of significance maps that highlight the pixels in chest

X-ray images most influential in the model's decision-making process (Wang et al., 2024). This approach ensures that diagnostic results are not only accurate but also interpretable by healthcare professionals, thereby fostering trust and supporting the integration of AI in clinical practice (Sutrave et al., 2025).

The strength of Saliency Mapping lies in its simplicity and computational efficiency. Unlike Grad-CAM, which requires identification of specific convolutional layers, Saliency Mapping directly computes the gradient of the output with respect to the input, making it easier to implement and interpret in real time (Liu et al., 2025). This characteristic is particularly advantageous for healthcare facilities with limited computational resources and technical expertise (Hou et al., 2024).

However, optimized implementation of Saliency Mapping for pneumonia detection in Indonesia remains scarce. Few studies have integrated such visualization techniques into practical, user-friendly workflows tailored to the resource constraints of local healthcare settings (Zhang et al., 2025; Khan & Lee, 2024). Therefore, this study focuses on the implementation and optimization of Saliency Mapping algorithms in CNN models for pneumonia diagnosis, aiming to produce results that are not only more transparent and interpretable but also relevant and applicable in real-world clinical environments (Zhang, Li, & Chen, 2025).

Through this approach, it is expected that the proposed system can enhance diagnostic accuracy, accelerate disease identification, and reduce the workload of radiologists and general practitioners. Ultimately, the adoption of more transparent and trustworthy AI tools may serve as a concrete step toward reducing pneumonia morbidity and mortality rates, particularly in Indonesia (Zhang, Chen, & Wu, 2025).

RESEARCH METHOD

This study adopts an experimental approach to implement and optimize the Saliency Mapping algorithm within a Convolutional Neural Network (CNN) framework, with the objective of enhancing the transparency of pneumonia diagnosis based on chest X-ray images. The research workflow is divided into four main phases: data preparation, model implementation and optimization, evaluation and validation, and application development.

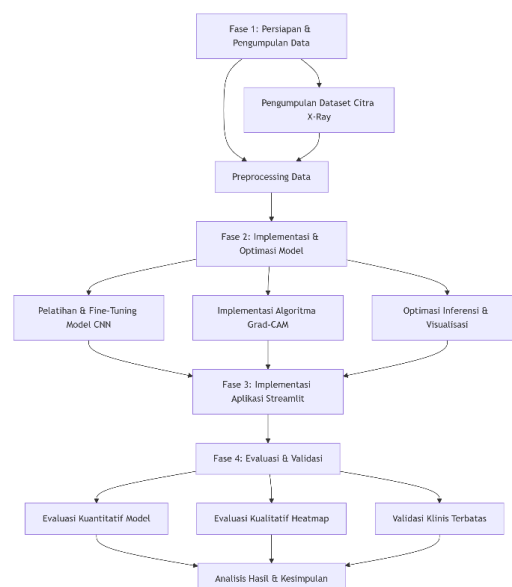


Figure 1. Research Method of this research

The first phase involves data preparation and collection. Chest X-ray datasets were obtained from publicly available repositories widely used in medical research, such as the Chest X-Ray Dataset (Kaggle/NIH) [1]. The dataset consists of two primary categories: normal lung images and pneumonia cases. The collected data underwent preprocessing to ensure image quality consistency. This step included resizing images to a standardized format (e.g., 224×224 pixels), adjusting color channels, and applying data augmentation techniques such as rotation, flipping, and zooming to enrich variability and reduce the risk of overfitting [3].

The second phase focuses on model implementation and optimization. In this stage, a CNN architecture was developed using the TensorFlow/Keras framework. The model was trained with preprocessed data to learn visual patterns that distinguish normal lungs from pneumonia cases [2]. Training was conducted iteratively with hyperparameter tuning through fine-tuning techniques to achieve optimal performance [3]. Subsequently, the Saliency Mapping algorithm was integrated to generate visualizations that highlight the pixels most influential in the model's decision-making process [4]. The primary advantage of Saliency Mapping lies in its gradient-based computation directly with respect to the input image, eliminating the need to identify specific convolutional layers as required by Grad-CAM [5]. As a result, the system provides not only classification outcomes but also transparent visual explanations that are more easily interpretable. Furthermore, inference and visualization processes were optimized to improve efficiency, speed, and accessibility for healthcare professionals.

