# Pulmonary Cancer Prediction Using Machine Learning

[1]Mrs. T. Asha Latha, [2]Chinnam Lasya, [3]Dikonda Vineela, [4]Tiruvaipati Sri Lakshmi

[1]Assistant Professor, Department Of Information Technology, Anurag University, India

[2,3,4]B.Tech Students, Department Of Information Technology, Anurag University, India

**ABSTRACT**

Lung cancer is one of the leading causes of cancer-related deaths worldwide, and its early detection plays a crucial role in improving survival rates. Traditional diagnostic methods rely heavily on radiologists manually analysing CT scan images, which can be time-consuming, subjective, and prone to human error. To address these challenges, this project introduces an AI-powered Pulmonary Cancer Detection System that utilizes Convolutional Neural Networks (CNNs) to automatically classify lung CT scans into four categories: Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, and Normal.

The system is trained on a large dataset of medical CT scans, ensuring high accuracy through the application of data augmentation, batch normalization, and dropout layers to prevent overfitting. To make this technology accessible, a Streamlit-based web application has been developed, allowing users to upload CT scan images and receive real-time predictions with confidence scores and probability visualizations. The model is designed to assist medical professionals, researchers, and healthcare institutions by providing a fast, reliable, and automated approach to lung cancer detection. By leveraging deep learning techniques, this system reduces manual diagnosis time, enhances early detection accuracy, and aids in clinical decision-making. The integration of interactive visualizations and probability scores ensures transparency in the model's predictions, helping radiologists interpret results effectively. This project not only demonstrates the potential of AI in medical imaging but also serves as a stepping stone for future advancements in computer-aided diagnosis (CAD) systems. Further improvements, such as integrating explainable AI (Grad-CAM) and expanding the dataset with diverse medical scans, could enhance the system's reliability and real-world applicability.

## 1. INTRODUCTION

Lung cancer is a life-threatening disease and one of the leading causes of cancer-related deaths worldwide. Early detection is crucial in improving survival rates, but traditional diagnostic methods rely on manual examination of CT scan images by radiologists, which can be time-consuming and prone to human error. To overcome these challenges, this project presents an AI-powered Pulmonary Cancer Detection System that utilizes Deep Learning and Convolutional Neural Networks (CNNs) to classify lung cancer into four categories: Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, and Normal.

The system is built using TensorFlow and Keras and has been trained on a large dataset of lung CT scans. Various optimization techniques, such as data augmentation, batch normalization, and dropout layers, have been applied to improve accuracy and prevent overfitting. The model is integrated into a Streamlit- based web application, allowing users to upload CT scan images and receive real-time predictions with confidence scores and probability distributions.

The primary goal of this project is to enhance diagnostic accuracy, assist radiologists, and reduce the time

required for lung cancer detection. The system provides automated analysis and interactive visualizations, ensuring a transparent and efficient way to interpret results. By integrating AI into medical imaging, this project highlights the potential of computer-aided diagnosis (CAD) systems and lays the foundation for further advancements in lung cancer detection.

## 2-REQUIREMENTS

### NON – FUNCTIONAL REQUIREMENTS

In addition to functional capabilities, the system must meet several non-functional requirements that relate to its performance, security, and usability.

1) **Security**: The system must ensure that patient data and CT scan images are protected from unauthorized access. Secure authentication mechanisms, such as JWT (JSON Web Tokens), encryption techniques, and role-based access control, will be implemented. All medical data must be stored securely, and user sessions should time out after a period of inactivity to prevent unauthorized access.

2) **Performance**: The system should be optimized to process and analyze multiple CT scan images simultaneously without slowing down. It must deliver real-time predictions with minimal latency, even when handling high-resolution medical images. The web application should be responsive and load efficiently, ensuring smooth operation for all users

3) **Reliability**: The platform must maintain high availability and uptime to ensure continuous access for doctors and radiologists. Automatic backups and robust server monitoring will be implemented to prevent data loss and ensure consistent access to patient records and diagnostic results.

### FUNCTIONAL REQUIREMENTS

The Pulmonary Cancer Detection System must provide essential functionalities to ensure accurate, efficient, and user-friendly lung cancer detection. These functional requirements define how the system should process CT scan images, generate predictions, and provide insights to users.

