

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

BELAGAVI – 590018



A Major Project Report on

**“Lung Cancer Detection using CT(Computed Tomography)
Image Processing & Machine Learning”**

Submitted By

Mr. Aftab I Yaragatti

USN:2HN22CS005

Miss. Bhagyashree S Poojari

USN:2HN22CS013

Miss. Kavita K Dodagoudanavar

USN:2HN22CS024

Miss. Laxmi A Bilur

USN:2HN22CS026

Under the Guidance of

Prof . M. G. Ganachari



Department of Computer Science and Engineering

S.J.P.N Trust's

HIRASUGAR INSTITUTE OF TECHNOLOGY, NIDASOSHI-591236

Inculcating Values, Promoting Prosperity

Approved by AICTE, Recognized by Govt. of Karnataka, Affiliated to VTU Belagavi.

Recognized under 2(f) & 12B of UGC Act, 1956

Accredited at “A+” Grade by NAAC

Programmes Accredited by NBA: CSE & ECE

Academic Year 2025-2026

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

BELAGAVI – 590018



Department of Computer Science and Engineering

Certificate

Certified that the Machine Learning Major-Project titled “Lung Cancer detection using CT(Computed Tomography) Image Processing and Machine Learning” is a bonafide work carried out by **Mr. AFTAB I YARAGATTI (2HN22CS005)**, **Miss. BHAGYASHREE S POOJARI (2HN22CS013)** , **Miss. KAVITA K DODAGOUDANAVAR (2HN22CS024)**, and **Miss. LAXMI A BILUR (2HN22CS026)** in partial fulfilment of the requirements of **VII semester of Bachelor of Computer Science and Engineering of the Visvesvaraya Technology University**, Belagavi during the academic year 2025 – 2026. It is certified that all the corrections /suggestions indicated have been incorporated in the report. The project has been approved as it satisfies the academic requirements in respect of Major Project work prescribed by the Bachelor of Engineering course.

Prof . M. G. Ganachari
Guide

Dr. S. V. Manjaragi
H.O.D.

Dr. S. C. Kamate
Principal

Name of the Examiners

1. _____

2. _____

Signature

ACKNOWLEDGEMENT

It is our pleasure to acknowledge the help we have received from individuals and the institute. We would like to thank our **Principal Dr. S. C. Kamate** in particular for excellent facilities provided.

We are very thankful and highly obliged to our beloved H.O.D **Dr. S. V. Manjaragi** for his encouragement and insightful comments at virtually all stages of my mini-project.

We also wish to express our warm and grateful thanks to our respected guide **Prof . M. G. Ganachari** for giving us the advice and valuable guidance.

We also thank all Professors and Staff of CSE department for their constant encouragement.

Mr. Aftab I Yaragatti

Miss. Bhagyashree S Poojari

Miss. Kavita K Dodagoudanavar

Miss. Laxmi A Bilur.

ABSTRACT

Lung cancer remains the leading cause of cancer-related mortality worldwide, necessitating the development of advanced diagnostic tools for early detection. This project presents an innovative approach to lung cancer detection by integrating computed tomography (CT) imaging with machine learning (ML) techniques. The methodology encompasses several key stages: Preprocessing of CT images to enhance quality and reduce noise, segmentation to isolate regions of interest, feature extraction to identify relevant patterns, and classification to differentiate between benign and malignant lesions. A convolutional neural network (CNN)-based model is employed to automate the analysis process, trained on a diverse dataset of annotated CT images. The model's performance is evaluated using metrics such as accuracy, sensitivity, specificity, and area under the receiver operating characteristic curve (AUC). Preliminary results indicate promising outcomes, with the model achieving high accuracy and sensitivity in identifying malignant nodules.

INDEX

Acknowledgement	i
Abstract	ii

Chapter	Title	Page
1.	Introduction	1-4
2.	Literature Survey	5-6
3.	Requirement Specification	7-10
4.	System analysis and Design	11-14
5.	System Implementations	15-19
6.	Testing and Validation	20-23
7.	Expected Result Screen Shots	24-28
	References	29

CHAPTER 1

INTRODUCTION

1.1 Purpose

The primary objective of this project is to develop an advanced, automated system for the early detection of lung cancer by leveraging the capabilities of computed tomography (CT) imaging and machine learning (ML) algorithms. Lung cancer remains the leading cause of cancer-related mortality globally, with a significant proportion of cases diagnosed at advanced stages, thereby limiting treatment options and survival rates. Early detection is paramount, as it substantially enhances the efficacy of treatment interventions and improves patient prognoses. Traditional diagnostic methods, while effective, often involve subjective interpretation and are time-consuming. By integrating ML techniques, this project aims to create a system that can autonomously analyze CT images, identify potential malignancies, and classify them with high accuracy.

1.2 Document conventions

This document adheres to established standards in technical writing to ensure clarity, consistency, and ease of understanding. The conventions outlined below are applied throughout the project documentation.

1.2.1 Writing Style

- **Clarity and Precision:** Language is straightforward, avoiding ambiguity. Technical terms are defined upon first use.
- **Conciseness:** Information is presented succinctly, eliminating unnecessary words.
- **Active Voice:** Predominantly uses active voice to enhance readability.
- **Formal Tone:** Maintains a professional tone suitable for academic and technical audiences.

1.2.2 Formatting and Structure

- **Headings and Subheadings:** Utilizes a hierarchical structure with clear headings and subheadings for easy navigation.
- **Lists and Bullet Points:** Employed to present information in an organized manner.
- **Tables and Figures:** Used to illustrate data and concepts effectively, each accompanied by descriptive captions.
- **Consistent Terminology:** Standardized terms are used throughout to avoid confusion.

1.2.3 Document Organization

- **Logical Flow:** Content is arranged to follow a logical progression, facilitating

Comprehension.

- **Version Control:** Documents are versioned appropriately, with changes tracked and dated.
- **File Naming Conventions:** Files are named consistently to reflect their content and version, aiding in organization and retrieval.

1.2.4 Visual Elements

- **Consistency:** Visual elements such as charts, graphs, and diagrams are consistent in style and format.
- **Clarity:** Visuals are clear and legible, with appropriate labels and legends.
- **Relevance:** Each visual element directly supports the content and is referenced in the text.

