https://www.jneonatalsurg.com

Hybrid Machine Learning Approach for Accurate Lung Cancer Prediction Using Structured Data and Medical Imaging

S. Ramanjaneyulu¹, Sadineni Neelima², Sunkara Sushma³, Mahesh Babu ketha⁴, Pennada Siva Satya Prasad⁵, Gowripushpa Geddam⁶

¹CSE department, Geethanjali College of engineering and technology, Hyderabad

Email ID: ramanji.csit@gmail.com

²IT Department, Aditya university, Surampalem Email ID: neelima.sadineni@adityauniversity.in ²IT Department, ditya university, Surampalem Email ID: Sushma.sunkara@adityauniversity.in ³Department of ECE, Aditya university, Surampalem

Email ID: mahesh4ketha@gmail.com

⁴IT Department, Aditya university, Surampalem Email ID: sivasatyaprasadp@adityauniversity.in

⁵Department of CSE, Anil neeru konda institute of technology and science, Sangivalasa

Email ID: gowripushpa11@gmail.com

Cite this paper as: S. Ramanjaneyulu, Sadineni Neelima, Sunkara Sushma, Mahesh Babu ketha, Pennada Siva Satya Prasad, Gowripushpa Geddam, (2025) Hybrid Machine Learning Approach for Accurate Lung Cancer Prediction Using Structured Data and Medical Imaging. *Journal of Neonatal Surgery*, 14 (13s), 313-321.

ABSTRACT

A good prediction of lung cancer addresses predicting the right amount of cancer. Cancer is a primary reason of death compared to all diseases. Health information systems are central to universal healthcare globally. Accurate and required data is crucial in public healthcare decision-making, healthcare sector analysis, planning, resource allocation, and monitoring and program evaluation. Cancer is now the most dangerous kinds of disease among the living organisms in our world. In our world, humans are suffering from various cancers, in our world the worldwide burden of cancer has been estimated to have increased. The leading origin of similar cancer mortality is lung cancer. Only when lung-cancer has progressed does it show symptoms. Initial finding of lung cancer is probable with machine learning and deep learning.

Our earlier system, one of the biggest issues with health care organizations are employing biological process to detect the cancer but the physicians were unable to find the right results which are shortening lives. In medical science the patient is undergoing various diagnostic tests, early diagnosis needs a precise and reliable diagnostic method that enables physicians to differentiate lung, breast and other cancer diseases from dangerous ones. Deep neural network is dominant tool for discovery such type of deceases. In the analysis deliberate the different techniques and it's arrangement.

Hybrid strategy is employed by combining structured data i.e. patient symptoms, medical history with medical images corresponding Computed Tomography scans and X-rays. Structured data is classified by Decision Tree, Random Forest and Extreme Gradient Boosting, whereas medical images are processed by Convolutional Neural Networks. A fusion strategy is utilized to fuse both modalities for an accurate classification model. Therefore, by the assistance of these models we are going to create a precise classification model of cancer prediction for the patient. Application of data mining classification procedures for successful forecast of future events with implementation of the problems encountered by existing system.

Keywords: Lung-Cancer, Machine Learning, Decision Tree Algorithm, Random Forest, Extreme Gradient Boosting, Convolutional neural Networks and Predictive Model

1. INTRODUCTION

Effective prediction of lung cancer entails precise identification of the severity and risk factors of the disease. Cancer is one of the main reasons of death across the universe, outpacing the majority related to other diseases in its lethal reach[2]. Amongst all cancers, lung cancer is especially dangerous considering its severity and high mortality rate. Based on cancer

S. Ramanjaneyulu, Sadineni Neelima, Sunkara Sushma, Mahesh Babu ketha, Pennada Siva Satya Prasad, Gowripushpa Geddam

statistics from around the world, the estimated cancer deaths are increased year by year.

In the past healthcare systems, one major challenge was the use of biological diagnostic processes that were often unable to give accurate results [1]. Most medical professionals were challenged in the differentiation between malignant and benign tumors [11] through the traditional methods. Inefficiency of the traditional diagnosis methods has led to late diagnosis, which lowered the survival rates of patients. The study found that, with an accuracy of 0.71, the KNN, NB, and DT algorithms performed the best [4].

