

AI-Driven Lungs and Breast Cancer Detection

Mrs. Harshitha B K*, M Ronith¹, Karthik N Hebbar¹, Rakshith Srinivasan¹, Shravan Karthik¹

*Assistant Professor, Department of Information Science and Engineering, Dayananda Sagar Academy of Technology and Management, Bengaluru, India.

¹ B.E. Students, Department of Information Science and Engineering, Dayananda Sagar Academy of Technology and Management, Bengaluru, India.

Abstract

Cancer, particularly lung and breast cancer, remains a significant global health challenge, underscoring the urgency for early detection methods to enhance survival rates. This paper centers on employing artificial intelligence (AI) techniques for the early detection of lung and breast cancer using deep learning and Convolutional Neural Network (CNN) models. This paper presents a comprehensive methodology for the development and evaluation of lung and breast cancer detection models using medical imaging data. Leveraging TensorFlow and a T4 GPU in Google Colab, the methodology navigates through key stages, including dataset collection, data pre-processing, model training, and model evaluation. For lung cancer detection, Chest CT-Scan images were collected and processed using data augmentation techniques, leading to the development of a robust EfficientNetB0 CNN architecture. The trained model demonstrated high accuracy and generalization ability, with promising results on validation and test datasets. For breast cancer detection, mammography images from the "cbis-ddsm-breast-cancer-image-dataset" were utilized. A VGG16 CNN architecture was tailored for binary classification, achieving an impressive accuracy on the test dataset.

Index Terms—Breast cancer, Lung cancer, Machine learning, Deep learning, Convolutional neural network, EfficientNetB0.

I. INTRODUCTION

Cancer remains one of the most pressing health challenges worldwide, with lung and breast cancer being among the leading causes of cancer-related mortality. Early detection of these cancers is paramount for improving patient outcomes and survival rates. In recent years, artificial intelligence (AI) techniques, particularly deep learning and Convolutional Neural Networks (CNNs), have emerged as powerful tools for medical image analysis and disease detection.

This paper focuses on the application of AI-driven approaches to facilitate early detection of lung and breast cancer using medical imaging data. Lung cancer, notorious for its aggressive nature and often late-stage diagnosis, necessitates robust detection methods to enable timely intervention. Similarly, breast cancer, affecting millions of individuals globally, requires accurate and efficient detection strategies to enhance treatment efficacy and patient outcomes.

To address these critical healthcare needs, we present a comprehensive methodology for developing and evaluating lung and breast cancer detection models. Our methodology encompasses key stages including dataset collection, data pre-processing, model training, and evaluation. For lung cancer detection, we employ Chest CT-Scan images, while mammography images are utilized for breast cancer detection.

In our approach, we meticulously pre-process the data, augmenting it to enhance model robustness and generalization. We design custom CNN architectures tailored to each cancer type, leveraging state-of-the-art models such as EfficientNetB0 for lung cancer and VGG16 for breast cancer detection. Through rigorous model training and evaluation, we demonstrate promising performance metrics, including high accuracy and generalization ability, underscoring the efficacy of our methodologies.

Furthermore, we discuss future enhancements to our models, including the integration of advanced CNN architectures, transfer learning with larger datasets, incorporation of clinical data, real-time deployment strategies, and continuous model monitoring and updating. These enhancements aim to further improve the accuracy, efficiency, and applicability of our cancer detection models in clinical practice, ultimately contributing to better patient care and outcomes.

II. METHODOLOGY

Lung cancer model: The methodology presented herein outlines a comprehensive approach for the development and evaluation of a lung cancer detection model utilizing Chest CT-Scan images. With the aim of addressing the critical need for accurate and efficient detection methods, this methodology meticulously navigates through key stages, from dataset collection to model evaluation. Leveraging TensorFlow as the primary framework and harnessing the computational power of a T4 GPU in Google Colab, each step is carefully crafted to ensure robustness and efficacy in the classification task.

1. Dataset Collection: In this method, a dataset of Chest CT-Scan images was collected for the purpose of lung cancer detection. The dataset comprises images categorized into four distinct classes: adenocarcinoma, large cell carcinoma, normal, and squamous cell carcinoma. These images serve as the foundational data for training and evaluating the proposed classification model. The data was sourced from Kaggle, a renowned platform for datasets and machine learning competitions, ensuring access to high-quality and diverse data for research and development purposes.