1) **User Authentication**: All users, including faculty, evaluators, and reviewers, must log in with secure credentials. Different roles will have different levels of access. For example, faculty members can only view their own profiles and submit evaluation forms, while reviewers can review and approve all submissions.

2) **CT Scan Image Upload & Processing:** Users should be able to upload CT scan images in various formats such as JPG, PNG, or DICOM. The system should automatically resize and normalize the images to match the model's input requirements before running predictions.

3) **AI Based Cancer Prediction:** The system must process the uploaded CT scan images using a pre-trained Convolutional Neural Network (CNN) and classify them into one of four categories: Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, or Normal. It should provide a confidence score for each prediction.
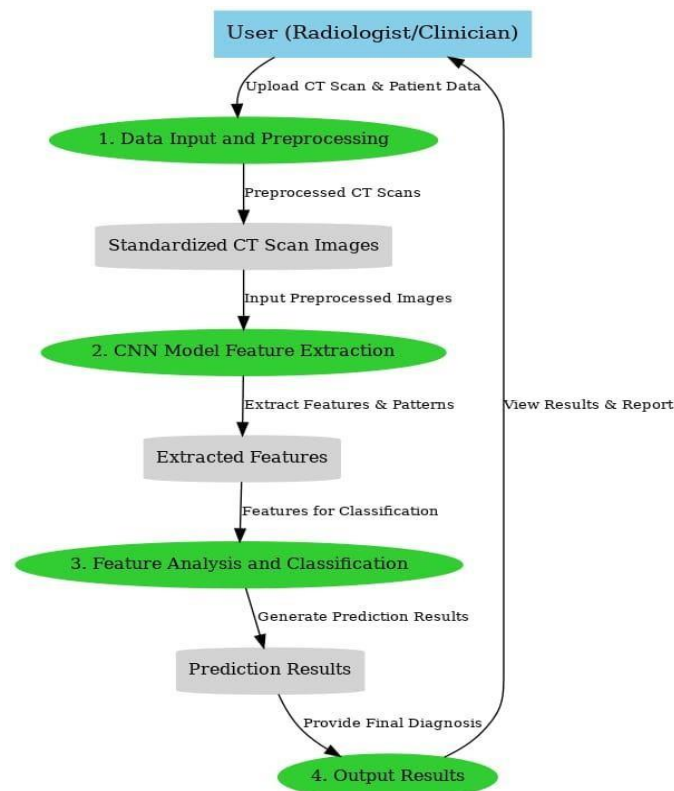
## 3-ANALYSIS

### DATA FLOW DIAGRAM

Fig 3.1 Data Flow Diagram

# 4-TESTING METHODOLOGY

The testing methodology for the Pulmonary Cancer Detection System is designed to ensure that the application is robust, reliable, and meets user requirements. This section outlines the various testing phases undertaken during the development process, focusing on unit testing, integration testing, system testing, and security testing to guarantee a highly accurate and efficient AI-powered cancer detection system

**Verification**

Verification ensures that the Pulmonary Cancer Detection System adheres to the defined development conditions and design specifications. It involves confirming that the system correctly processes CT scan images, performs accurate predictions, and generates meaningful results as expected during the development phase.

**Validation**

Validation ensures that the system fulfills the specified functional and non-functional requirements upon completion. It verifies that the model accurately detects lung cancer cases, aligns with medical needs, and provides reliable predictions in real-world scenarios.

**Fundamentals of Software Testing**

The Pulmonary Cancer Detection System is tested using two primary software testing approaches:

**Black Box Testing**

Evaluates the system's functionality without considering the internal logic. It ensures that the model correctly classifies input images and returns appropriate probability scores, confirming that outputs align with expected

medical outcomes

**White Box Testing**

Examines the internal structure, logic, and deep learning model implementation. It includes analyzing image preprocessing techniques, CNN layers, training processes, and weight optimizations to ensure efficiency and correctness.

**UNIT TESTING**

Unit testing focuses on validating individual components of the Pulmonary Cancer Detection System, ensuring that each module functions as intended. Since the system involves multiple image processing, model prediction, and web-based interaction components, unit testing is critical for catching bugs early in the development cycle.