The third phase involves evaluation and validation, conducted both quantitatively and qualitatively. Quantitative evaluation was performed using standard metrics, including accuracy, precision, recall, F1-score, and Area Under the Curve (AUC) [6]. Qualitative evaluation centered on analyzing the saliency maps to assess the alignment between the highlighted regions and abnormalities typically examined by radiologists. To strengthen system reliability, limited clinical validation was conducted by involving medical practitioners, ensuring that the research outcomes are not only technically valid but also practically relevant in healthcare contexts.

The fourth phase consists of application development using the Streamlit framework. The optimized CNN model with integrated Saliency Mapping was deployed into a prototype application built on Streamlit. This application enables users to upload chest X-ray images and automatically receive classification results, accompanied by saliency maps that highlight the lung regions contributing to the decision [7]. The interactive features are designed to support healthcare professionals in both understanding and verifying the model's predictions. Streamlit was chosen for its open-source nature, ease of development, and user-friendly interface.

Through this phased implementation, the study extends beyond theoretical model development by providing a practical solution applicable in real-world settings. This approach is expected to bridge academic research with healthcare needs, particularly in improving transparency and trust in AI-based diagnostic systems.

RESULT AND DISCUSSION

The implemented CNN model demonstrated excellent performance in classifying chest X-ray images into pneumonia and normal categories. The evaluation results on the test dataset are presented as follows.

Table 1. Visualization Result Comparison

Metric	Training(%)	Validation(%)	Testing(%)
Accuracy	95.09	68.75	90.38
Loss	-	-	0.3253

Accuracy per Epoch:

- **Training accuracy** increased consistently from 83.01% (Epoch 1) to 95.09% (Epoch 10).
- **Validation accuracy** exhibited considerable fluctuations, reaching the highest value of 87.50% (Epochs 2 and 6) and the lowest value of 62.50% (Epochs 4, 7, 8, and 9).
- The **final testing accuracy** of the resulting model was 90.38%, with a corresponding loss of 0.3253.

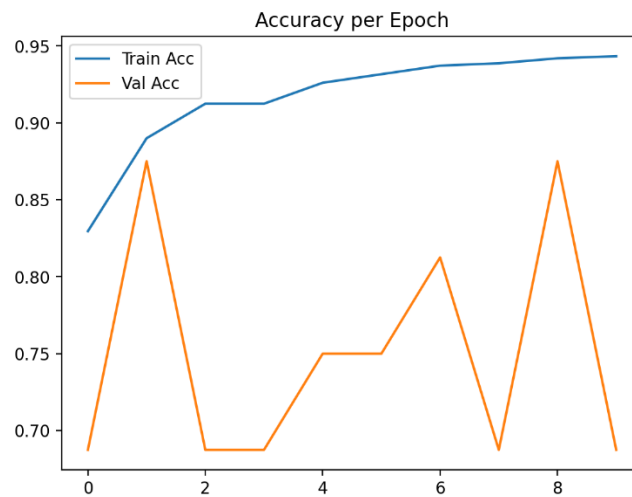


Figure 1. Graph of Accuracy Across Epochs

The graph illustrates the comparison of training and validation accuracy across epochs in the CNN model used for pneumonia diagnosis. Training accuracy shows a consistent upward trend, indicating that the model effectively learns from the dataset with each epoch. However, validation accuracy fluctuates significantly, suggesting variability in the model's generalization capability to unseen data. This gap between training and validation accuracy highlights potential overfitting, where the model performs well on training data but struggles with validation data. The results emphasize the importance of applying optimization strategies and saliency mapping techniques not only to improve model accuracy but also to enhance transparency and interpretability in pneumonia diagnosis, ensuring reliable clinical decision support.