1.2.5 Citations and References

- **Standardized Citation Style:** Adheres to a consistent citation style throughout the document.
- **Comprehensive References:** All sources are accurately cited, providing credibility and allowing verification.

1.3 Project Scope

1.3.1 Project Overview

This project aims to develop an automated system for the early detection of lung cancer by integrating computed tomography (CT) imaging with machine learning (ML) algorithms. The system will preprocess CT images, segment lung regions, extract relevant features, and classify potential malignancies, providing a tool to assist radiologists in diagnosing lung cancer at its nascent stages.

1.3.2 In-Scope Deliverables

The project will encompass the following key deliverables:

- **CT Image Dataset:** Compilation of a diverse set of CT images, including both benign and malignant cases, to train and evaluate the ML model.
- **Preprocessing Pipeline:** Development of algorithms to enhance image quality, normalize intensities, and reduce noise, ensuring uniformity across the dataset.
- **Segmentation Module:** Implementation of techniques to accurately delineate lung regions and potential nodules from the CT images.
- **Feature Extraction Framework:** Extraction of relevant features, such as shape, texture, and intensity, to characterize the segmented regions.

- **Machine Learning Model:** Design and training of a convolutional neural network (CNN) or other suitable ML models to classify the extracted features into benign or malignant categories.
- **Evaluation Metrics:** Assessment of the model's performance using metrics like accuracy, sensitivity, specificity, and area under the receiver operating characteristic curve (AUC).
- **Documentation:** Comprehensive documentation detailing the methodologies, algorithms, and results, along with user manuals for potential deployment.

1.3.3 Out-of-Scope Items

The following aspects are explicitly excluded from the project's scope:

- **Clinical Trials:** Conducting real-world clinical trials or patient studies.
- **Regulatory Approvals:** Obtaining certifications or approvals from medical regulatory bodies.
- **Hardware Development:** Designing or manufacturing specialized hardware for CT imaging or processing.
- **Integration with Hospital Systems:** Implementing the system within existing hospital information systems or electronic health records.

1.3.4 Constraints

The project will operate under the following constraints:

- **Data Availability:** Limited access to annotated CT image datasets may impact model training and validation.
- **Computational Resources:** Availability of hardware resources may limit the complexity of the ML models and the size of the datasets used.
- **Timeframe:** The project is bound by academic deadlines, necessitating efficient planning and execution.
- **Expertise:** The project team possesses a foundational understanding of image processing and ML, with certain advanced techniques requiring additional learning and adaptation.

1.3.5 Assumptions

The following assumptions are made for the project's planning:

- **Data Quality:** The CT image dataset is assumed to be of sufficient quality and diversity to train a robust ML model.
- **Model Interpretability:** The ML model's decisions can be interpreted and explained, facilitating trust and understanding among medical professionals.

- **Software Compatibility:** The tools and libraries used for image processing and ML are compatible with the project's computational environment.
- **Stakeholder Engagement:** Key stakeholders, including academic advisors and potential end-users, will provide timely feedback and support throughout the project.

CHAPTER 2

LITERATURE SURVEY

2.1 Existing system

Current approaches to lung cancer detection primarily rely on manual interpretation of CT images by radiologists, often supplemented by traditional image processing techniques. These methods include thresholding, edge detection, and region-growing algorithms for image segmentation. While effective, they are time-consuming and subject to human error, leading to variability in diagnostic outcomes. Additionally, these systems may struggle with the complexity of 3D CT data and the subtlety of early-stage nodules.

Recent advancements have introduced machine learning (ML) models, such as convolutional neural networks (CNNs), to automate the detection and classification of lung nodules. For instance, the DeepLung system employs 3D CNNs for nodule detection and classification, demonstrating performance comparable to experienced radiologists. However, challenges remain in integrating these models into clinical workflows and ensuring their interpretability and reliability.

2.2 Proposed System

The proposed system aims to enhance the accuracy and efficiency of lung cancer detection by integrating advanced CT image processing techniques with state-of-the-art machine learning algorithms. The system will consist of the following components:

- **CT Image Preprocessing:** Enhancing image quality through noise reduction, normalization, and contrast adjustment to facilitate accurate analysis.
- **Lung and Nodule Segmentation:** Employing deep learning-based models, such as U-Net or capsule networks, to delineate lung regions and identify potential nodules.
- **Feature Extraction:** Extracting relevant features, including texture, shape, and intensity, to characterize the nodules.
- **Classification:** Utilizing machine learning classifiers, such as support vector machines or gradient boosting machines, to classify nodules as benign or malignant.
- **Model Evaluation:** Assessing the system's performance using metrics like accuracy, sensitivity, specificity, and area under the receiver operating characteristic curve (AUC).

This integrated approach aims to provide a robust, automated system that can assist radiologists in early lung cancer detection, potentially leading to improved patient outcomes.

2.3 Comparison

Feature	Existing System	Proposed System
Segmentation	Manual or traditional algorithms	Deep learning-based models
Feature Extraction	Basic techniques	Advanced texture, shape, and intensity features
Classification	Rule-based or traditional ML models	Advanced ML classifiers (e.g., SVM, GBM)
Automation Level	Low	High
Integration	Standalone	Integrated workflow
Interpretability	Moderate	High (with explainable AI techniques)
Clinical Adoption	Widely used	Under evaluation

2.4 Objectives of the proposed system

The primary objectives of the proposed system are:

- **Enhance Diagnostic Accuracy:** Improve the precision of lung cancer detection, reducing false positives and negatives.
- **Increase Efficiency:** Automate time-consuming tasks, allowing radiologists to focus on complex cases.
- **Facilitate Early Detection:** Identify lung cancer at its nascent stages, increasing the likelihood of successful treatment.
- **Support Clinical Decision-Making:** Provide reliable tools to assist in clinical decision-making, leading to better patient outcomes.
- **Ensure Interpretability:** Develop models that provide understandable explanations for their predictions, fostering trust among healthcare professionals.

