

A Data Pre-processing Strategy Utilizing Adaptive Masking for the Classification of Pediatric Pneumonia Using VGG-16

Yoshi Inne Herawati¹, Basuki Rahmat², Hendra Maulana³

¹Informatika, Fakultas Ilmu Komputer, UPN "Veteran Jawa Timur"

*^{2,3}Magister Informatika, Fakultas Ilmu Komputer, UPN "Veteran Jawa Timur"
basukirahmat.if@upnjatim.ac.id*

ABSTRACT

Pneumonia is still a leading cause of death in children, especially in areas with limited medical resources. This study aims to test several pre-processes to find the best set of pre-processes that can be applied to the children's chest X-ray dataset by applying adaptive masking, histogram equalisation, CLAHE and Gaussian blur. Then, childhood pneumonia is classified using a CNN architecture, namely VGG-16. By applying these pre-processing methods, this study is divided into several scenarios. The highest accuracy was obtained from scenario 1, which used a combination of adaptive masking, histogram equalization and Gaussian blur, resulting in an accuracy of 94%. Scenario 2 uses histogram equalization and Gaussian blur with an accuracy of 92%. Then Scenario 3 uses a histogram equalization replacement for CLAHE with a combination of adaptive masking, CLAHE and Gaussian blur with 93% accuracy. Finally, scenario 4 uses a combination of CLAHE and Gaussian blur methods with 91% accuracy. In addition, this research also addresses the challenges posed by unbalanced data sets and the need for highly accurate detection tools.

Keywords: Pneumonia, adaptive masking, Classification, pre-processing, VGG-16.

INTRODUCTION

Pneumonia is a disease of the lungs caused by inflammation due to bacterial, fungal or viral infections (Wati, Irsyad, & Rivan, 2020). This disease is still a serious concern, especially considering the high mortality rate among young children due to pneumonia. According to data from the Badan Pusat Statistik in 2018, East Nusa Tenggara province reported a total of 3,529 cases of pneumonia in children under the age of five. This number is low compared to East Java Province in 2022, which recorded a total of 92,118 cases of pneumonia. Even in 2019, Unicef Indonesia reported that pneumonia was the leading cause of death among children, with the percentage reaching 36%.

With such a significant number, one of the causes of delayed diagnosis is the limited number of health workers. The limited number of health workers available to make an initial diagnosis can lead to potential delays in patient treatment. Typically, there are several stages in the diagnostic process, such as taking blood samples and taking X-rays (Nuraeni & Rahmawati,

2019). The limited number of medical personnel who can perform the tests or interpret the results of X-rays is one of the main factors that can delay the diagnosis, which can affect the delay in patient treatment. Therefore, there is a need for technology that can assist in diagnosis, enabling healthcare professionals to detect the presence of pneumonia in chest X-rays at an early stage. In its application, to achieve high accuracy results, it is necessary to handle the data set well before the researcher carries out the training process.

There are three relevant studies that form the basis of this study. First, research by Morteza Heidari et al (2020) developed a CNN-based CAD scheme to detect COVID-19 pneumonia from chest X-ray image, achieving 94.5% accuracy. They used pre-processing in the form of histogram equalisation, bilateral low-pass filtering and diaphragm removal. Second, Agata Gielczyk et al (2022) evaluated six pre-processing scenarios on X-ray images and found that the combination of adaptive masking, histogram equalisation and Gaussian blur produced the highest accuracy of 98.3%. Thirdly, research by Gaurav Labhane et al (2020) compared several CNN architectures to detect pneumonia in children, with VGG-16 and Inception-V3 showing the best accuracy of 98%. Based on these studies, this research focuses on improving X-ray images by using the best pre-processing method and VGG-16 architecture for classification of childhood pneumonia.