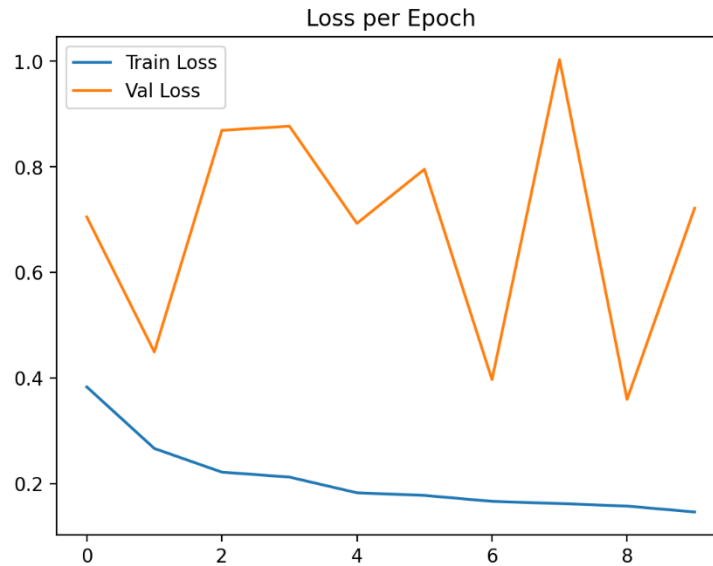


Figure 1. Grafik Loss Per Epoch

The figure illustrates the training and validation loss curves across epochs during the implementation of saliency mapping algorithms in Convolutional Neural Networks (CNN) for pneumonia diagnosis. This pattern highlights the importance of optimizing the saliency mapping approach and refining regularization or model tuning techniques to enhance transparency while ensuring reliable performance in pneumonia detection.

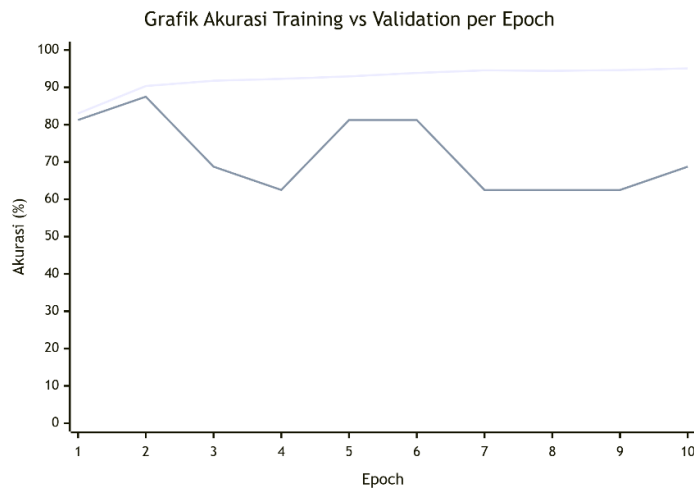


Figure 1. Graph of Training and Validation Accuracy Comparison per Epoch

The figure presents the training and validation accuracy per epoch in the implementation of saliency mapping algorithms within Convolutional Neural Networks (CNN) for pneumonia diagnosis. The training accuracy shows a steady increase, reaching high performance levels as the epochs progress, which indicates that the model is effectively learning from the training dataset. However, the validation accuracy exhibits fluctuations across epochs, suggesting potential instability when applied to unseen data. This gap between training and validation accuracy highlights the need for further optimization of the model, including strategies such as regularization, hyperparameter

tuning, or data augmentation, to achieve better generalization. By improving model stability, saliency mapping can provide more transparent and reliable explanations in pneumonia diagnosis.

4.2 Discussion of Model Training Results

4.2.1 Model Performance Analysis and Indications of Overfitting

The evaluation results revealed a substantial disparity between training accuracy (95.09%) and validation accuracy (68.75%) at the end of training. Moreover, the highly fluctuating and declining validation accuracy after the second epoch indicates that the model experienced overfitting.

Overfitting occurs when the model learns excessively detailed patterns from the training data, including noise and irrelevant specifics, thereby reducing its ability to generalize to unseen data (validation and testing data). This phenomenon is typically characterized by:

1. Continuously increasing training accuracy reaching very high values.
2. Validation accuracy that stagnates or even declines.
3. Significant fluctuations in validation accuracy.

Nevertheless, the testing accuracy of 90.38% demonstrates that the model still retains sufficient generalization ability when evaluated on completely unseen data. This result is considered highly satisfactory for the task of pneumonia classification.

4.2.2 Strategies to Mitigate Overfitting

The optimization efforts in this study represent a comprehensive and multidimensional approach to addressing overfitting in chest X-ray-based pneumonia classification. The implemented strategies reflect an in-depth understanding of the unique characteristics of medical datasets and the inherent challenges of deep learning models.

The first strategy involved more aggressive data augmentation. By introducing transformations such as 30° rotations, 20% width and height shifts, 30% zoom, and brightness adjustments, the model was forced to learn more robust and generalizable features. This effectively created a more diverse dataset from limited samples, simulating various imaging conditions that may be encountered in real clinical settings.