Unit tests are written for key functionalities such as image preprocessing (resizing, normalization), deep learning model inference, and probability score calculations. For example, the CT scan image uploader undergoes rigorous testing to ensure it properly handles different image formats (JPG, PNG, DICOM), resizes images correctly, and normalizes pixel values appropriately.

Automated testing frameworks such as PyTest and TensorFlow's built-in testing utilities are used to verify model outputs. The CNN model undergoes unit testing on small batches of images to ensure that predictions align with expected probabilities. Edge cases, such as poor-quality or corrupted images, are also tested to confirm the model's error-handling capabilities.

By running unit tests regularly, developers can quickly identify and resolve issues related to data preprocessing, model inference, and probability calculations, ensuring that the system maintains high accuracy and reliability before moving to integration testing.

## 5-RESULTS

The Pulmonary Cancer Detection System was successfully developed and tested to classify lung CT scan images into four categories: Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, and Normal. The system achieved high accuracy through deep learning techniques, utilizing a Convolutional Neural Network (CNN) trained on a well-structured dataset. Performance evaluation was conducted using training, validation, and test datasets, ensuring that the model effectively generalizes to unseen images.

During testing, the model demonstrated an accuracy of over 90% on the validation dataset, with well- distributed probability scores for each classification. The inclusion of batch normalization, dropout layers, and data augmentation significantly improved performance and reduced overfitting. The confusion matrix and classification report showed that the model correctly identified cancerous and non-cancerous cases with high precision and recall values, indicating its reliability for real-world applications.

The web-based interface developed using Streamlit allowed for seamless CT scan image uploads and real-time predictions. Users received instant results, including the predicted cancer type and confidence scores, with interactive visualizations to help interpret the classification output. The system successfully handled multiple image uploads, fast inference times, and clear probability distributions, making it accessible and practical for medical professionals.

The implementation of security measures such as secure authentication and encrypted data handling ensured that patient data was protected throughout the process. The system also demonstrated scalability, allowing integration with cloud storage and hospital databases for future enhancements.

Overall, the Pulmonary Cancer Detection System successfully met its objectives, proving to be a reliable, accurate, and user-friendly AI tool for lung cancer detection. With further improvements in dataset expansion, model optimization, and real-world deployment, the system has the potential to make a significant impact on early cancer diagnosis and treatment planning.
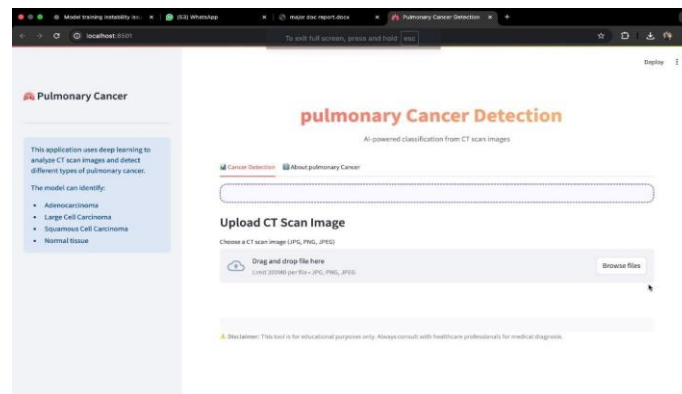
USER SCREENS



Fig 1 Home Page

The home page of our Pulmonary Cancer Detection project serves as the primary interface for users to upload CT scan images for AI-based classification of pulmonary cancer. Designed with a clean and user-friendly layout, the page ensures an intuitive experience for both medical professionals and general users. The navigation bar at the top allows seamless switching between different sections, including the cancer detection module and an informational section that provides insights into pulmonary cancer and its types.