Based on the consideration of several literature studies carried out in relation to the pre-processing of X-ray image datasets and CNN methods, the author chose to improve X-ray images as the main focus of research using the scenario of the best pre-processing method found in previous research by Agata Gielczyk and her friends, and tried to add the CLAHE method. In addition, the author uses one of the CNN architectures, namely VGG-16, in order to obtain accuracy results with the best value. VGG or Visual Geometry Group itself is one of the architectures developed after AlexNet, which can be used for feature extraction in the convolution layer. (Saputro et al., 2022). This method consists of functional, dropout, flatten, batch normalisation, dense, and activation with a total of 16 layers (Setiawan et al., 2022). With this research, the author hopes to contribute to the accurate classification of childhood pneumonia.

METHODS

In this study, three scenarios are carried out to determine the effect of the pre-processing method used on the accuracy of the resulting model. The following is a breakdown of the four research scenarios. A research flowchart is shown on figure 1 below.

1. The first scenario was carried out using the VGG-16 algorithm with adaptive masking, histogram equalisation and Gaussian blur methods at the pre-processing stage.
2. The second scenario experiments with the VGG-16 algorithm using histogram equalisation and gaussian blur methods in the pre-processing stage.
3. In the third scenario, the author used the VGG-16 algorithm by changing the histogram equalisation method with the CLAHE method in the pre-processing stage, namely adaptive masking, CLAHE and Gaussian blur.
4. Fourth, experiments were conducted using the VGG-16 algorithm with CLAHE and Gaussian blur methods.

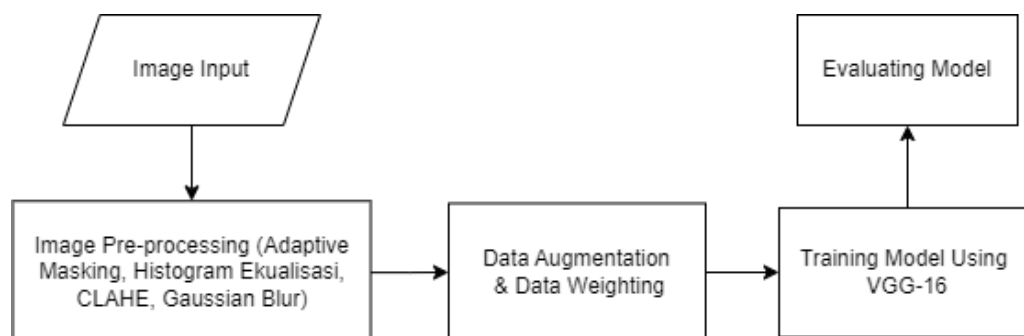


Figure 1. Methods

Data Collection

The dataset used in this study was drawn from public data uploaded by Paul Mooney to the Kaggle platform under the title 'Chest X-Ray Images (Pneumonia)'. This data was obtained from the Guangzhou Women and Children's Medical Centre, which contains chest X-rays of children aged 1-5 years as part of routine patient care. The dataset consists of a total of 5,856 X-rays, which are divided into two labels: normal and pneumonia infected with bacteria or viruses.

Table 1. Pneumonia Dataset

Normal	Pneumonia Virus	Pneumonia Bacteria



As illustrated in the image above, the lungs of a healthy individual typically exhibit a clear appearance, devoid of any discernible patches. This is characterised by the entire network of blood vessels and the diaphragm (the muscle separating the chest and abdomen), which is clearly visible. In contrast, lungs diagnosed with pneumonia caused by viruses or bacteria tend to exhibit blurry spots covering the lung area, as illustrated in Table 1. Subsequently, the data was divided into three distinct categories: training data, test data, and validation data. This division was implemented for the purpose of training and testing the pneumonia classification model. The total number of data points included in the training data set was 5,216, while the test data set comprised 640 data points. The validation data set, on the other hand, consisted of a total of 16 data points. The training data set, which was the largest of the three, was necessary for the training process, as it required a substantial amount of data to achieve high levels of accuracy.