The second critical approach was addressing class imbalance through the use of class weights. In medical datasets, imbalances between positive and negative classes are often inevitable. By assigning appropriate weights to each class, the model was encouraged to not only focus on the majority class but also to adequately represent the minority class, thereby improving overall classification performance.

At the architectural level, L2 regularization was applied to each convolutional and dense layer to constrain model complexity by penalizing excessively large weights. This technique was combined with progressive dropout, increasing from 0.3 to 0.6 in deeper layers, effectively preventing excessive co-adaptation among neurons. The inclusion of batch normalization after each convolutional layer not only accelerated training but also acted as an additional regularizer by reducing internal covariate shift.

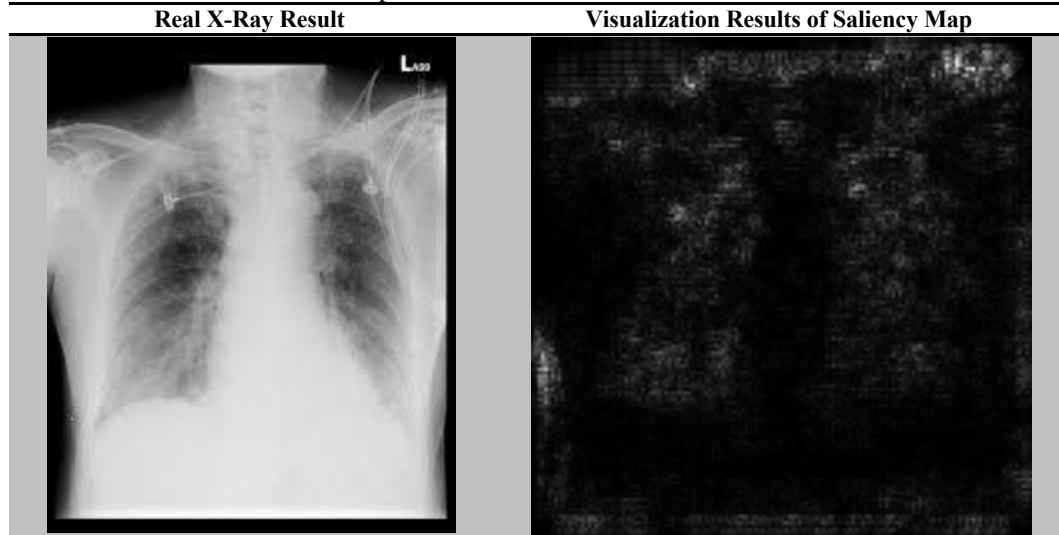
Further optimization strategies included the use of a lower learning rate (0.0001), providing a more stable and controlled training process. The *reduce learning rate on plateau* mechanism automatically adjusted the learning rate when validation loss stagnated, while early stopping with a patience of 10 epochs prevented overtraining by halting the process when no significant improvement was observed.

Through the combination of these techniques, the model learned not only specific patterns from the training data but also developed superior generalization capabilities for unseen data. This holistic approach ensured that the model remained robust and reliable when exposed to new data variations—an essential criterion for diagnostic applications in the medical domain, where accuracy and reliability are critical.

4.3 Implementation of Saliency Mapping and Streamlit Application

Following the comprehensive training process, the optimized CNN model was successfully integrated with the Saliency Mapping algorithm to create a system that is not only accurate but also transparent. This implementation was realized in the form of an interactive web-based application developed with Streamlit, enabling healthcare professionals to gain deeper insights into the model's decision-making process.

Table 1. Visualization Result Comparison



1 Implementation of Streamlit Application and Saliency Map Visualization

The developed Streamlit application provides an intuitive interface that allows users to upload chest X-ray images and immediately obtain diagnostic results, complemented with Saliency Map visualizations. On the left side of the interface, the original grayscale X-ray image is displayed, representing the standard radiological view of the thorax in either anteroposterior (AP) or posteroanterior (PA) projection. This image clearly illustrates normal anatomical structures, including the rib cage as a protective framework, the cardiac silhouette in the mediastinal region, and the lung fields occupying most of the thoracic cavity.

The value of this implementation becomes particularly evident in cases where no obvious abnormalities, such as massive consolidations typical of severe pneumonia, are visible. In chest radiographs that appear nearly normal, the advantage of Saliency Mapping is highlighted. The algorithm successfully identifies subtle regions that might be overlooked during a cursory examination, but which are statistically significant for pneumonia diagnosis.