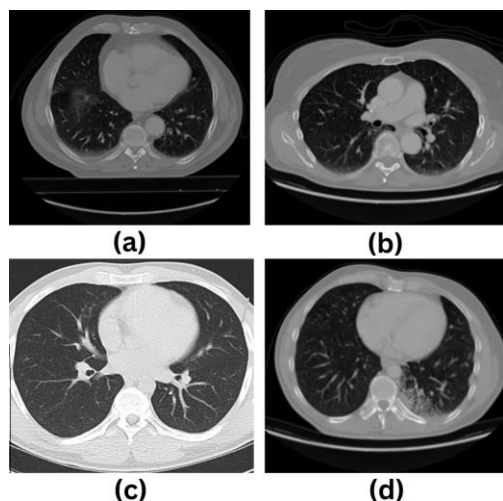


Fig. 2.1. Lung cancer dataset

In the above Fig 2.1, provided images (a, b, c, d), each represents a distinct classification within the dataset. Image (a) is classified as adenocarcinoma, image (b) as large cell carcinoma, image (c) as normal, and image (d) as squamous cell carcinoma. These classifications denote different types of conditions or states within the dataset

2. Data Pre-processing: In this method, a meticulous approach to data pre-processing was adopted to ensure the suitability of the images for subsequent analysis. Initially, the images were loaded and resized to a uniform dimension of 224x224 pixels using the OpenCV library. This step standardizes the image dimensions, ensuring consistency throughout the dataset. Subsequently, the dataset was split into training and testing sets, with 20% reserved for testing. This partitioning ensures that the model's performance can be accurately assessed on unseen data, thereby gauging its ability to generalize.

3. Data Augmentation: Data augmentation techniques were employed to enhance the variability of the training data, thereby improving the model's robustness and generalization ability. Techniques such as rotation, width and height shifting, shear, zooming, and horizontal flipping were applied using the ImageDataGenerator from TensorFlow to generate effectively increasing the diversity of the dataset. This augmentation process helps mitigate overfitting by exposing the model to a wider range of variations present in the data, ultimately leading to better performance on unseen samples.

4. CNN Architecture: The convolutional neural network (CNN) architecture employed in this work is constructed around the EfficientNetB0 model, which is pre-trained on ImageNet. This choice is strategic, as EfficientNet models are renowned for their superior performance in terms of accuracy and efficiency. Through a series of layer manipulations, the architecture was tailored to suit the classification task, ultimately culminating in a multi-class classification output layer with softmax activation. This architecture leverages the innovative scaling method of EfficientNet models, optimizing both depth, width, and resolution for enhanced performance.

5. Model Training: Following the construction of the CNN architecture, the model was trained using the compiled configuration. The training process involved iterating through the training dataset for multiple epochs while updating the model's weights to minimize the defined loss function. During training, the model parameters were optimized using the Adam optimizer, which adaptively adjusts the learning rate for each parameter. This adaptive optimization scheme helps accelerate convergence and improve model performance.

To monitor the training progress and ensure effective convergence, several callbacks were employed. TensorBoard was utilized for visualization, allowing for real-time tracking of key metrics such as loss and accuracy. ModelCheckpoint was utilized to save the model with the best performance on the validation set, ensuring that the optimal model weights were preserved for future

use. Additionally, ReduceLRonPlateau dynamically adjusted the learning rate during training based on the validation loss, facilitating convergence by fine-tuning the learning process.

Throughout the training process, detailed metrics such as loss and accuracy were logged, providing insights into the model's performance. By analyzing these metrics, adjustments to the model architecture or hyperparameters could be made to improve performance further.

6. Model Evaluation: Once the model completed training, its efficacy was evaluated on the held-out test dataset to assess its ability to generalize to unseen data. During evaluation, the trained model processed the test images, and predictions were compared against ground truth labels to compute performance metrics. Key evaluation metrics included accuracy and loss, which provide a quantitative measure of the model's performance. Accuracy represents the proportion of correctly classified instances out of the total number of instances, while loss quantifies the discrepancy between predicted and true labels.