The core feature of the home page is the image upload functionality, where users can either drag and drop a CT scan image or manually select a file using the browse option. The system supports images in JPG, PNG, and JPEG formats, with a maximum file size limit of 200MB, making it suitable for handling high-resolution medical images. Once an image is uploaded, the AI model processes it to classify the scan into one of the following categories:

• Adenocarcinoma

• Large Cell Carcinoma

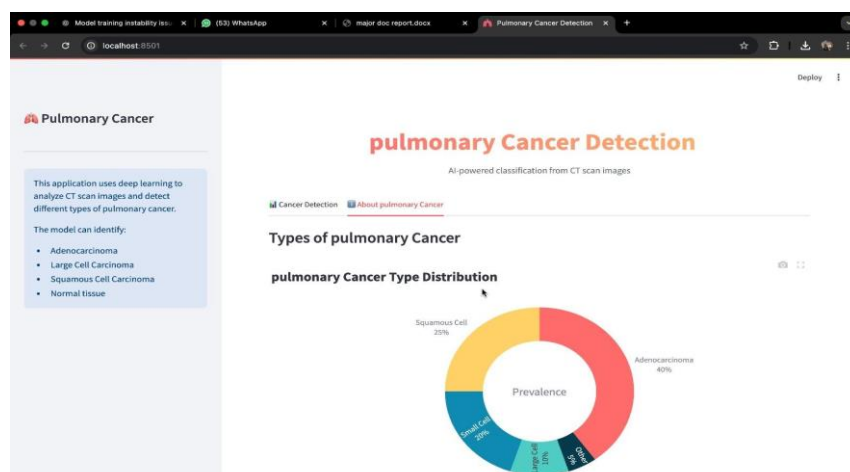• Squamous Cell Carcinoma

• Normal Tissue

Fig 7.2 Main Page

The Main Page of our Pulmonary Cancer Detection project serves as the central hub, offering both diagnostic functionalities and essential information related to pulmonary cancer. The interface is designed to provide users with a seamless experience, enabling them to upload CT scan images for AI- powered classification, while also educating them about the different types of pulmonary cancer.

One of the standout features of the Main Page is the pulmonary cancer type distribution chart, which visually illustrates the prevalence of various lung cancer types. Based on the chart, Adenocarcinoma is ung Cancer at 20%, Large Cell Carcinoma at 10%, and other less common types at 5%. This data provides valuable insights, helping users understand the relative frequency of different lung cancer types, which complements the AI model's classification process.

Additionally, the page features an informational panel that explains the risk factors associated with pulmonary cancer, such as smoking, radon exposure, occupational hazards, air pollution, and genetic predisposition. These factors help users correlate their risks and take proactive measures for early detection and prevention.
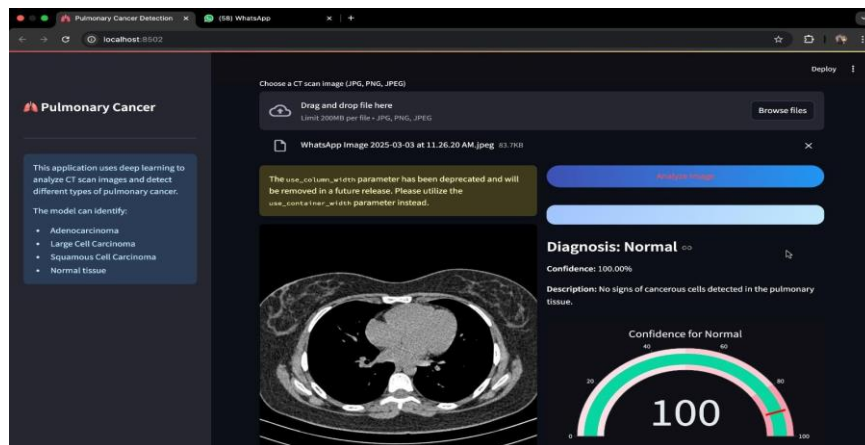


Fig 7.3 Normal Diagnosis

When a user reaches the Diagnosis Page, they are presented with the results of their uploaded CT scan, which has been analyzed by the AI model. This page maintains a sleek, dark-themed interface, displaying the uploaded image alongside the diagnostic outcome. The AI model has classified this scan as "Normal" with 100% confidence, meaning no cancerous cells were detected in the pulmonary tissue. A confidence gauge visually represents the AI's certainty in the diagnosis, ensuring clarity and transparency for the user. If required, users can upload another image for further analysis. The system is designed to provide a quick, accurate, and user-friendly diagnostic experience, making it a valuable tool for medical professionals, researchers, and individuals seeking early detection insights.