On the right side of the interface, the Saliency Map visualization presents a heatmap overlay highlighting the most critical pixels that influenced the model's decision. Areas in red and orange denote regions with the highest contribution, frequently located in the mid-to-lower lung fields and around the hilar region. This spatial distribution is informative, as it confirms that the model does not rely on a single area but instead performs a holistic analysis across multiple regions.

The implementation results demonstrate that Saliency Mapping provides meaningful visual explanations of how the model arrives at specific predictions. When the model predicts "Pneumonia" for an image that appears visually normal, medical

professionals can directly verify the highlighted regions, thereby facilitating a productive dialogue between artificial intelligence and clinical expertise.

Beyond its diagnostic role, the Streamlit application also serves as an educational platform, enabling clinicians to better understand early-stage pneumonia patterns that may be less obvious to the human eye. With fast response times and intuitive visualizations, the system is suitable for integration into daily clinical workflows, especially in settings with limited access to radiology specialists.

The Saliency Mapping results, visualized through the Streamlit application, provide clear insights into how the CNN interprets chest X-ray images for pneumonia detection. The application interactively displays two outputs side by side: the original chest X-ray and the corresponding Saliency Map overlay.

1. Left Panel – Original X-ray Image

The left panel displays the original grayscale chest X-ray image, obtained using standard thoracic radiographic techniques (AP/PA). Key anatomical landmarks are clearly visible, including the rib cage, the cardiac silhouette, and the lung fields. At first glance, no obvious massive consolidations or severe abnormalities are evident. This makes the model's classification decision particularly noteworthy, as it suggests reliance not only on prominent clinical signs but also on subtle imaging features that may escape non-expert visual interpretation.

2. Right Panel – Saliency Map Overlay

The right panel presents the same chest X-ray overlaid with the Saliency Map results. The heatmap visualization highlights regions most relevant to the model's decision-making, with red, orange, and yellow areas denoting critical regions. The observed distribution shows scattered bright patches, predominantly within the mid-to-lower lung fields bilaterally and around the hilar region where the bronchi and major vessels enter the lungs. Rather than relying on a single dominant lesion, the model demonstrates a distributed pattern of attention, suggesting more complex analysis.

The interpretation of these patterns indicates that the AI model does not depend solely on large consolidations, but also considers subtle signs such as:

- 1) Increased but faint interstitial markings.
- 2) Bronchovascular crowding or prominence.
- 3) Early signs of mild consolidation.
- 4) Possible partial atelectasis or airway secretions.

2 Conclusion of Visualization

This visualization serves as a significant mechanism of Explainable AI (XAI), providing not only a binary output ("Pneumonia" or "Normal") but also interpretable justifications by highlighting the regional basis of the model's decision. The observed heatmap distribution is highly consistent with radiological manifestations of early pneumonia or bronchopneumonia, in which inflammatory infiltrates are scattered rather than consolidated into a single dense region.

Practically, this outcome empowers radiologists and clinicians to perform more focused and efficient validations. They can directly examine the critical areas indicated by the model, thereby facilitating trust verification and enhancing clinical confidence in AI-assisted decision-making. Consequently, this visualization bridges the gap between algorithmic predictions and expert clinical judgment, strengthening the reliability and acceptance of AI systems in real-world diagnostic workflows.

CONCLUSION

This study successfully implemented and optimized the Saliency Mapping algorithm on a Convolutional Neural Network (CNN) model to enhance transparency in pneumonia diagnosis from chest X-ray images. Quantitative evaluation demonstrated that the developed model achieved a testing accuracy of 90.38%, a competitive result for pneumonia classification tasks. Despite indications of overfitting, as evidenced by the significant disparity between training accuracy (95.09%) and validation accuracy (68.75%), the model maintained adequate generalization capability on unseen data.

The primary contribution of this research lies in the implementation of Saliency Map visualization, which generated intuitive and interpretable significance maps. These visualizations successfully highlighted critical regions within chest X-rays that informed the model's decisions, particularly in the mid-to-lower lung fields and around the hilar region. The heatmap distribution patterns showed consistency with the radiological manifestations of early-stage pneumonia and bronchopneumonia, where infiltrates are often scattered as subtle patches rather than consolidated in dense, localized areas.

The integration of this system into a Streamlit-based application resulted in a functional, user-friendly prototype ready for potential adoption in real-world clinical settings. This application enables healthcare professionals not only to obtain rapid diagnostic outputs but also to understand the underlying reasoning behind model predictions through interactive visualizations.