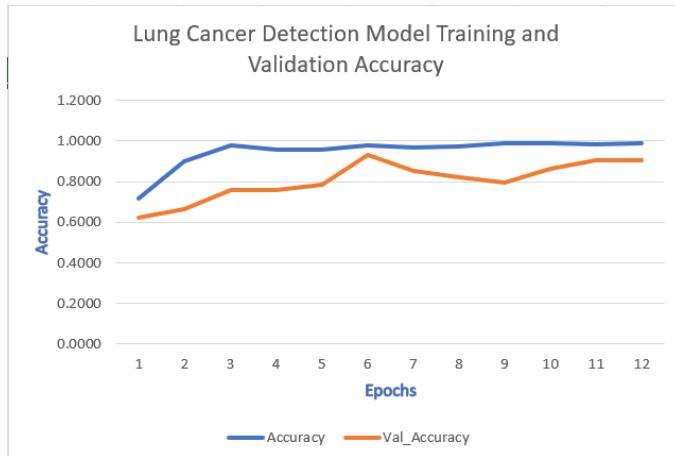


Fig. 2.2. Lung Cancer Detection Model Training and Validation Accuracy Graph

acquisition to model evaluation. By leveraging the TensorFlow framework and harnessing the computational prowess of a T4 GPU within Google Colab, each stage is meticulously designed to guarantee resilience and effectiveness in the classification endeavour.

1. Dataset Collection: In this method, we utilized the "cbis-ddsm-breast-cancer-image-dataset" for training and evaluating our breast cancer detection model. This dataset is a well-known resource in the field of medical imaging, containing a diverse collection of mammography images annotated with pathology labels, making it suitable for training machine learning models for breast cancer detection tasks. By leveraging this dataset, the model benefits from a rich source of annotated medical images, enabling it to learn complex patterns indicative of breast cancer presence.

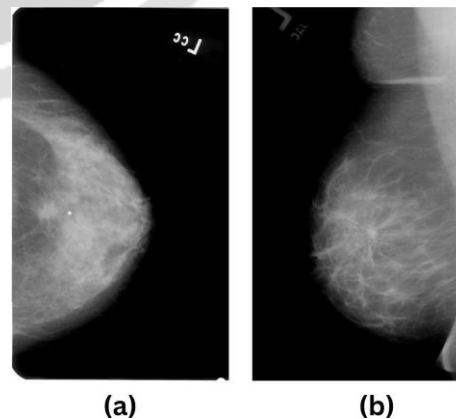
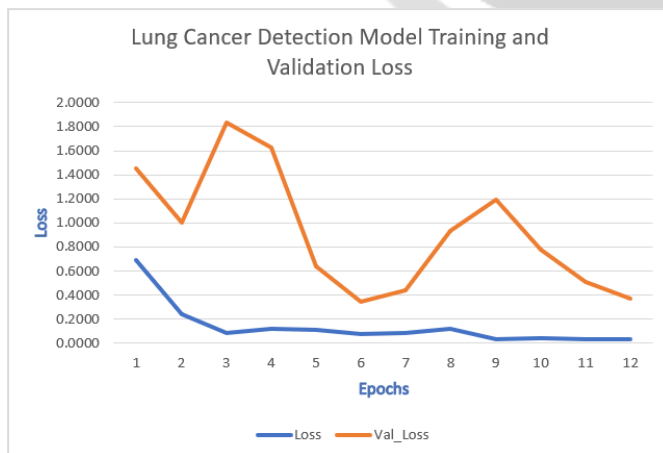


Fig. 2.3. Lung Cancer Detection Model Training and Validation Loss Graph dataset

Fig. 2.4. Breast cancer

Figures 2.2 and 2.3 represent the graphs showing the training/validation loss and accuracy, respectively, of the lung cancer detection model. They provide visual insights into the model's performance during training, aiding in the interpretation of its learning progress and generalization abilities.

In this method, after training for 12 epochs, the model achieved a training accuracy of approximately 95% and a training loss of approximately 0.0296. During validation, the model attained an accuracy of around 90.54% with a validation loss of about 0.3688. These metrics indicate that the model learned effectively from the training data and generalized well to unseen data.

Breast cancer model: The methodology detailed herein presents a comprehensive framework for the creation and assessment of a breast cancer detection model utilizing mammography images. With a primary objective of addressing the critical necessity for precise and efficient detection methods in breast cancer diagnosis, this methodology intricately guides through pivotal stages, ranging from dataset

In the above Fig 2.4, provided images (a, b), are classified into two categories: malignant and benign. Image (a) is classified as malignant, indicating the presence of harmful cancerous cells, while image (b) is classified as benign, suggesting the absence of cancerous growth or non-threatening cellular abnormalities. These classifications are crucial in medical diagnosis, guiding treatment decisions and prognosis for patients.

2. Data Pre-processing: The subsequent definition of the image processor function underscores the importance of data pre-processing, as it handles the reading and processing of images from the specified dataset directory. This function, equipped with GPU configuration capabilities, optimizes computational efficiency, a crucial aspect when dealing with large-scale image datasets.