Image Pre-processing



Figure 2. Image Pre-processing

Resize Image

Resizing the image is the process of changing the size of the image to achieve harmony in size. Therefore, in order to adapt to the training needs, the author changes all the existing image data to a size of 224×224 , where the size is the appropriate size for the input of the architecture to be used, namely VGG-16.

Contrast Enhancing

- a. Histogram Equalization

After the image resize process has successfully resized the original image, the next step is to increase the intensity of the image in order to reach the details of the image that were previously difficult to see due to the low contrast. One method that can effectively improve image quality is histogram equalisation.

b. CLAHE

In this study, the author used 8 grids with a Clip Limit of 3 to avoid excessive contrast enhancement while maintaining image quality and ensuring that relevant details are preserved.

c. Gaussian Blur

Gaussian blurring is a way of improving an image by reducing unnecessary detail and noise in an image. This process is necessary to make it easier for the system to recognise existing patterns. Therefore, the author applies Gaussian blur to reduce noise. In this case, the author uses a Gaussian kernel with a size of 5×5 to perform the calculation process.

d. Adaptive Masking

Adaptive masking is one of the ways in which the image can be enhanced to suit the requirements. The use of adaptive masking is often different in each case because it is required to suit the problems that exist in the project or case. Below is a visualisation of the adaptive masking steps that are performed.

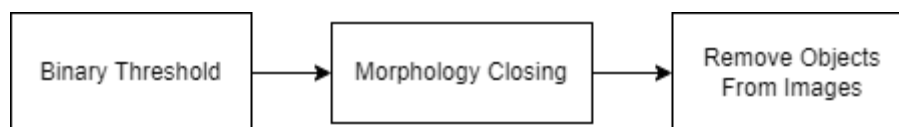


Figure 3. Adaptive Masking

In this case, the adaptive masking used is to remove the diaphragm object from the chest image by performing several steps as shown in Figure 3.

- a. Apply binary thresholding after finding the minimum and maximum values of pixel intensity.
- b. Use morphological closure to create a suitable mask for the shape of the diaphragm.
- c. The last step is to remove the membrane from the original image using the mask created in the previous step.

Data Augmentation & Data Weighting

Prior to model training, data augmentation and weighting are performed to address imbalances in the dataset. Augmentation aims to create new variations of the existing data without physically duplicating it, using techniques such as rotation, mirroring and zooming. This helps the model to recognise more diverse data and reduces the risk of overfitting. Data weighting is then applied so that classes with less data are still noticed by the model, reducing bias towards dominant classes.

Training Model VGG-16

In this study, the VGG-16 model is employed as the primary tool for feature extraction from the pre-processed training data. The proposed architecture comprises 13 convolutional layers and three fully connected layers, which are responsible for the fully connected layer stage and are illustrated in Figure 4. Prior to this, VGG-16 had been trained on the ImageNet dataset, and thus was not trained from scratch, but rather used as a pre-trained model. This accelerates the training process, as the underlying features have already been identified from previous training.