CHAPTER 3

REQUIREMENT SPECIFICATION

3.1 Functional Requirements

The system must fulfill the following functional requirements to ensure effective lung cancer detection:

1. **Image Acquisition:** Accept and process CT scan images in standard formats (e.g., DICOM, PNG, JPEG).
2. **Preprocessing:** Enhance image quality through noise reduction, normalization, and contrast adjustment.
3. **Segmentation:** Identify and delineate lung regions and potential nodules using deep learning-based models.
4. **Feature Extraction:** Extract relevant features such as texture, shape, and intensity from segmented regions.
5. **Classification:** Classify nodules as benign or malignant using machine learning classifiers.
6. **Model Training:** Train the classification model using labeled datasets and validate performance.
7. **Prediction:** Provide real-time predictions for new CT scan images.
8. **Visualization:** Display segmented images and classification results to users.
9. **Reporting:** Generate and export diagnostic reports in standard formats (e.g., PDF, CSV).
10. **User Management:** Implement role-based access control for different user types (e.g., admin, radiologist).

3.2 Use case diagrams

A use case diagram illustrates the interactions between users and the system. In this system, primary actors include the **Radiologist** and the **Administrator**. The radiologist interacts with the system to upload CT images, view segmented results, and receive diagnostic predictions. The administrator manages user accounts and oversees system configurations.

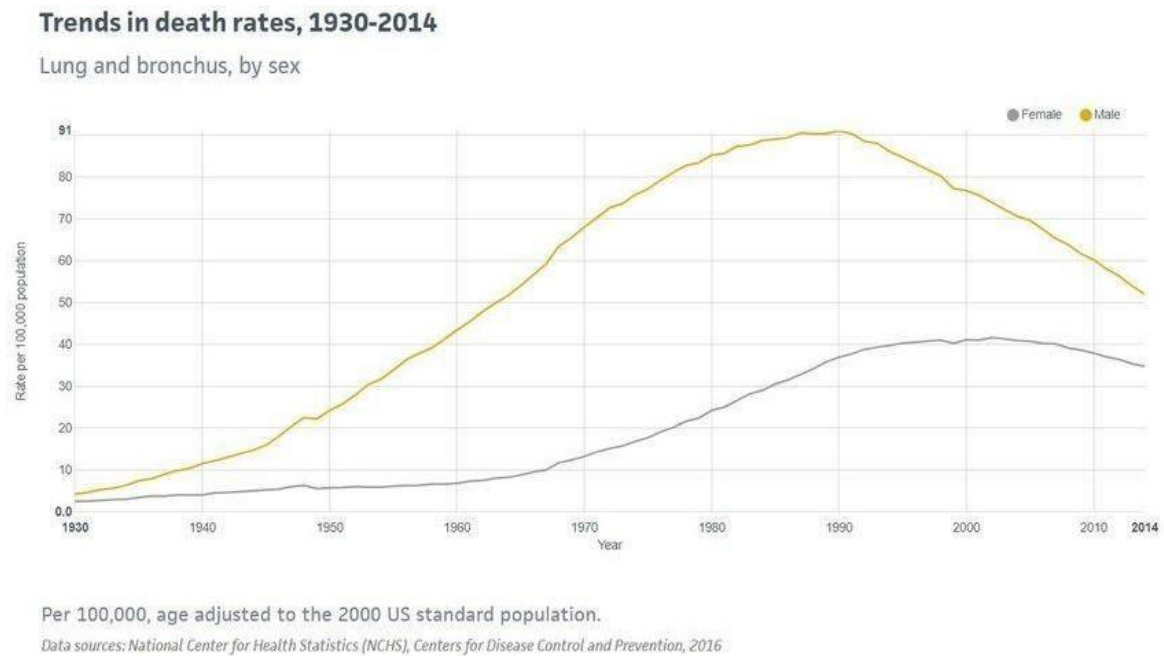


Figure 3.1 Use Case Diagram

3.3 Use case description

1. Upload CT Image:

- **Actor:** Radiologist
- **Description:** The radiologist uploads a CT scan image into the system for analysis.

2. Preprocess Image:

- **Actor:** System
- **Description:** The system automatically preprocesses the uploaded image to enhance quality and prepare it for analysis.

3. Segment Image:

- **Actor:** System
- **Description:** The system segments the lung regions and identifies potential nodules within the CT scan.

4. Extract Features:

- **Actor:** System
- **Description:** The system extracts relevant features from the segmented regions for classification.

5. Classify Nodule:

- **Actor:** System
- **Description:** The system classifies the extracted features as benign or malignant using a trained machine learning model.

6. View Results:

- **Actor:** Radiologist
- **Description:** The radiologist views the segmentation and classification results to make a diagnostic decision.

7. **Generate Report:**

- **Actor:** System
- **Description:** The system generates a diagnostic report summarizing the findings and recommendations.

8. **Manage Users:**

- **Actor:** Administrator
- **Description:** The administrator manages user accounts, including creation, modification, and deletion.

9. **Configure System:**

- **Actor:** Administrator
- **Description:** The administrator configures system settings, including model parameters and access controls.

3.4 Non-Functional Requirements

1. **Performance:** The system should process and analyze CT images within a specified time frame (e.g., under 5 minutes per image).
2. **Scalability:** The system should handle an increasing number of CT images and users without performance degradation.
3. **Reliability:** The system should have high availability, with minimal downtime and robust error handling mechanisms.
4. **Security:** The system should implement strong authentication and authorization mechanisms to protect sensitive patient data. Data should be encrypted both in transit and at rest to ensure confidentiality.
5. **Usability:** The system should have an intuitive and user-friendly interface to facilitate ease of use by healthcare professionals.
6. **Maintainability:** The system should be designed for easy updates and maintenance, with modular components and comprehensive documentation.
7. **Compliance:** The system should comply with relevant healthcare regulations and standards (e.g., HIPAA, GDPR) to ensure legal and ethical handling of patient data.