Moving forward, the script loads the breast cancer image dataset, a fundamental step in any machine learning task, followed by the encoding of labels into a binary format, essential for binary classification tasks like distinguishing between malignant and benign cases. Furthermore, the dataset is intelligently partitioned into training and validation subsets, ensuring the model's ability to generalize to unseen data.

3. Data Augmentation: Data augmentation techniques were employed using Keras' ImageDataGenerator to enrich the training data with diverse transformations such as rotation, shifting, and flipping. These augmentation strategies not only expand the dataset but also enhance the model's robustness against variations in input images, essential for real-world applicability.

4. CNN Architecture: The model architecture was meticulously designed, leveraging the VGG16 convolutional neural network, a pre-trained model known for its effectiveness in image classification tasks. VGG16 stands out for its simplicity and interpretability, comprising a series of convolutional and pooling layers followed by fully connected layers. This architecture's deep yet straightforward structure facilitates efficient feature extraction, making it suitable for a wide range of image classification tasks.

Moreover, the hierarchical feature representations learned by VGG16 are particularly advantageous for medical imaging tasks, where subtle patterns and textures can signify disease presence. By leveraging these learned features, the model can effectively discern between benign and malignant breast tissue, contributing to accurate cancer diagnosis. To tailor the VGG16 model to the specific requirements of breast cancer detection, custom fully connected layers and a softmax output layer were added. These layers enable the model to make binary classifications, distinguishing between malignant and benign cases with precision.

5. Model Training: In the training phase, the carefully crafted VGG16 convolutional neural network architecture is optimized to discern between benign and malignant breast tissue, contributing to accurate cancer diagnosis. The process begins with compiling the model, where we specify appropriate loss functions and optimizers to guide the learning process. For our breast cancer detection task, binary cross-entropy serves as the loss function, facilitating the differentiation between two classes – malignant and benign. The Adam optimizer is chosen for its effectiveness in adjusting the model's parameters to minimize the loss function.

To enhance model generalization and prevent overfitting, several techniques are employed. One such method is early stopping, implemented as a callback mechanism. Early stopping monitors the validation loss during training and halts the training process when the validation loss begins to increase, indicating that the model's performance on unseen data is deteriorating. By terminating training at this optimal point, early stopping helps prevent the model from memorizing the training data and promotes better generalization to new data.

Moreover, regularization techniques such as dropout were applied to further mitigate overfitting. Dropout randomly disables

a fraction of neurons during training, forcing the model to learn more robust and generalized representations of the data.

The model’s performance metrics and training history are meticulously logged and analyzed throughout the training process. This comprehensive tracking allows for a thorough understanding of the model’s learning dynamics and convergence behaviour, enabling adjustments to be made as necessary.

6. Model Evaluation: Following the model training phase, rigorous evaluation is conducted to assess the performance and generalization ability of the breast cancer detection model. This evaluation phase is crucial for validating the model’s efficacy in real-world scenarios and ensuring its reliability in clinical settings.

The trained model is evaluated on an independent test dataset, distinct from the training and validation sets, to provide an unbiased assessment of its performance. The test dataset consists of unseen mammography images, representative of the diverse cases encountered in clinical practice. Performance metric such is computed to quantify the model’s accuracy. The accuracy measures the proportion of correctly classified instances out of all instances, providing an overall assessment of the model’s correctness.