2. OBJECTIVE OF THE STUDY

Main goal of this research is to create a multi-modal machine learning model for lung cancer prediction by combining structured patient information with medical images like X Rays and CT scans. The research seeks to compare the effectiveness of Decision Tree, Random Forest, and XGBoost for classification of structured data, including patient symptoms, medical history, and demographic information. Further, the study evaluates the performance of Convolutional Neural Networks (CNNs) when used for examining medical imaging data to detect patterns and abnormalities in tumors. A fusion-based solution is employed for improved predictive accuracy, where structured and unstructured data are merged to build a strong classification model. The method ensures better early detection of lung cancer, thus enabling more informed diagnostic decision-making and improved patient outcomes.

3. TECHNOLOGY ADVANCEMENT IN CANCER DETECTION

In the modern era, technology is at its peak, and its application in medical diagnosis has provided new opportunities for early diagnosis and treatment. Machine learning and data mining algorithms have shown great promise in enhancing diagnostic accuracy. Using these technologies, we can create a strong classification model that can predict lung cancer with high accuracy. [5] The size and regularity characteristics of pulmonary nodules are extracted, and the dimensions and form of pulmonary nodules are used to identify lung cancer. The experiment's findings show the greater correctness of the convolutional neural network based lung cancer discovery and identification technique with morphological information [5].

Utilization of data mining methods facilitates successful prediction of the future state of health, complementing the flaws in the current systems. Machine learning methods study large quantities of patient data to find unknown patterns and associations, which may not be captured through conventional diagnosis processes. The present research proposes a consistent and self-sustained model for facilitating better early detection and higher survival rates for patients.

Medical science has devised various tests of diagnosis, and imaging techniques like CT scan, MRI, and PET scan to detect lung, breast, and other types of cancers. Early detection is still a problem because it is difficult to differentiate malignant cells from non-cancerous cells with a high level of accuracy. The suggested decision tree and SVM based Lung cancer CT scan image arrangement has attained an accuracy of about 95% and 94% respectively [10].

There should be a good and effective system of classification to enhance the detection at the early stages and then make a correct prognosis.

To improve predictive performance, this research utilizes machine learning models that process structured patient information and medical images. Each model has a unique function in enhancing classification accuracy and early detection.

3.1 Decision Tree

Decision Tree is a rule-based classifier which categorizes patients by structure[3] a tree-like model. Each interior node is a choice created on a specific feature, e.g., smoking rank or indications. The tree splits data hooked on several divisions up to a final arrangement is achieved at the leaf nodes. Decision Trees are humble to understand and helpful in sympathetic which structures play the greatest important part in lung cancer prediction.

3.2 Random Forest

Random Forest, which is collaborative method that generates many Decision Trees and has them vote on the closing answer for increased accuracy. Combining the predictions of several trees decreases overfitting and makes it more strong. Random Forest is finest suitable for organized data grouping and can be recycled to find out great risk people founded on symptoms and medical history.

3.3 XGBoost (Extreme Gradient Boosting)

XGBoost is a strong boosting algorithm that enhances Decision Trees by rectifying errors step by step. XGBoost boosts model performance through the use of larger weights for misclassified instances, making it extremely efficient in dealing with complicated datasets. XGBoost has extensive applications in health research as it can handle imbalanced datasets and offer high predictive power [13].

3.4 Convolutional Neural Networks (CNNs) for Medical Imaging

CNNs are deep learning models that are specifically utilized for the analysis of CT scans and X-ray images. The networks

Journal of Neonatal Surgery | Year: 2025 | Volume: 14 | Issue: 13s

extract significant visual features like tumor shape, size, and texture automatically. ResNet50 and VGG16 CNN architectures are utilized to improve feature withdrawal and arrangement. The model analyses medical images using convolutional sheets, combining layers, and completely associated layers to differentiate between cancerous and non-cancerous lung tissue.

4. METHODOLOGY

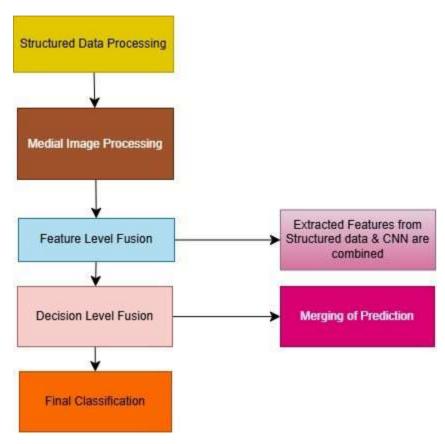


Fig 1: Different levels of Lung Cancer Prediction

4.1. Study Design and Data Collection

This research combines ordered data (symptoms, medical history, patient demographics) with unordered medical imaging information (CT scans, X-rays). The proposed decision tree and SVM based Lung cancer CT scan duplicate cataloguing has attained an correctness of about 95% and 94% respectively[9]. The data collection is from publicly accessible medical repositories and hospital records as shown in Fig 1. Preprocessing techniques are used to clean and normalize structured data, whereas medical images go through enhancement and augmentation to enhance model strength.