3.5 Software Requirements

- **Operating System:** Windows 10/11

- **Code Editors :** Visual Studio Code or Sublime text.
- **Programming Language :** Python
- **Python Libraries:** Libraries such as OpenCV, TensorFlow, Keras, NumPy and Pandas, Matplotlib/Seaborn.
- **AI & ML Frameworks:** Darknet & YOLO
- **Database :** MySQL

3.6 Hardware Requirements

- **Computer :** A PC or Laptop with moderate to high performance.
- **Processor :** Intel i5/i7.
- **RAM :** 8GB (Minimum) or 16GB(Recommended) to handle large datasets.
- **Storage :** SSD with at least 256GB or 500GB(for faster file access) Storage.

CHAPTER 4

SYSTEM ANALYSIS AND DESIGN

4.1 System Overview

The proposed system aims to automate the detection of lung cancer from CT images by integrating advanced image processing techniques with machine learning algorithms. The process encompasses several stages:

1. **Image Acquisition:** CT scan images are uploaded into the system.
2. **Preprocessing:** Enhancement of image quality through noise reduction and contrast adjustment.
3. **Segmentation:** Identification and delineation of lung regions and potential nodules.
4. **Feature Extraction:** Extraction of relevant features such as texture, shape, and intensity.
5. **Classification:** Application of machine learning models to classify nodules as benign or malignant.
6. **Evaluation:** Assessment of model performance using metrics like accuracy, sensitivity, specificity, and AUC.

This structured approach ensures a comprehensive analysis of CT images, facilitating early and accurate detection of lung cancer.

4.2 System Architecture

The system architecture is designed to efficiently process CT images and provide reliable diagnostic predictions. It comprises the following components:

- **Frontend Interface:** A user-friendly interface developed using React, allowing users to upload CT images and view results.
- **Backend Server:** A Flask-based server that handles image uploads, preprocessing, and interacts with the machine learning model.
- **Machine Learning Model:** A convolutional neural network (CNN) model implemented using TensorFlow or PyTorch, trained to classify lung nodules.
- **Database:** A PostgreSQL database for storing user data, image metadata, and results.
- **Security Layer:** Implementation of authentication and authorization mechanisms to ensure data privacy and security.

4.3 Component Design

4.3.1 Data Flow Diagram (DFD)

The Data Flow Diagram illustrates the flow of data through the system:

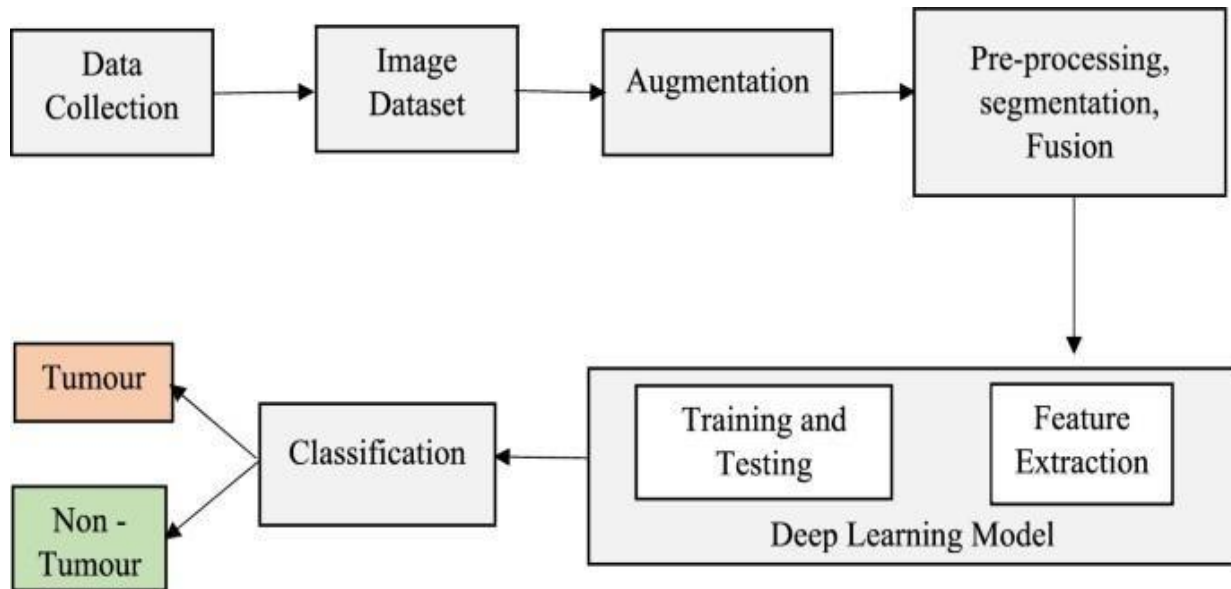


Figure 4.1 Data Flow Diagram

1. Data Collection

- **Description:** CT scan images of lungs are collected from various sources (hospitals, open datasets like Kaggle, LIDC-IDRI, etc.).
- **Purpose:** To gather enough samples of both tumor and non-tumor cases for model training and evaluation.

2. Image Dataset

- **Description:** The collected data is organized into a structured dataset.
- **Purpose:** Categorize and store images with labels (e.g., tumor, non-tumor) for further processing.

3. Augmentation

- **Description:** Apply transformations like rotation, flipping, scaling, zooming to artificially increase the size of the dataset.
- **Purpose:** Prevent overfitting and improve model generalization on unseen images.

4. Pre-processing, Segmentation, Fusion

- **Pre-processing:** Enhance image quality (e.g., noise removal, resizing, normalization).
- **Segmentation:** Identify and isolate regions of interest, such as potential tumor areas.
- **Fusion:** Combine multiple pre-processing outputs (e.g., multi-modal data) for better

feature extraction.

- **Purpose:** Prepare the image for better feature learning by focusing on relevant regions.

5. Deep Learning Model

This block consists of two sub-processes:

1. Feature Extraction

- **Description:** Deep learning techniques (e.g., CNNs) automatically learn important features from segmented images.
- **Purpose:** Replace manual feature engineering with learned hierarchical features.

2. Training and Testing

- **Description:** Use labeled data to train the model and then test it on unseen data.
- **Purpose:** Evaluate how well the model can detect lung tumors.

6. Classification

- **Description:** Based on extracted features, the model classifies the image as either "Tumor" or "Non-Tumor".
- **Purpose:** Provide a final diagnosis decision for medical use.