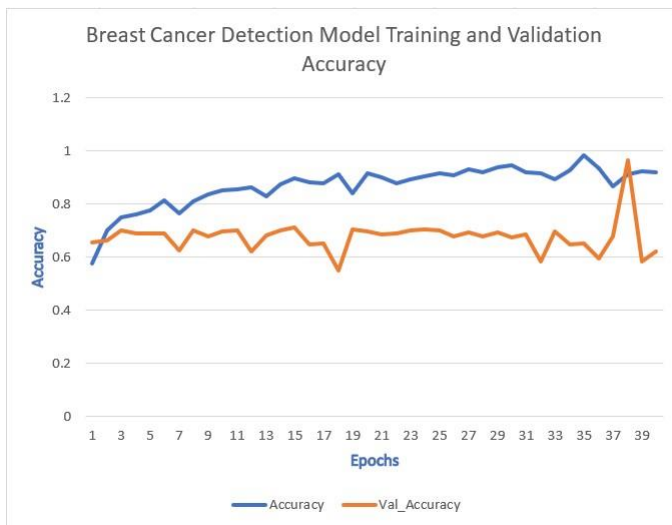


Fig. 2.5. Breast Cancer Detection Model Training and Validation Accuracy Graph

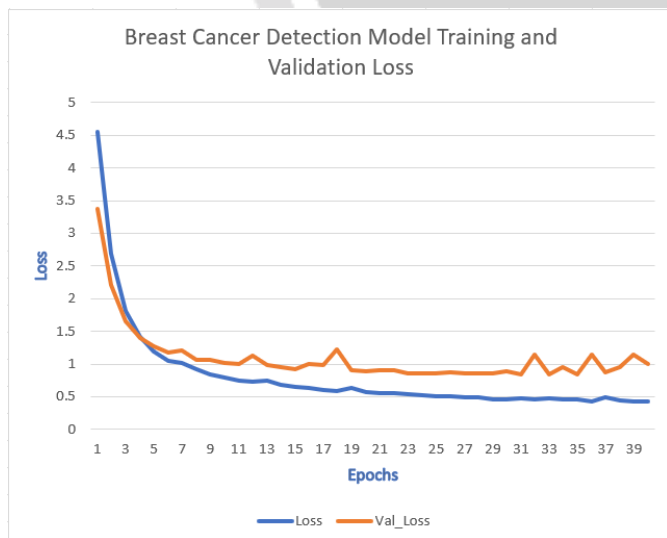


Fig. 2.6. Breast Cancer Detection Model Training and Validation Loss Graph

Figures 2.5 and 2.6 depict the training/validation accuracy and loss, respectively, of the breast cancer detection model. These graphs offer visual representations of the model’s performance during training, allowing for an understanding of its learning trends and capacity to generalize to new data.

In our implementation, the breast cancer detection model achieves an impressive accuracy of 93% on the test dataset, indicating

its ability to correctly classify the majority of cases. This high accuracy level underscores the effectiveness of the chosen model architecture, data pre-processing techniques, and training strategies.

III. RESULTS

After completing the model training phase, our lung cancer detection model demonstrated promising performance across various evaluation metrics. The model achieved a training accuracy of approximately 95% and a training loss of approximately 0.0296 after training for 12 epochs. This high training accuracy indicates that the model effectively learned from the training data and successfully captured the underlying patterns associated with lung cancer detection. During validation, the model attained an accuracy of around 90.54% with a validation loss of about 0.3688. While the validation accuracy is slightly lower than the training accuracy, it still demonstrates the model's ability to generalize to unseen data. The validation loss, although higher than the training loss, remains relatively low, indicating that the model did not overfit excessively during training. These results highlight the efficacy of the chosen model architecture, data pre-processing techniques, and training strategies in developing a lung cancer detection model with high accuracy and generalization ability.

Upon evaluation, our breast cancer detection model showcased impressive performance, demonstrating its capability to accurately classify mammography images into malignant and benign cases. The model achieved an accuracy of 93% on the independent test dataset, indicating its proficiency in correctly identifying the majority of breast cancer cases. This high accuracy level underscores the effectiveness of the VGG16 convolutional neural network architecture, data pre-processing methods, and training strategies employed in the model development process. Overall, the results validate the robustness and reliability of the breast cancer detection model, positioning it as a promising tool for aiding in early breast cancer diagnosis and improving patient outcomes.

IV. FUTURE ENHANCEMENTS

1. Integration of Advanced CNN Architectures: Explore the integration of more advanced convolutional neural network (CNN) architectures, such as DenseNet, ResNet, or EfficientNet, for both lung and breast cancer detection models. These architectures may offer improved performance and feature extraction capabilities, potentially enhancing the models' accuracy and generalization ability.

2. Transfer Learning with Larger Datasets: Utilize transfer learning with larger and more diverse datasets to fine-tune pre-trained models on a broader range of features relevant to lung and breast cancer detection. Fine-tuning pre-trained models on domain-specific data can expedite model convergence and potentially yield better performance.

3. Incorporation of Clinical Data: Integrate additional clinical data, such as patient demographics, medical history, and biomarker information, into the model training process. By incorporating relevant clinical features, the models may gain deeper insights into disease characteristics and improve their predictive accuracy.

4. Real-Time Deployment: Develop strategies for real-time deployment of the models in clinical settings, ensuring efficient inference on new patient data. Optimization techniques, such as model quantization and pruning, can help reduce model size and inference latency, facilitating seamless integration into existing healthcare systems.