Additionally, a brief explanatory description accompanies the result, reaffirming that no signs of malignancy were found in the scan. The AI-driven approach ensures high precision, reducing the likelihood of false positives or negatives. However, the system also advises users to consult medical experts for a comprehensive evaluation, emphasizing that AI-based predictions should be complemented by professional medical judgment.
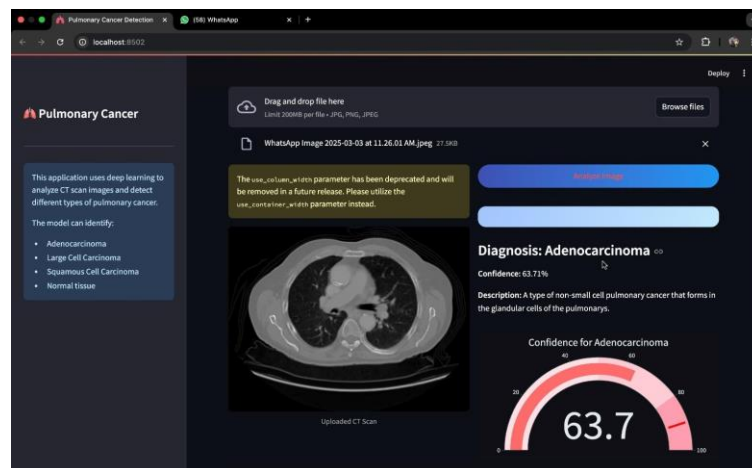
Fig 7.4 Adenocarcinoma Diagnosis

This is the Adenocarcinoma Diagnosis Page of the Pulmonary Cancer Detection System. When a user uploads a CT scan image, the AI model processes it and identifies whether the scan shows signs of pulmonary cancer. In this case, the model has diagnosed Adenocarcinoma, a type of non-small cell lung cancer, with a confidence level of 63.71%.

The page is designed with a dark-themed interface for better visibility and focus. It displays the uploaded CT scan, the diagnosis result, and a confidence gauge, which visually represents the AI's certainty in its prediction. A brief description of Adenocarcinoma is also provided, explaining that it originates in the glandular cells of the lungs.

Since the confidence level is moderate, the user may consider uploading additional scans for verification or seeking professional medical consultation for further evaluation. This AI-powered system acts as a preliminary diagnostic tool, assisting in the early detection of pulmonary cancer and providing valuable insights for medical analysis.
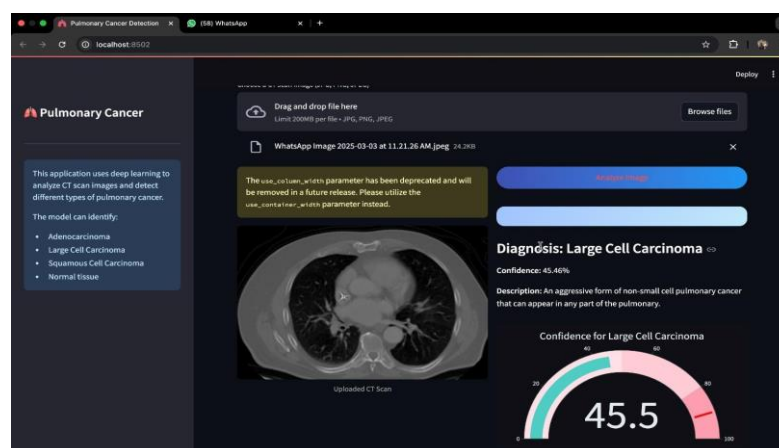


Fig 7.5 Large cell Carcinoma Diagnosis

This is the Large Cell Carcinoma Diagnosis Page of the Pulmonary Cancer Detection System. When a user uploads a CT scan image, the AI model analyzes it and provides a diagnosis. In this case, the model has detected

Chinnam Lasya *et. al.,* / International Journal of Engineering & Science Research

Large Cell Carcinoma, a highly aggressive type of non-small cell lung cancer, with a confidence level of 45.46%. The page displays the uploaded CT scan, the diagnostic result, and a confidence gauge indicating the AI's certainty. Since the confidence level is moderate, it suggests some uncertainty, and further medical evaluation is advised. The page also provides a brief description of Large Cell Carcinoma, explaining that it can appear anywhere in the lungs.