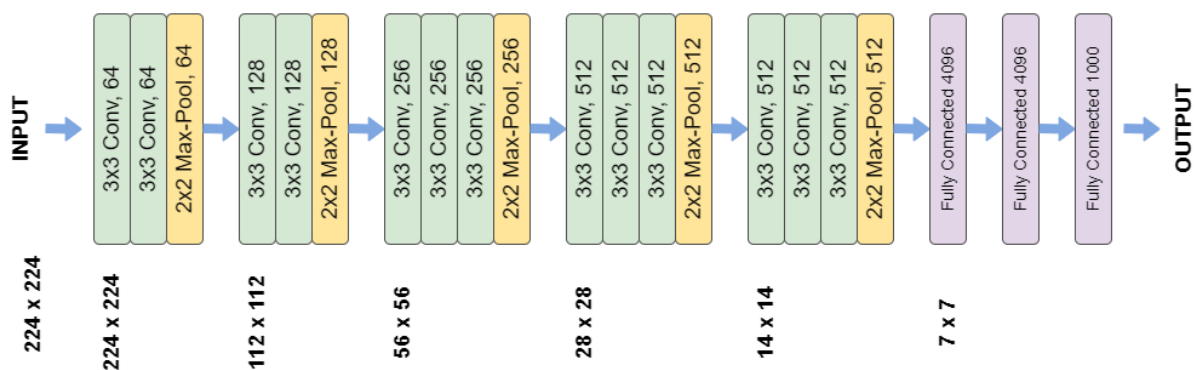


Figure 4. VGG-16 Architecture

During the research, the authors made some adjustments to the fully connected layer to make the model more suitable for specific needs. The model was trained in 10 epochs, with three repetitions of the training process, to achieve stable and optimal accuracy. The combination of convolutional and fully connected layers of VGG-16 enables the model to recognise complex patterns in images and produce accurate classifications.

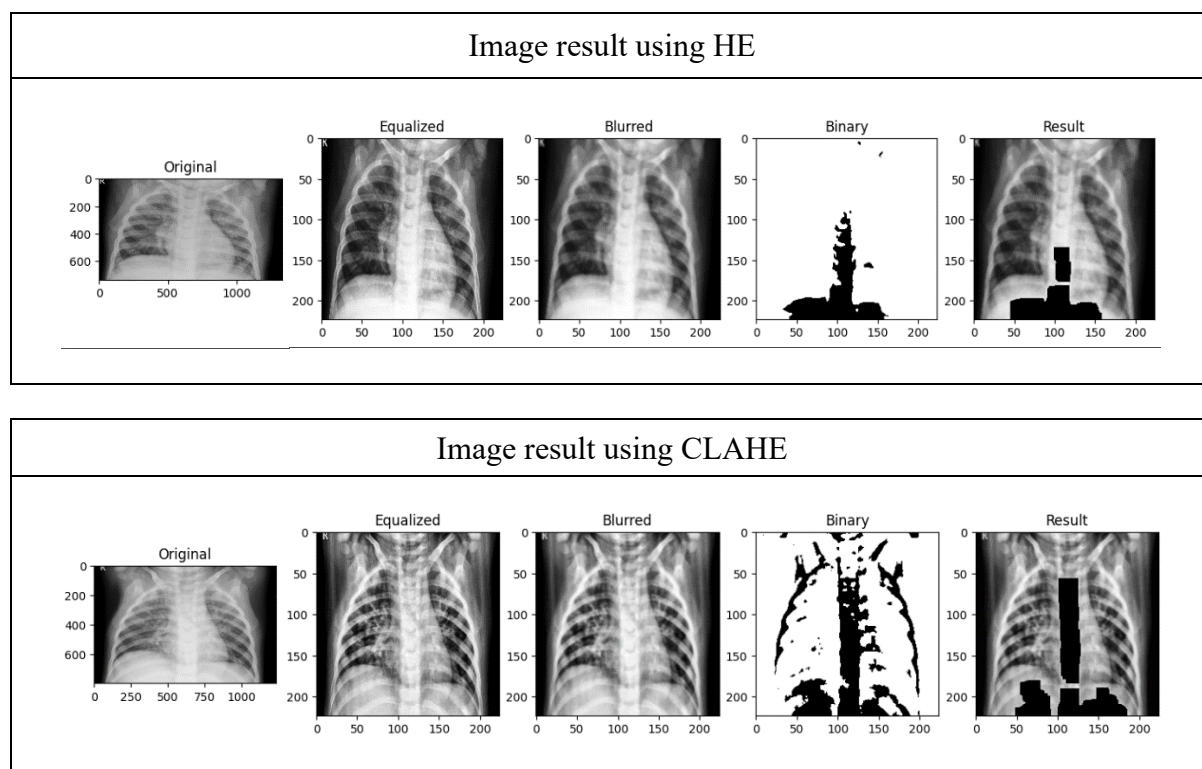
Evaluation

In the testing phase, a series of procedures are conducted, including the input of the pre-prepared test data into the trained VGG-16 model. Following the completion of the testing process, the prediction outcomes will be contrasted with the original data set. Then, the number of accurate or incorrect predictions will be quantified using the confusion matrix.