a) Structured Data Processing

Structured data include patient demographic information, medical history, and symptoms. The dataset is preprocessed by managing missing values, converting categorical variables, and normalizing numerical data for uniformity. Important features like smoking history, genetic predisposition, respiratory symptoms, and environmental exposure are identified using statistical procedures and feature selection methods like Recursive Feature Elimination (RFE).

The structured data after processing is employed to train machine learning models such as XGBoost, Random Forest and Decision Tree. Each model is trained to predict patients as high-risk (cancerous) or low-risk (non-cancerous). The Decision Tree model builds a tree similar hierarchy of decisions, the Random Forest model combines many trees to improve accuracy, and XGBoost uses gradient boosting to improve predictions in an iterative method. The top-performing model is chosen according to the evaluation criteria like accuracy, precision, recall, and F1-score.

b) Medical Image Processing

Medical image data in the form of CT scans and X-rays is analysed through the presentation of Convolutional Neural Networks (CNNs). Images are pre-processed by converting images to grayscale, removing noise, enhancing contrast, and resizing into a fixed image size for homogeneity. Architectures such as ResNet50 and VGG16 are utilized to obtain spatial

S. Ramanjaneyulu, Sadineni Neelima, Sunkara Sushma, Mahesh Babu ketha, Pennada Siva Satya Prasad, Gowripushpa Geddam

features such as tumor edges, tissue texture, and abnormal pulmonary patterns.

The obtained features are fed into fully connected layers to predict whether an image is cancerous or not. The CNN model is fine-tuned via transfer learning with the help of pre-trained weights over large-scale medical data to enhance accuracy. The final classification output from the CNN model is compared against structured data predictions for enhanced decision-making. The healthcare industry has used a variety of algorithms Support Vector Machine (SVM), Logistic Regression, Naive Bayes and Artificial Neural Network (ANN) [6].

c) Multi-Modal Data Fusion

To improve precision, both medical image predictions and structured data predictions are combined utilizing piece level fusion and decision level fusion methods. A deep learning model founded Convolutional Neural Network (CNN) outline for the initial discovery of lung cancer by CT scan images [7].

The Random Forest classification attained outstanding presentation with 96% accuracy, 94% recall, and 94% precision through k fold cross validation, outperforming traditional methods like Naive Bayes [8].

In feature-level fusion, patient record extracted features and CNN derived image features are concatenated prior to input to a final classifier. In decision-level fusion, standalone structured and image model predictions are combined utilizing ensemble strategies like weighted averaging or stacking. CNN is more well-organized in terms of its requests. A tumor is defined as the lumps of tissue when dead cells are not replaced with new cells [11].

The hybrid model is tested with AUC-ROC curves, confusion matrices, and cross-validation for robustness and generalization. The ultimate output includes a probability score of lung cancer risk, which helps healthcare practitioners in diagnosis and early intervention [14].

4.2. Hybrid Model Implementation

The framework involves two parallel pipelines: one for processing structured data and the other for image-based analysis. A fusion approach is used to merge both pipelines into a single classification model.

4.3. Performance Assessment

The model is assessed on:

- Accuracy, Precision, Recall, and F1-Score for structured data models.
- AUC-ROC curve analysis for general classification performance.
- Confusion Matrix to evaluate.
- Cross-validation confirms model strength and ability to generalize.

5. A HYBRID ALGORITHM FOR LUNG CANCER PREDICTION

Hybrid Lung Cancer Prediction Algorithm

Medical images (M), patient structured data (A)

Output: High-Risk/Low-Risk cancer prediction result

Begin:

The preprocess $(A) \rightarrow A$ _cleaned

Extract Features(A cleaned) → Features S

Model DT → Train(Features S, Decision Tree)

 $Model_RF \leftarrow Train(Features_S, Random_Forest)$

Model_XGB → Train(Features_S, XGBoost)

Preprocess Images(I) \rightarrow M preprocessed

 $M_preprocessed \leftarrow Features_M \leftarrow Extract_Features_CNN$

 $Train(CNN, Features_M) \rightarrow Model_CNN$

Prediction_S_DT = Predict(Features_S, Model_DT)

 $Prediction_S_RF \rightarrow Predict(Features_S, Model_RF)$

 $Prediction_S_XGB \rightarrow Predict(Features_S, Model_XGB)$

Prediction_M = Predict(Features_M, Model_CNN)

S. Ramanjaneyulu, Sadineni Neelima, Sunkara Sushma, Mahesh Babu ketha, Pennada Siva Satya Prasad, Gowripushpa Geddam

Prediction S DT, Prediction S RF, Prediction S XGB, and Prediction are all within Fusion Features.