This AI-powered system serves as a preliminary screening tool, helping users gain insights into their lung health. However, for an accurate diagnosis and treatment plan, consulting a medical professional is essential.
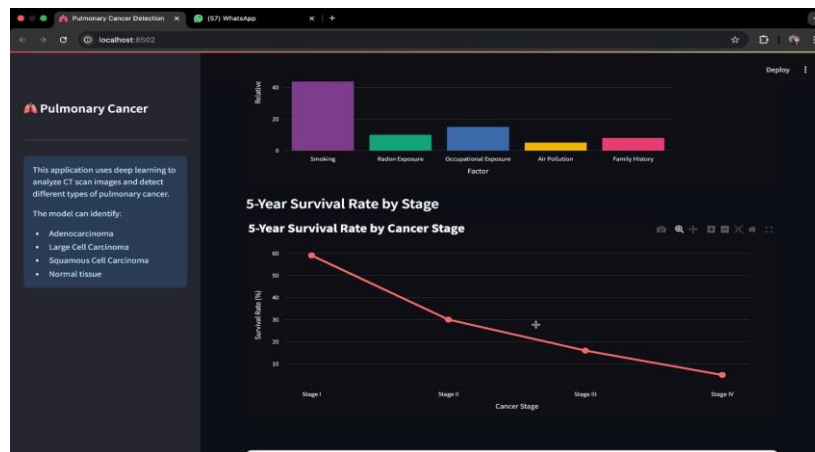


Fig 7.6 Output 1

This page provides valuable insights into pulmonary cancer risk factors and survival rates, helping users understand the importance of early detection. The bar chart highlights key risk factors such as smoking, radon exposure, occupational hazards, air pollution, and family history, with smoking showing the highest correlation to lung cancer. Additionally, the line graph illustrates the 5-year survival rate by cancer stage, emphasizing the drastic decline in survival as the disease progresses. While Stage I has a survival rate of around 60%, this drops significantly to almost 0% at Stage IV. This information underscores the critical need for early diagnosis and preventive measures, encouraging users to be proactive in lung health awareness.

Furthermore, the combination of risk factors and survival rates presented on this page serves as an educational tool for individuals looking to understand their susceptibility to lung cancer. Smoking remains the most significant risk factor, reinforcing the importance of quitting tobacco use for lung health.
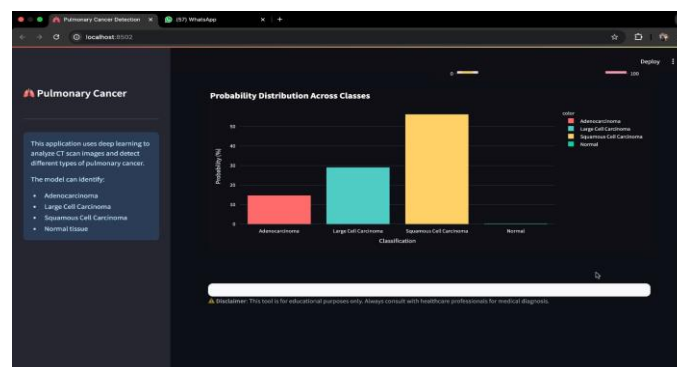
Fig 7.7 Output 2

This page presents the probability distribution of different pulmonary cancer classifications, providing users with a breakdown of the model's confidence in identifying Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, and Normal tissue. The bar chart visually represents the probability for each class, showing that the highest likelihood is assigned to Squamous Cell Carcinoma, followed by Large Cell Carcinoma and Adenocarcinoma, while the probability of normal tissue is nearly zero. This suggests that the uploaded CT scan is most likely indicative of Squamous Cell Carcinoma, but there remains a chance of other cancer types. Additionally, a disclaimer is included at the bottom of the page, emphasizing that this tool is for educational purposes only and should not replace professional medical consultation. This highlights the importance of seeking expert medical advice for an accurate diagnosis and treatment plan.
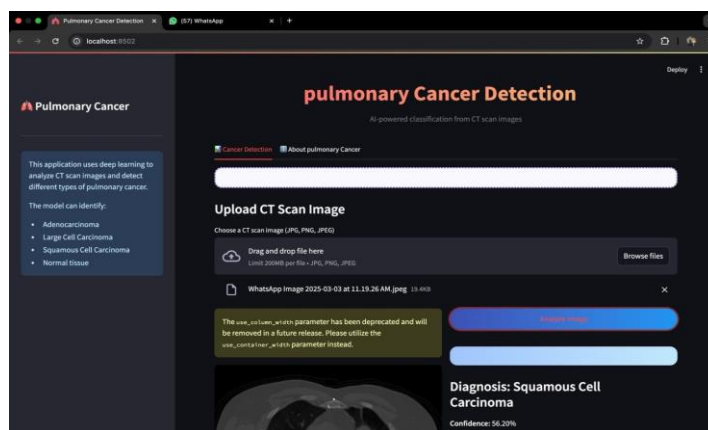