RESULTS AND DISCUSSIONS

Image Pre-processing

Table 2. Image Pre-processing Result



The process of data pre-processing comprises a series of essential steps that are undertaken in order to prepare the image in question for subsequent classification. The image is then resized to 224x224 pixels, after which contrast enhancement is performed using histogram equalisation and CLAHE. A Gaussian blur is subsequently applied with the objective of reducing noise. Subsequently, adaptive masking is conducted through binary thresholding to differentiate the primary objects, followed by morphological closing. Subsequently, the diaphragm area is excised, leaving only the lung area for subsequent analysis. This process enhances the image quality, thereby optimising the accuracy of the classification model.

Evaluation

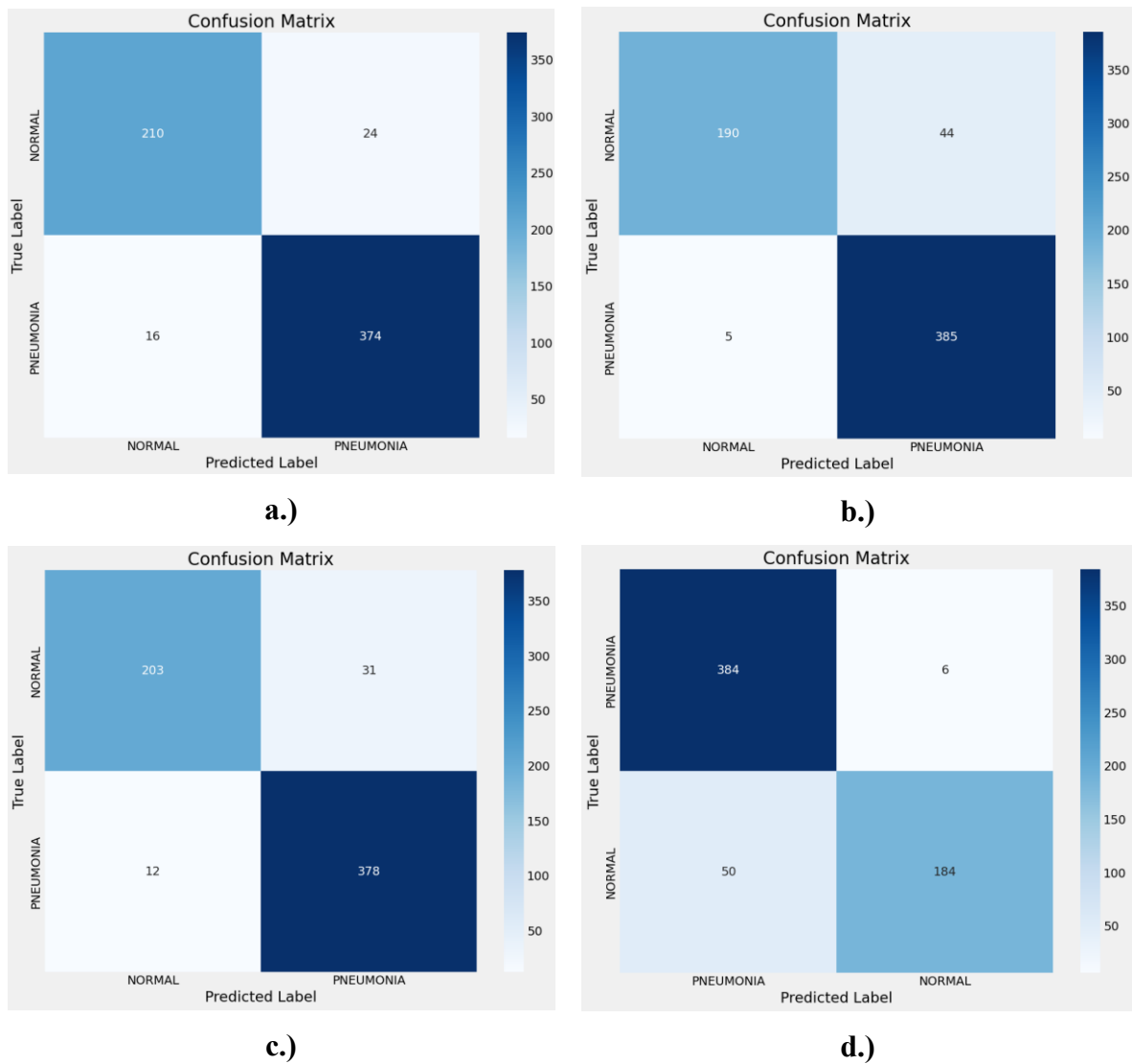


Figure 5. Confussion Matrix

(a. Adaptive Masking, HE, and Gaussian Blur; b. HE and Gaussian Blur; c. Adaptive Masking, CLAHE, and Gaussian Blur; d. CLAHE and Gaussian Blur)

Table 3. Model Evaluation

Test Scenarios	Methods	Class	Accuracy	Precision	Recall	F1-Score
1	Adaptive Masking, HE, & Gaussian Blur	Pneumonia	94%	93%	90%	91%
		Normal		94%	96%	95%
2	HE & Gaussian Blur	Pneumonia	92%	90%	99%	94%
		Normal		97%	81%	89%

3	Adaptive Masking, CLAHE, & Gaussian Blur	Pneumonia	93%	92%	97%	95%
		Normal		94%	87%	90%
4	CLAHE & Gaussian Blur	Pneumonia	91%	88%	98%	93%
		Normal		97%	79%	87%

The four test scenarios, which were conducted, demonstrated variations in results based on the pre-processing methods that were applied prior to classification using VGG-16. In the initial test scenario, which employed adaptive masking, histogram equalisation and Gaussian blur, the model attained a 94% accuracy rate with 93% precision for the normal class and 94% for the pneumonia class. The recall for the pneumonia class reached 96%, while that for the normal class was 90%. In the second test scenario, histogram equalisation and Gaussian blur were employed, resulting in an accuracy of 92%. Although the precision