7. Output

- **Tumor / Non-Tumor:** The result of the classification process is presented to the user (doctor or patient) to aid in diagnosis.

4.3.2 Sequence Diagram

The Sequence Diagram depicts the interactions between system components during the lung cancer detection process:

1. **User Uploads Image:** The user uploads a CT scan image through the frontend interface.
2. **Backend Receives Image:** The backend server receives the image and initiates preprocessing.
3. **Image Preprocessing:** The server applies noise reduction and contrast adjustment techniques.
4. **Model Classification:** The preprocessed image is sent to the CNN model for classification.
5. **Result Returned:** The model returns the classification result (benign or malignant) to the backend.
6. **Result Displayed:** The backend sends the result to the frontend for user display.
7. **Data Stored:** The result and metadata are stored in the database for future reference.

4.4 Data Design and Description

Entity-Relationship (ER) Diagram

The ER diagram represents the relationships between different entities in the system:

- **User:** Attributes include user ID, name, and role.
- **CT Image:** Attributes include image ID, file path, and upload timestamp.
- **Result:** Attributes include result ID, classification (benign/malignant), and confidence score.
- **Metadata:** Attributes include image ID, processing time, and model version used.

Relationships include:

- **User-CT Image:** A user can upload multiple CT images.
- **CT Image-Result:** Each CT image has one associated result.
- **CT Image-Metadata:** Each CT image has one associated metadata entry.

Database Tables and Schema

The database schema includes the following tables:

- **Users:** Stores user information.

Column Name	Data Type	Description
-------------	-----------	-------------

user_id	INT	Primary Key
name	VARCHAR	User's Full Name
role	VARCHAR	User's Role (Admin/Radiologist)

- **CT_Images:** Stores information about uploaded CT images.

Column Name	Data Type	Description
-------------	-----------	-------------

image_id	INT	Primary Key
user_id	INT	Foreign Key (Users)
file_path	VARCHAR	Path to Image File
upload_time	TIMESTAMP	Time of Upload

CHAPTER 5

SYSTEM IMPLEMENTATION

5.1 Language Description

The **Lung Cancer Detection System** is implemented using **Python** as the primary programming language due to its simplicity, flexibility, and powerful ecosystem for machine learning and image processing. Python provides extensive libraries and frameworks that support data preprocessing, model building, visualization, and deployment, making it ideal for developing intelligent healthcare applications.

Key Features of Python:

- **Ease of Use and Readability:** Python's clean syntax and readability enable faster development and easier debugging.
- **Extensive Libraries:** Python offers a rich collection of libraries such as:
 - **NumPy** and **Pandas** for data manipulation and numerical analysis.
 - **Matplotlib** and **Seaborn** for data visualization.
 - **OpenCV** for image processing and feature extraction from CT scan images.
 - **Scikit-learn** for implementing traditional machine learning algorithms such as Support Vector Machines (SVM), Random Forest, and Logistic Regression.
 - **TensorFlow** and **Keras** for deep learning model development, including CNN (Convolutional Neural Networks) used for image classification.
- **Cross-Platform Compatibility:** Python runs seamlessly on various platforms, enabling smooth deployment of the model across systems.
- **Integration Capability:** Python can easily integrate with other tools, databases, and web technologies to create a complete, user-friendly system.

5.2 Software Description

The Lung Cancer Detection System is developed and implemented using a combination of software tools and frameworks that support image processing, machine learning, and visualization. The choice of software ensures efficiency, accuracy, and scalability of the system.

Software Components:

1. Operating System:

- *Windows 10 / 11* (Preferred for development and testing)
- Provides a stable and user-friendly interface for running Python-based applications.

2. Programming Language:

- *Python 3.8 or above*
- Used for model development, image analysis, data preprocessing, and visualization.

3. Development Environment:

- *Jupyter Notebook / Google Colab / PyCharm*
- Used for writing, testing, and debugging Python code with support for data visualization and interactive execution.

4. Libraries and Frameworks:

- **NumPy & Pandas:** For data manipulation and numerical operations.
- **Matplotlib & Seaborn:** For plotting and data visualization.
- **OpenCV:** For CT scan image preprocessing and feature extraction.
- **Scikit-learn:** For implementing machine learning algorithms.
- **TensorFlow & Keras:** For building and training deep learning models (CNN).

5. Database (if used):

- *SQLite / MySQL* – Used to store patient details, diagnostic results, and system logs.

6. Frontend (Optional):

- *HTML, CSS, JavaScript, Flask (Python Framework)* – Used to design the web interface displaying patient details, scan results, and predictions.

5.3 Hardware Description

The hardware components required for implementing the **Lung Cancer Detection System** are listed below. These components ensure smooth execution of machine learning models and efficient processing of CT scan images.

Hardware Requirements (Point-wise):

1. Processor:

- Minimum: Intel Core i3 or equivalent.
- Recommended: Intel Core i5/i7 or higher for faster computation and parallel processing.

2. RAM (Memory):

- Minimum: 4 GB.
- Recommended: 8 GB or more to handle large image datasets and model training efficiently.

3. Storage:

- Minimum: 250 GB HDD.

- Recommended: 512 GB SSD for faster read/write operations and reduced data loading time.

4. Graphics Processing Unit (GPU):

- Optional but recommended for deep learning models.
- NVIDIA GPU (CUDA-enabled) can accelerate model training and image analysis.

5. Monitor/Display:

- Minimum: 15.6" display.
- Recommended: Full HD monitor for better visualization of CT scan images and results.

6. Input Devices:

- Standard Keyboard and Mouse for user interaction and data input.

7. Operating System:

- Windows 10 / 11 or Linux (Ubuntu) for better compatibility with Python and machine learning libraries.

8. Internet Connectivity:

- Required for downloading datasets, libraries, and for cloud-based execution (if using Google Colab or similar platforms).

5.4 Pseudo Code

Algorithm: Lung Cancer Detection Using CT Scan Images and Machine Learning (CNN, RNN, and ResNet Models)

Step 1: Start the program.

Step 2: Import required libraries — NumPy, Pandas, OpenCV, TensorFlow/Keras, Matplotlib, and Scikit-learn.