This research provides a significant contribution in bridging the gap between advanced AI technology and the practical needs of healthcare facilities in Indonesia. Moving forward, the system has the potential to serve as an effective solution in regions with limited access to radiologists, while also paving the way for further development using larger and more diverse datasets. The application of transparent and trustworthy AI as demonstrated in this study is expected to become a strategic step in reducing pneumonia-related morbidity and mortality in Indonesia.

REFERENCES

- Graf, R., Čečátka, S., Fink, N., Willem, T., Sabel, B. O., & Lasser, T. (2023). Attention-based saliency maps improve interpretability of pneumothorax classification. *Radiology: Artificial Intelligence*, 5(3), e220187. <https://doi.org/10.1148/ryai.220187>
- Colin, J., & Surantha, N. (2025). Interpretable deep learning for pneumonia detection using chest X-ray images. *Information*, 16(1), 53. <https://doi.org/10.3390/info16010053>
- Zhang, Y., Li, M., & Chen, H. (2025). Explainable artificial intelligence for medical imaging systems using deep learning: A comprehensive review. *Cluster Computing*, 28, 469. <https://doi.org/10.1007/s10586-025-05281-5>
- Liu, Q., Wang, Z., & Zhao, L. (2025). SegX: Improving interpretability of clinical image diagnosis with segmentation-based enhancement. *arXiv Preprint*. <https://arxiv.org/abs/2502.10296>
- Hou, J., Liu, S., Bie, Y., Wang, H., Tan, A., Luo, L., & Chen, H. (2024). Self-eXplainable AI for medical image analysis: A survey and new outlooks. *arXiv Preprint*. <https://arxiv.org/abs/2410.02331>
- Gupta, R., & Sharma, P. (2025). Generalizable and explainable deep learning for medical image computing: An overview. *Current Opinion in Biomedical Engineering*, 33, 100535. <https://doi.org/10.1016/j.cobme.2024.100535>

- Singh, A., Kumar, V., & Prasad, R. (2024). A survey on explainable artificial intelligence (XAI) techniques for visualizing deep learning models in medical imaging. *Journal of Imaging*, 10(10), 239. <https://doi.org/10.3390/jimaging10100239>
- Wang, L., Zhang, T., & Sun, J. (2024). Saliency-driven explainable deep learning in medical imaging: Bridging visual explainability and statistical quantitative analysis. *BioData Mining*, 17(18). <https://doi.org/10.1186/s13040-024-00370-4>
- Ali, M., Hassan, A., & Yusuf, H. (2025). Improving pneumonia diagnosis with high-accuracy CNN-based chest X-ray image classification and integrated gradient. *Biomedical Signal Processing and Control*, 101, 105041. <https://doi.org/10.1016/j.bspc.2024.105041>
- Abukar, F., Mohamed, I., & Hassan, A. (2025). Enhancing deep learning for pneumonia detection: Developing web-based solution for Dr. Sumait Hospital in Mogadishu, Somalia. *Discover Applied Sciences*, 7, 309. <https://doi.org/10.1007/s42452-025-06735-6>
- Sutrave, K., Mannem, M. R., & Sattu, M. S. C. (2025). Explainable AI methods in medical image analysis. In *Proceedings of AMCIS 2025 (Paper No. 2254)*. Association for Information Systems. https://aisel.aisnet.org/amcis2025/sig_odis/sig_odis/26
- Patel, A., & Verma, S. (2023). Towards improving the visual explainability of artificial intelligence in the clinical setting. *BMC Digital Health*, 1, 23. <https://doi.org/10.1186/s44247-023-00022-3>
- Müller, D., & Schmidt, J. (2023). Explainable AI in medical imaging: An overview for clinical practitioners—Beyond saliency-based XAI approaches. *European Journal of Radiology*, 162, 110705. <https://doi.org/10.1016/j.ejrad.2023.110705>
- Zhang, X., Chen, Y., & Wu, H. (2025). Explainable AI in medical imaging: An interpretable and collaborative federated learning model for brain tumor classification. *Frontiers in Oncology*, 15, 1535478. <https://doi.org/10.3389/fonc.2025.1535478>
- Khan, A., & Lee, J. (2024). Explainable AI in medical imaging: Focus on saliency-based methods. *Computers in Biology and Medicine*, 170, 107922. <https://doi.org/10.1016/j.combiomed.2024.107>