for the normal class was high at 97%, the recall was low at 81%. In contrast, the pneumonia class demonstrated a precision of 90% and a remarkably high recall of 99%. The third test scenario, which employed adaptive masking, CLAHE, and Gaussian blur, demonstrated an accuracy rate of 93%, with a precision of 94% for the normal class and 92% for the pneumonia class, and a recall rate of 87% and 97%, respectively. The fourth test scenario, which employed solely CLAHE and Gaussian blur, yielded an accuracy rate of 91%. The precision of the normal class was 97%, while that of the pneumonia class was 88%. The recall of the normal class was 79%, while that of the pneumonia class was 98%. Of the four scenarios, the first test scenario demonstrated the optimal performance, exhibiting a balance of high precision and recall in both classes.

CONCLUSION

This study examines the classification of pneumonia in children using chest X-ray images with the adaptive masking method and several other pre-processing techniques. Of the four combination scenarios tested, the combination of adaptive masking, histogram equalisation, and Gaussian blur was found to provide the most optimal results, with an accuracy of 94%. The data augmentation and weighting processes proved effective in addressing the imbalance of the dataset, enabling the model to learn a more diverse range of data shapes and avoid neglecting classes with limited data.

Despite the application of the CLAHE method for image enhancement, the results demonstrated that it generated a greater amount of noise, which ultimately led to a reduction in the model's accuracy in comparison to histogram equalisation. The combination of test scenarios with CLAHE demonstrated slightly reduced accuracy, with rates of 93% and 91%, respectively, in comparison to the combination utilising regular histogram equalisation. Furthermore, the adaptive masking method was demonstrated to consistently enhance model accuracy across all scenarios, irrespective of whether CLAHE or histogram equalisation was employed.