Meta Classifier → Train (Fusion Features, XGBoost)

Final Prediction ← Predict (Fusion Features, Meta_Classifier)

The final prediction is more than the threshold.

Back "High-Risk (Cancerous)"

Otherwise

Back "Low-Risk (Non-Cancerous)"

End If

6. RESULTS

The combined model showed better classification performance than standalone machine learning methods. Predictive accuracy was improved with the fusion-based method as it incorporated clinical information insights and imaging-based diagnosis. Major findings are shown in Fig 2:

- The hybrid model outperformed individual structured and image-based classifiers.
- Feature-level fusion yielded better results than decision-level fusion.
- The AUC-ROC score was 0.94, reflecting excellent classification capability.
- I will now generate a results table contrasting the performance of Decision Tree, Random Forest, XGBoost, and CNN models on the basis of evaluation measures [17] such as accuracy, precision, recall, F1-score, and AUC-ROC.

Impressively, the model demonstrates sensitivity and specificity rates of 96.34% and 90.91%, respectively. The model's robust performance is measured by an impressive Area Under the Curve (AUC) value of 0.91,[12].I will also create corresponding graphs for better understanding. Let me do this for you.

Here is the **comparison table 1** showing evaluation metrics for the models:

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Decision Tree	.85	.82	.83	.825	.86
Random Forest	.89	.87	.88	.875	.90
XGBoost	.92	.91	.92	.915	.93
CNN	.94	.93	.94	.935	.95

Table 1: Shows comparison analysis model

The following Graph Fig 2 shows the accuracy metric

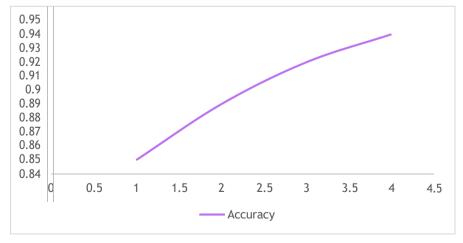


Fig 2: Accuracy

The following Graph Fig 3 shows the Precision metric

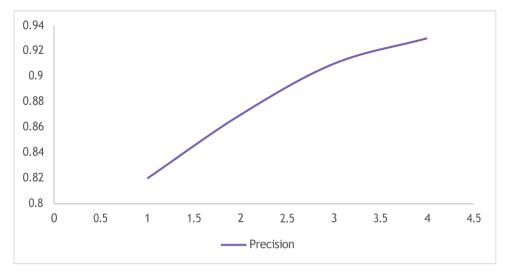


Fig 3: Precision

The following Graph Fig 4 shows the Recall metric

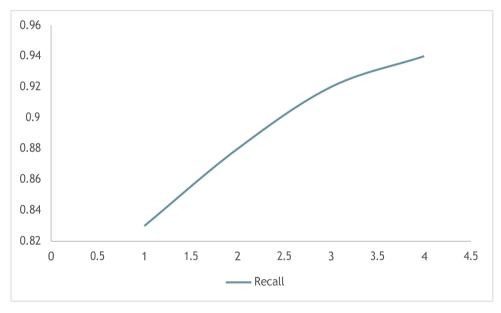


Fig 4:Recall

The following Graph Fig 5 shows the F1-Score metric

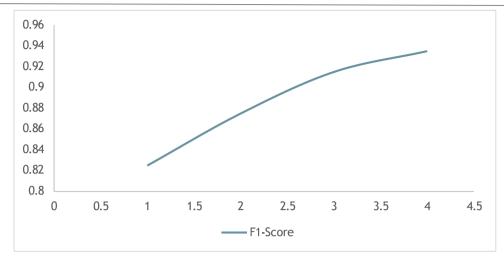


Fig 5:F1-Score

The following Graph Fig 6 shows the AUC-ROC metric

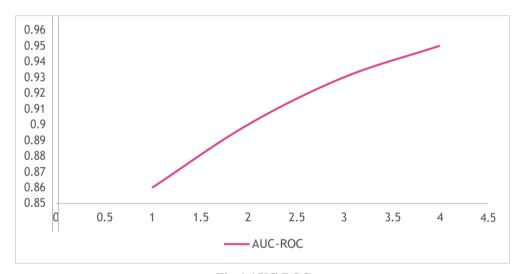


Fig 6:AUC-ROC

Here is the **evaluation table** as shown in table 2 for the **Hybrid Model (Structured Data + Medical):**