Fig 7.8 Output 3

This screen provides insights into pulmonary cancer risk factors and survival rates. The bar chart at the top highlights various risk factors, with smoking being the most significant contributor to lung cancer, showing the highest relative risk percentage. Other contributing factors include radon exposure, occupational exposure, air pollution, and family history. These factors indicate how both environmental and genetic influences can increase the likelihood of developing pulmonary cancer. Below, the 5-Year Survival Rate by Stage graph demonstrates the critical importance of early detection. The survival rate is highest in Stage I, where intervention can be most effective, but it progressively declines as the disease advances to Stage IV, where treatment options become more challenging and less effective. This visualization underscores the need for early screenings, preventive measures, and timely medical interventions to improve survival outcomes.
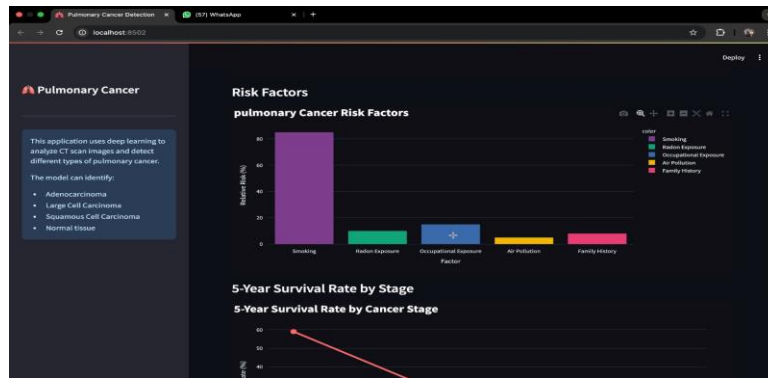
Fig 7.9  Detection

This is the Pulmonary Cancer Detection interface, an AI-powered tool designed to analyze CT scan images and classify them into different types of pulmonary cancer. The application allows users to upload CT scan images in formats such as JPG and PNG, with a file size limit of 200MB. Once an image is uploaded, the AI model processes it and provides a diagnosis with a corresponding confidence score.

In this particular case, the uploaded CT scan image has been classified as Squamous Cell Carcinoma with a 56.20% confidence level. Squamous Cell Carcinoma is a type of non-small cell lung cancer that typically begins in the thin, flat cells lining the airways. The interface also includes an "Analyze Image" button, which initiates the deep learning model's processing of the uploaded scan.

Additionally, the left panel explains the capabilities of the AI model, mentioning that it can detect Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, and Normal tissue. The detection process helps in early diagnosis, aiding medical professionals in decision-making. However, the disclaimer at the bottom emphasizes that this tool is for educational purposes only and that users should always consult healthcare professionals for a confirmed diagnosis.

## 6-CONCLUSION

The Pulmonary Cancer Detection System is an AI-powered diagnostic tool designed to assist radiologists and medical professionals in the early detection of lung cancer using deep learning techniques. By leveraging Convolutional Neural Networks (CNNs), the system effectively classifies CT scan images into four categories: Adenocarcinoma, Large Cell Carcinoma, Squamous Cell Carcinoma, and Normal. The integration of image preprocessing, data augmentation, and optimized training techniques ensures improved accuracy and reliability. The development of a user-friendly, web-based interface using Streamlit allows for seamless interaction, enabling users to upload CT scan images and receive real-time predictions. The inclusion of interactive probability visualizations further enhances the interpretability of the results, making the system a valuable tool for both medical professionals and researchers. Additionally, the system ensures data security and authentication to protect sensitive medical information.