In light of these findings, the researcher proposes that future research should incorporate additional techniques for addressing unbalanced datasets, with the aim of optimising model training. Additionally, it is recommended to investigate the potential of alternative CNN architectures, such as ResNet50, MobileNet, or analogous architectures, to further enhance the accuracy of pneumonia classification models.

REFERENCES

- Wati, R. A., Irsyad, H., & Rivani, M. E. A. R. (2020). Klasifikasi Pneumonia Menggunakan Metode Support Vector Machine. *Jurnal Algoritme*, 1(1), 21–32. <https://doi.org/10.35957/algoritme.v1i1.429>
- Badan Pusat Statistik Provinsi Jawa Timur. (2022). Jumlah Jenis Penyakit Malaria, TB Paru, Pneumonia, Kusta Menurut Kabupaten/Kota di Provinsi Jawa Timur 2022. Diakses dari <https://jatim.beta.bps.go.id/id/statistics-table/1/MzAwMSMx/jumlah-jenis-penyakit-malaria-tb-paru-pneumonia-kusta-menurut-kabupaten-kota-di-provinsi-jawa-timur-2022.html>
- Badan Pusat Statistik Provinsi Nusa Tenggara Timur. (2018). Jumlah Kasus Penyakit Pneumonia Menurut Jenis Penyakit 2018. Diakses dari <https://ntt.beta.bps.go.id/id/statistics-table/1/NzE0IzE=/jumlah-kasus-penyakit-pneumonia-menurut-jenis-penyakit-2018.html>
- UNICEF Indonesia. Kesehatan. Diakses dari <https://www.unicef.org/indonesia/id/kesehatan>

- Nuraeni, T., & Rahmawati, A. (2019). Pneumonia Pada Balita Dan Penanganan Yang Tepat. *Seminar Nasional Kesehatan Masyarakat UMS*, 147–151. <https://publikasiilmiah.ums.ac.id/xmlui/handle/11617/11862>
- Heidari, M., Mirniaharikandehi, S., Khuzani, A. Z., Danala, G., Qiu, Y., & Zheng, B. (2020). Improving the performance of CNN to predict the likelihood of COVID-19 using chest X-ray images with preprocessing algorithms. *International Journal of Medical Informatics*, 144(September), 104284. <https://doi.org/10.1016/j.ijmedinf.2020.104284>
- Labhane, G., Pansare, R., Maheshwari, S., Tiwari, R., & Shukla, A. (2020). Detection of Pediatric Pneumonia from Chest X-Ray Images using CNN and Transfer Learning. *Proceedings of 3rd International Conference on Emerging Technologies in Computer Engineering: Machine Learning and Internet of Things, ICETCE 2020, February*, 85–92. <https://doi.org/10.1109/ICETCE48199.2020.9091755>
- Gielczyk, A., Marciniak, A., Tarczewska, M., & Lutowski, Z. (2022). Pre-processing methods in chest X-ray image classification. *PLoS ONE*, 17, 1–11. <https://doi.org/10.1371/journal.pone.0265949>
- Setiawan, D., Widodo, S., Ridwan, T., & Ambari, R. (2022). Perancangan Deteksi Emosi Manusia berdasarkan Ekspresi Wajah Menggunakan Algoritma VGG16. *e-Proceeding of Engineering*, 11(01), 1–11.
- Saputro, A., Mu'min, S., Moch. Lutfi, & Putri, H. (2022). Deep Transfer Learning Dengan Model Arsitektur Vgg16 Untuk Klasifikasi Jenis Varietas Tanaman Lengkeng Berdasarkan Citra Daun. *JATI (Jurnal Mahasiswa Teknik Informatika)*, 6(2), 609–614. <https://doi.org/10.36040/jati.v6i2.5456>