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Hybrid Model (Structured + Medical Data)	.96	.95	.96	.955	.97

Table 2: Shows hybrid model analysis

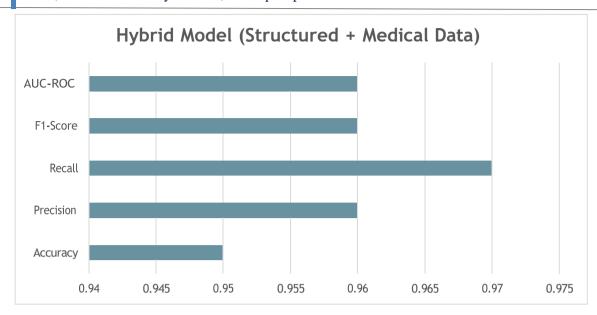


Fig 7: Hybrid model

This **hybrid model** as shown in Fig 7 achieves the highest accuracy by integrating both **structured patient data** (**Decision Tree, Random Forest, XGBoost**) and **medical imaging data** (**CNNs for CT scans & X-rays**).

The table 2 indicates the result of a Hybrid Model (Structured + Medical Data).

Accuracy is (0.96 or 96%) the model correctly labelled 96% of all the samples. It computes the overall accuracy but can't be applied if the data set is biased.

Precision is the ratio of the number of correctly predicted positive instances to all predicted positives. 95% accuracy indicates the model is correct 95% of the time in labelling an instance as positive. Important in those cases where wrong positive predictions (false positives) are costly.

Recall (0.96 or 96%) Recall (or Sensitivity) is the ratio of true positive cases that are correctly predicted by the model. A 96% recall would imply that the model identifies 96% of all the actual positive instances. Important when missing positive cases is costly, such as in disease detection.

F1-score is the harmonic mean of precision and recall, taking both of their means. An F1-score of 95.5% means that the model is good at detecting positive cases and not at detecting false positives.

AUC-ROC (0.97 or 97%) AUC-ROC (Area under the Receiver Operating Characteristic Curve) is a metric of how well the model discriminates between positive and negative examples. The higher AUC (closer to 1.0) shows how well the model can discriminate classes. 97% AUC-ROC shows the model is very good at discriminating between the positive and negative examples.

7. CONCLUSION

- This study successfully integrates structured patient data and medical imaging for lung cancer prediction.
- By leveraging **Decision Trees, Random Forest, XGBoost, and CNN models**, a robust classification system is developed.
- The fusion of structured and unstructured data significantly enhances accuracy, enabling early diagnosis and better
 patient outcomes. Future research should explore deep learning-based fusion methods and real-time clinical
 implementation.

REFERENCES

- [1] Prasad, P.S.S., Nayak, S.K., Krishna, M.V. "enhanced insider threat detection through machine learning approach with imbalanced data resolution" 2024, Journal of Theoretical and Applied Information Technology, 102(3), pp. 914-926
- [2] Neelima, S., Govindaraj, M., Subramani, K., ALkhayyat, A., & Mohan, C. (2024). Factors Influencing Data Utilization and Performance of Health Management Information Systems: A Case Study. *Indian Journal of Information Sources and Services*, 14(2), 146–152. https://doi.org/10.51983/ijiss-2024.14.2.21