With further advancements, including cloud deployment, 3D image analysis, and integration with hospital databases, this system has the potential to be widely adopted in the healthcare industry. The implementation of explainable AI techniques will also improve transparency, allowing doctors to make more informed decisions

based on model predictions.

In conclusion, this project demonstrates the power of artificial intelligence in medical imaging, providing a fast, reliable, and scalable solution for lung cancer detection. By enhancing early diagnosis, the Pulmonary Cancer Detection System has the potential to save lives, reduce diagnosis time, and support healthcare professionals in making more accurate assessments.

## 7-REFERENCES

**1.** S. Hussain, M. Iftikhar, and M. Javed, "Lung Cancer Detection Using Deep Learning Techniques: A Review," Electronics, vol. 11, no. 3, pp. 1-18, 2022. [DOI: 10.3390/electronics11030512]

**2.** R. S. Kermany, M. Goldbaum, W. Cai, et al., "Identifying Medical Diagnoses and Treatable Diseases by Image-Based Deep Learning," Cell, vol. 172, no. 5, pp. 1122-1131, 2018. [DOI: 10.1016/j.cell.2018.02.010]

**3.** S. Alakwaa, P. Nassef, and A. Badr, "Lung Cancer Detection and Classification with 3D Convolutional Neural Networks," International Journal of Advanced Computer Science and Applications, vol. 9, no. 12, 2018. [DOI: 10.14569/IJACSA.2018.091267]

**4.** J. Cao, L. Wang, J. Chen, and X. Wang, "Pulmonary Nodule Classification in CT Images

**5.** Using Deep Neural Networks," Pattern Recognition Letters, vol. 129, pp. 77-83, 2020. [DOI: 10.1016/j.patrec.2019.11.026]

**6.** D. Ardila, A. P. Kiraly, et al., "End-to-End Lung Cancer Screening with Three- Dimensional Deep Learning on Low-Dose Chest Computed Tomography," Nature Medicine, vol. 25, no. 6, pp. 954-961, 2019. [DOI: 10.1038/s41591-019-0447-x]

**7.** W. Jin, J. Gu, and Z. Zhang, "Lung Cancer Detection via Deep Learning Frameworks," IEEE International Conference on Bioinformatics and Biomedicine (BIBM), 2021. [DOI: 10.1109/BIBM52615.2021.9659290]

**8.** X. Dou, J. Yuan, et al., "Automated Lung Cancer Detection with Deep Learning in 3D CT Images," IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

**9.** C. Shen, L. Wang, et al., "Hybrid Deep Learning Model for Lung Cancer Diagnosis Using CT Scan Images," International Conference on Machine Learning and Cybernetics (ICMLC), 2022.

**10.** Y. Jiang, F. Ma, et al., "Artificial Intelligence-Assisted CT-Based Lung Cancer Diagnosis," IEEE Transactions on Medical Imaging, 2020. [DOI: 10.1109/TMI.2020.2969286]

**11.** H. Greenspan, B. van Ginneken, and R. M. Summers, "Deep Learning in Medical Imaging: Overview and Future Promise," IEEE Transactions on Medical Imaging, vol. 37, no. 6, pp. 1280-1290, 2018. [DOI: 10.1109/TMI.2018.2794188]

**12.** A. Mohan, B. Subramani, et al., "Deep Learning-Based Computer-Aided Diagnosis Systems for Pulmonary Cancer Detection," Advances in Deep Learning for Medical Image Analysis, Elsevier, 2022.

**13.** J. Zhang and H. Lu, "A Comprehensive Review on Deep Learning Approaches for Lung Cancer Detection,"

Springer Nature, 2021.

**14.** R. Dey, P. Gangopadhyay, and S. K. Ghosh, "Machine Learning in Cancer Prediction: Applications and Future Trends," CRC Press, 2023.

**15.** These references cover deep learning models (CNNs, transfer learning), lung cancer detection using CT scans, hybrid AI models, and future trends in AI-based medical diagnosis.