5. Continuous Model Monitoring and Updating: Implement mechanisms for continuous model monitoring and updating to adapt to evolving healthcare data and clinical practices. Regular retraining of the models with updated datasets can help maintain their performance over time and ensure their relevance in clinical decision-making.

V. CONCLUSION:

In conclusion, this paper has presented a comprehensive methodology for the early detection of lung and breast cancer using artificial intelligence techniques, specifically deep learning and Convolutional Neural Networks (CNNs). Lung and breast cancer are significant global health challenges, and

improving early detection methods is crucial for enhancing survival rates and patient outcomes.

Our methodology navigates through key stages, including dataset collection, data pre-processing, model training, and evaluation, leveraging TensorFlow and GPU acceleration for efficient model development. For lung cancer detection, we utilized Chest CT-Scan images and developed a robust EfficientNetB0 CNN architecture, achieving high accuracy and generalization ability. Similarly, for breast cancer detection, mammography images were employed, and a VGG16 CNN architecture tailored for binary classification demonstrated impressive accuracy on the test dataset.

The results obtained from our models underscore their efficacy in accurately identifying cancerous conditions, showcasing promising performance across various evaluation metrics. Moreover, our discussion on future enhancements highlights opportunities to further improve the models' accuracy, efficiency, and applicability in clinical settings. Integration of advanced CNN architectures, transfer learning with larger datasets, incorporation of clinical data, real-time deployment strategies, and continuous model monitoring and updating are among the avenues for future research and development.

Overall, our methodology and results contribute to advancing the field of cancer detection through AI-driven approaches, offering promising prospects for early diagnosis and improved patient care. By harnessing the power of artificial intelligence and medical imaging data, we aim to make significant strides towards combating cancer and ultimately enhancing public health on a global scale.

REFERENCES

- 1) Munir, K.; Elahi, H.; Ayub, A.; Frezza, F.; Rizzi, A. Cancer Diagnosis Using Deep Learning: A Bibliographic Review. *Cancers* 2019, 11, 1235, <https://doi.org/10.3390/cancers11091235>.
- 2) Tanzila Saba, Recent advancement in cancer detection using machine learning: Systematic survey of decades, comparisons and challenges, *Journal of Infection and Public Health*, Volume 13, Issue 9, 2020, Pages 1274-1289, ISSN 1876-0341, <https://doi.org/10.1016/j.jiph.2020.06.033>.
- 3) Asuntha, A., Srinivasan, A. Deep learning for lung Cancer detection and classification. *Multimed Tools Appl* 79, 7731–7762 (2020), <https://doi.org/10.1007/s11042-019-08394-3>.
- 4) Suren Makaju, P.W.C. Prasad, Abeer Alsadoon, A.K. Singh, A. Elchouemi, Lung Cancer Detection using CT Scan Images, *Procedia Computer Science*, Volume 125, 2018, Pages 107-114, ISSN 1877-0509, <https://doi.org/10.1016/j.procs.2017.12.016>.
- 5) Thakur, S.K., Singh, D.P. & Choudhary, J. Lung cancer identification: a review on detection and classification. *Cancer Metastasis Rev* 39, 989–998(2020), <https://doi.org/10.1007/s10555-020-09901-x>.
- 6) Anum Masood, Bin Sheng, Ping Li, Xuhong Hou, Xiaoe Wei, Jing Qin, Dagan Feng, Computer- Assisted Decision Support System in Pulmonary Cancer detection and stage classification on CT images, *Journal of Biomedical Informatics*, Volume 79, 2018, Pages 117-128, ISSN 1532-0464, <https://doi.org/10.1016/j.jbi.2018.01.005>.
- 7) Mambou, S.J.; Maresova, P.; Krejcar, O.; Selamat, A.; Kuca, K. Breast Cancer Detection Using Infrared Thermal Imaging and a Deep Learning Model. *Sensors* 2018, 18, 2799, <https://doi.org/10.3390/s18092799>.
- 8) Selvathi, D., Aarthi Poornila, A. (2018). Deep Learning Techniques for Breast Cancer Detection Using Medical Image Analysis. In: Hemanth, J., Balas, V. (eds) *Biologically Rationalized Computing Techniques for Image Processing Applications*. Lecture Notes in Computational Vision and Biomechanics, vol 25, https://doi.org/10.1007/978-3-319-61316-1_8.
- 9) Bhise, S., Gadekar, S., Gaur, A.S., Bepari, S. and Deepmala Kale, D.S.A., 2021. Breast cancer detection using machine learning techniques. *Int. J. Eng. Res. Technol*, 10(7), pp.2278-0181.