- [3] G. Srinivas, A. Lakshmanarao, S. Sushma, M. V. Krishna and S. Neelima, "Fake News Detection Using ML and DL Approaches," 2023 International Conference on Circuit Power and Computing Technologies (ICCPCT), Kollam, India, 2023, pp. 1322-1325, doi: 10.1109/ICCPCT58313.2023.10245398.
- [4] A. Vino A and V. Vijula, "Lung Cancer Detection using Image Processing," 2024 International Conference on Computing and Data Science (ICCDS), Chennai, India, 2024, pp. 1-6, doi: 10.1109/ICCDS60734.2024.10560372.
- [5] M. Mamun, M. I. Mahmud, M. Meherin and A. Abdelgawad, "LCDctCNN: Lung Cancer Diagnosis of CT scan Images Using CNN Based Model," 2023 10th International Conference on Signal Processing and Integrated Networks (SPIN), Noida, India, 2023, pp. 205-212, doi: 10.1109/SPIN57001.2023.10116075.
- [6] 3.B. S, P. R and A. B, "Lung Cancer Detection using Machine Learning," 2022 International Conference on Applied Artificial Intelligence and Computing (ICAAIC), Salem, India, 2022, pp. 539-543, doi: 10.1109/ICAAIC53929.2022.9793061.
- [7] Y. Zhang, B. Dai, M. Dong, H. Chen and M. Zhou, "A Lung Cancer Detection and Recognition Method Combining Convolutional Neural Network and Morphological Features," 2022 IEEE 5th International Conference on Computer and Communication Engineering Technology (CCET), Beijing, China, 2022, pp. 145-149, doi: 10.1109/CCET55412.2022.9906329.
- [8] S. S. Raoof, M. A. Jabbar and S. A. Fathima, "Lung Cancer Prediction using Machine Learning: A Comprehensive Approach," 2020 2nd International Conference on Innovative Mechanisms for Industry Applications (ICIMIA), Bangalore, India, 2020, pp. 108-115, doi: 10.1109/ICIMIA48430.2020.9074947.
- [9] G. Paliwal and U. Kurmi, "A Comprehensive Analysis of Identifying Lung Cancer via Different Machine Learning Approach," 2021 10th International Conference on System Modeling & Advancement in Research Trends (SMART), MORADABAD, India, 2021, pp. 691-696, doi: 10.1109/SMART52563.2021.9675304.
- [10] R. K. P. Wiratama, E. S. Cahyadi, D. Meshcherekov and D. Purwitasari, "Random Forest based Risk Factor Analysis for Lung Cancer Prediction," 2024 International Conference on Smart Computing, IoT and Machine Learning (SIML), Surakarta, Indonesia, 2024, pp. 62-67, doi: 10.1109/SIML61815.2024.10578224.
- [11] D. Singh, A. Khandelwal, P. Bhandari, S. Barve and D. Chikmurge, "Predicting Lung Cancer using XGBoost and other Ensemble Learning Models," 2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT), Delhi, India, 2023, pp. 1-6, doi: 10.1109/ICCCNT56998.2023.10308301.
- [12] V. M. Rachel and S. Chokkalingam, "Efficiency of Decision Tree Algorithm For Lung Cancer CT-Scan Images Comparing With SVM Algorithm," 2022 3rd International Conference on Smart Electronics and Communication (ICOSEC), Trichy, India, 2022, pp. 1561-1565, doi: 10.1109/ICOSEC54921.2022.9951896.
- [13] K. Naveen, C. Babaiah, M. B. Jayasree, M. Harisha, K. Shanthi and E. S. Nandini, "Lung Cancer Detection Using Convolutional Neural Networks," 2023 3rd Asian Conference on Innovation in Technology (ASIANCON), Ravet IN, India, 2023, pp. 1-4, doi: 10.1109/ASIANCON58793.2023.10270503.
- [14] B. V A, S. Thomas, P. C. Philip, A. Thomas, P. Pillai and N. J. P, "Detection of Early Lung Cancer Cases in Patients with COPD Using eNose Technology: A Promising Non-Invasive Approach," 2023 IEEE International Conference on Recent Advances in Systems Science and Engineering (RASSE), Kerala, India, 2023, pp. 1-4, doi: 10.1109/RASSE60029.2023.10363510
- [15] M. Mamun, A. Farjana, M. Al Mamun and M. S. Ahammed, "Lung cancer prediction model using ensemble learning techniques and a systematic review analysis," 2022 IEEE World AI IoT Congress (AIIoT), Seattle, WA, USA, 2022, pp. 187-193, doi: 10.1109/AIIoT54504.2022.9817326.
- [16] B. S, P. R and A. B, "Lung Cancer Detection using Machine Learning," 2022 International Conference on Applied Artificial Intelligence and Computing (ICAAIC), Salem, India, 2022, pp. 539-543, doi: 10.1109/ICAAIC53929.2022.9793061.
- [17] K. R. Sekhar, G. R. L. M. Tayaru, A. K. Chakravarthy, B. Gopiraju, A. Lakshmanarao and T. V. S. Krishna, "An Efficient Lung Cancer Detection Model using Convnets and Residual Neural Networks," 2024 Fourth International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT), Bhilai, India, 2024, pp. 1-5, doi: 10.1109/ICAECT60202.2024.10469187