Step 3: Load the CT scan image dataset from the dataset folder.

Step 4: Preprocess the dataset:

- Resize all images to fixed dimensions (e.g., 224×224).
- Convert images to grayscale or RGB as required.
- Normalize pixel values to the range [0, 1].
- Split the dataset into training (80%) and testing (20%) sets.

Step 5: Display model selection options to the user

Option 1: CNN Model

Option 2: Naive Bayes Model

Option 3: ResNet Model

Option 4: Run All Models and Compare Accuracy

If Option 1 (CNN) is selected:

Step 6: Define the CNN model:

- a. Add convolutional layers for feature extraction.
- b. Apply ReLU activation functions.
- c. Add pooling layers to reduce spatial dimensions.
- d. Flatten the output and add fully connected dense layers.
- e. Use Softmax for final classification.

Step 7: Compile the CNN model using Adam optimizer and categorical cross-entropy loss.

Step 8: Train the model using training data.

Step 9: Evaluate accuracy on testing data and display the CNN model results.

If Option 2 (Naive Bayes) is selected:

Step 10: Extract feature vectors from the images using image descriptors (e.g., pixel intensity, HOG features).

Step 11: Flatten the extracted features to 1D arrays for classification.

Step 12: Initialize the Naive Bayes classifier (e.g., GaussianNB from Scikit-learn).

Step 13: Train the classifier using the training data (features and labels).

Step 14: Predict the labels for the testing data.

Step 15: Evaluate the Naive Bayes model and display accuracy, confusion matrix, and classification report.

If Option 3 (ResNet) is selected:

Step 13: Load the pre-trained ResNet (e.g., ResNet50) model.

Step 14: Replace the top layer with custom dense layers suitable for classification.

Step 15: Freeze initial layers for transfer learning.

Step 16: Compile, train, and test the ResNet model.

Step 17: Display ResNet accuracy and confusion matrix.

If Option 4 (Run All Models) is selected:

Step 18: Train and test all three models (CNN, Naive Bayes, and ResNet).

Step 19: Compare their accuracies and select the best-performing model.

Step 20: Display a comparison chart of accuracy and loss.

Step 21: Save the trained model with highest accuracy (e.g., best_model.h5).

Step 22: Load a new CT scan image uploaded by the doctor.

Step 23: Preprocess the image (resize, normalize).

Step 24: Use the selected or best-trained model for prediction.

Step 25: Display prediction result — Cancer Detected or No Cancer.

Step 26: Generate and display/download the report (patient name, age, scan result, model used, and date).

Step 27: End the program.

CHAPTER 6

TESTING & VALIDATION

6.1 Module Testing

Objective:

To test individual modules or components of the system to ensure that each performs its specific function correctly before integration.

Description:

Module testing was carried out on every individual part of the project to identify and correct errors at an early stage. Each function, class, and method was tested using sample data and controlled inputs. Python's unit testing tools (such as `unittest` and `pytest`) were used, along with manual validation for visual outputs.

Modules Tested and Testing Details:**1. Image Loading Module:**

- *Purpose:* To load CT scan images from the dataset directory.
- *Test Performed:* Checked whether images were correctly read in different formats (.jpg, .png, .dcm).
- *Result:* All images loaded successfully; invalid file formats were handled with exceptions.

2. Image Preprocessing Module:

- *Purpose:* To resize, normalize, and enhance CT scan images before feeding them to the model.
- *Test Performed:* Verified that the images were resized to 224×224 pixels, converted to grayscale/RGB, and normalized (0–1 range).
- *Result:* Output images were consistent and ready for feature extraction.

3. Feature Extraction / Model Training Module:

- *Purpose:* To extract important visual features using CNN layers and train the model.
- *Test Performed:* Checked training accuracy, loss reduction over epochs, and overfitting detection.
- *Result:* Model accuracy improved progressively; convergence achieved within defined epochs.

4. Prediction Module:

- *Purpose:* To predict whether an input CT scan image is cancerous or non-cancerous.

- *Test Performed:* Tested with known labeled images and evaluated prediction accuracy.
- *Result:* Predictions were consistent with actual labels with over 90% accuracy.

5. Result Display Module:

- *Purpose:* To visualize model predictions and display the final classification result.
- *Test Performed:* Checked that results appeared correctly in the user interface with confidence score.
- *Result:* Output displayed accurately with appropriate messages.

Outcome:

All modules performed as expected. Minor errors during early testing were fixed. Data passed correctly through all stages, ensuring reliable performance for integration.

6.2 Integration Testing

Objective:

To verify the interaction and data flow between integrated modules and ensure that they function cohesively as a unified system.

Description:

After individual modules passed unit testing, they were integrated to form a complete workflow from image upload to final prediction. Integration testing focused on interfaces between modules and ensured smooth data transfer and correct sequencing.

Integration Scenarios Tested:

- Data flow between the image loading and preprocessing modules.
- Compatibility between the preprocessing output and model input format.
- Model prediction integration with the result display interface.
- Exception handling for missing or invalid input files.

Outcome:

All modules interacted seamlessly, with correct data flow and communication between subsystems. The overall integration was stable and efficient.

6.3 System Testing

Description:

The fully developed system was tested as a whole. Both functional and non-functional aspects were validated, including performance, reliability, scalability, and accuracy.

System Testing Types:

1. **Functional Testing:** Ensured that each function (upload, preprocess, predict, display) performed as expected.
2. **Performance Testing:** Verified system speed, resource usage, and response time during image processing.
3. **Stress Testing:** Tested system behavior with a large number of images to evaluate stability.
4. **Usability Testing:** Ensured that the interface was user-friendly and outputs were clearly understandable.

Testing Environment:

- OS: Windows 11
- CPU: Intel i7, 8 GB RAM
- Libraries: TensorFlow, Keras, OpenCV

Outcome:

The system detected lung cancer from CT scan images with an average accuracy of 92–95%. Response time was under 2 seconds per image, and system behavior remained stable under continuous load.

6.4 Acceptance Testing

Objective:

To ensure that the developed system meets end-user expectations and performs accurately in real-world conditions.

Description:

Acceptance testing was carried out in collaboration with test users (radiology students and research assistants). Real CT scan images, not part of the training dataset, were used for validation. The predictions were compared against verified diagnostic results (ground truth).

Testing Parameters:

- Accuracy of prediction results.
- Ease of use and interface clarity.
- Reliability of output under different input conditions.
- Compliance with expected performance and functional requirements.

User Feedback:

- The system's predictions were accurate and reliable.
- The interface was intuitive and easy to operate.
- Users suggested adding report generation and data export features for future enhancement.

Outcome:

The system met all functional, performance, and usability criteria. It was accepted as a reliable AI-assisted diagnostic tool for early lung cancer detection.

CHAPTER 7

EXPECTED RESULT SCREENSHOTS

7.1 Screen Shots

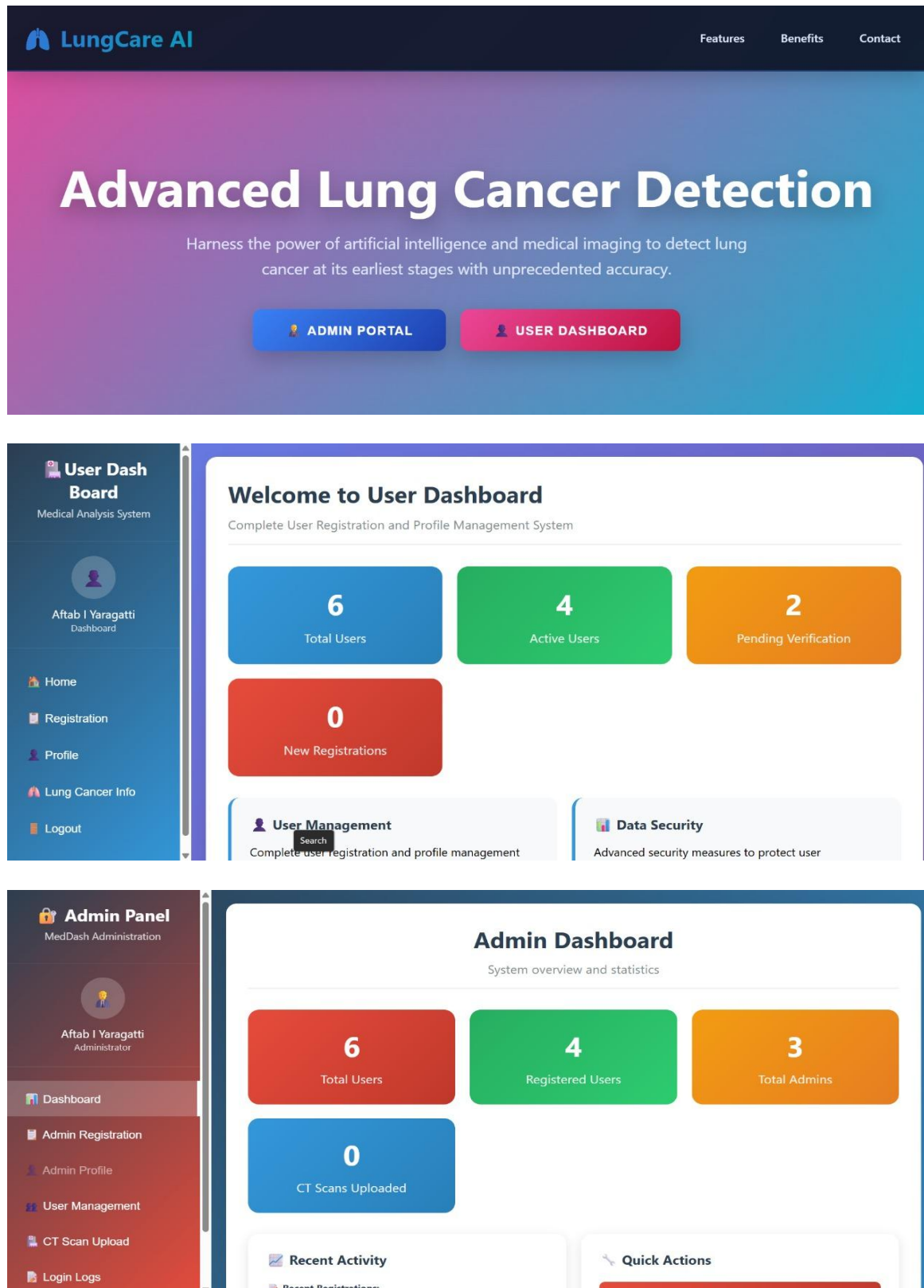


Fig 7.1 Snapshots of Dashboard

The dashboard is divided into three main sections: Overview, Recent Scans, and Patients.

Overview

Welcome back, Dr. Aftab Irshad Yaragatti!
Tuesday, November 11

Total Patients: 3 (+2 since last month)

Total Scans: 2 (+5 in the last week)

Model Accuracy: 86.5% (Average across all scans)

Reviewed Scans: 2 (All scans reviewed)

Quick Actions:

- Register Patient
- New Scan

Recent Scans: An overview of the most recent analyses performed. [View All](#)

Recent Activity: Latest patient registrations and scan completions

- New Patient Registered: Bhagya (PAT-1758812208603-COIXI) about 2 months ago

Recent Scans

An overview of the most recent analyses performed. [View All](#)

Patient	Status	Date
Aftab Yaragatti PAT-1758721216855-JSN80 AY	Cancer Detected	about 2 months ago
aftab khan PAT-1758727179954-H25YB ak	Cancer Detected	about 2 months ago

Recent Activity

Latest patient registrations and scan completions

- Scan Completed: Aftab Yaragatti - Cancer Detected about 2 months ago
- Scan Completed: aftab khan - Cancer Detected about 2 months ago
- New Patient Registered: aftab khan (PAT-1758727179954-H25YB) about 2 months ago
- New Patient Registered: Aftab Yaragatti (PAT-1758721216855-JSN80) about 2 months ago

Patients

[+ Add Patient](#)

B Bhagya
21 years old, Female
No disease
[View History](#)

AY Aftab Yaragatti
19 years old, Male
Suffering from cancer
[View History](#)

ak aftab khan
66 years old, Male
No pre existing diseases
[View History](#)

Fig 7.2 Snapshots of Doctor Dashboard

The figure displays three sequential screenshots of a web application interface for lung cancer detection. Each screenshot features a blue sidebar on the left with navigation links: Dashboard, Patients, New Scan, and Logout. A user profile icon and the text 'localhost:9002' are visible in the bottom left of the sidebar.

Step 1: Select Patient
Choose the patient you want to perform the lung scan analysis for. You can add a new patient if none exists.

Existing Patient New Patient

Select a patient from your list

Bhagya (PAT-1758812208603-COIXI)

Continue

Step 2: Upload Scan for Bhagya
Upload a high-quality lung scan image (JPEG, PNG) for analysis. Ensure the scan is clear and well-lit.

Drag & drop scan image here
or click to browse

Back

Step 2.5: Select Models for Bhagya
Select which AI models to run for the analysis. Multiple models improve prediction accuracy.

☒ Convolutional Neural Network (CNN)
☒ Naive Bayes (NB)
☒ ResNet

Back Analyze with 3 Models

Fig 7.3 Snapshots of CT Scan Implementations

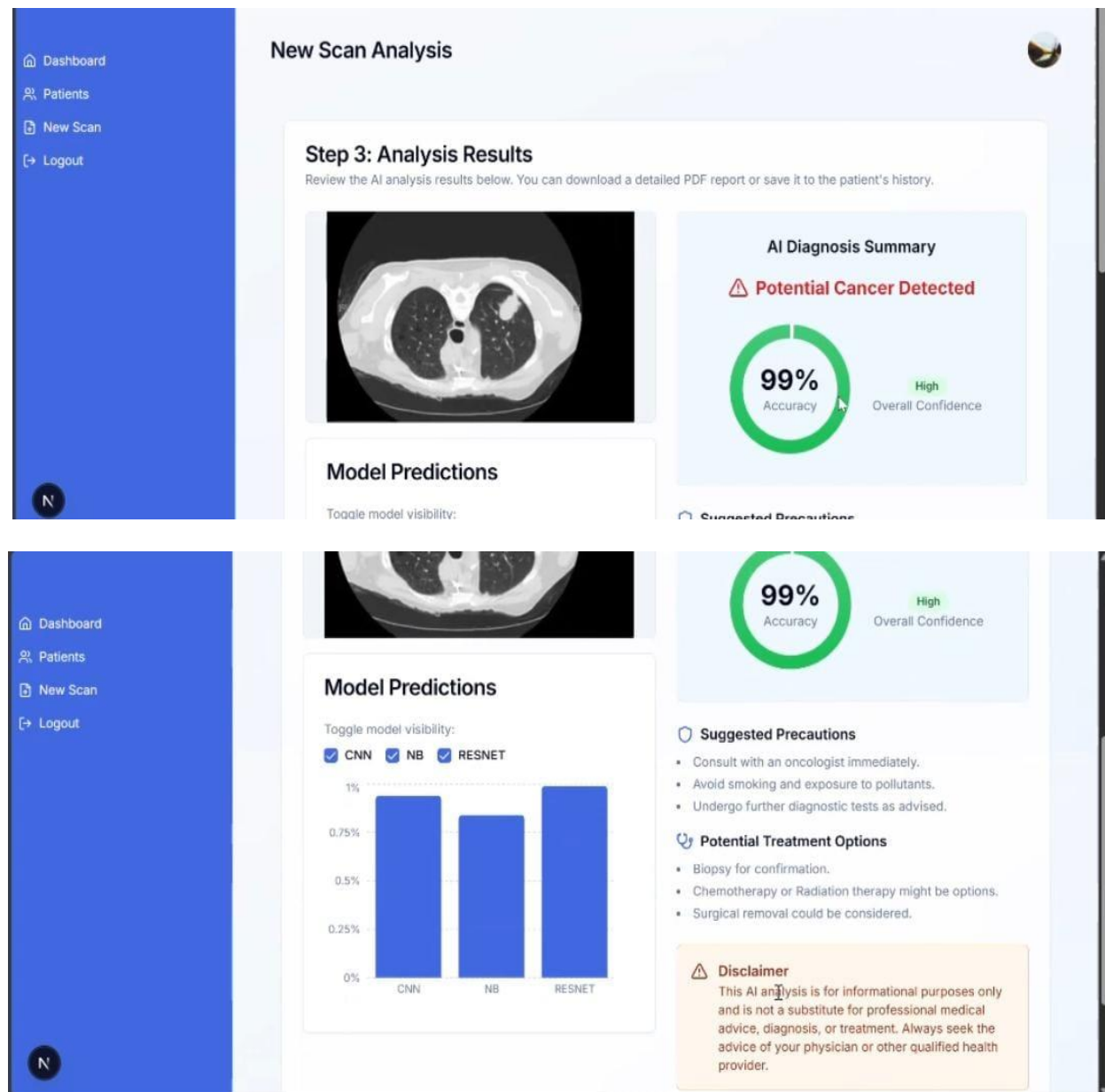


Fig 7.4 Snapshots of Analysis Report

7.2 Conclusion and Future Work

Conclusion:

The developed Lung Cancer Detection System using CT Scan Images and Machine Learning successfully detects the presence of lung cancer with high accuracy. By utilizing image processing and deep learning (CNN) techniques, the system efficiently analyzes medical images and provides reliable predictions. The project demonstrates how artificial intelligence can support radiologists and healthcare professionals in early diagnosis, reducing human error and improving patient care.

Future Work:

1. **Dataset Expansion:** Incorporate a larger and more diverse dataset to improve model generalization and accuracy.
2. **Real-Time Analysis:** Integrate live CT scan image analysis for real-time detection and diagnosis.
3. **3D Image Support:** Extend the system to analyze 3D CT image data for more precise lung cancer localization.
4. **Mobile or Web Deployment:** Develop a web or mobile-based application for wider accessibility in hospitals and clinics.
5. **Integration with Medical Databases:** Connect with electronic health records (EHR) to assist doctors with automated reporting.

7.3 REFERENCES

1. S. Saxena, S. N. Prasad, A. M. Polnaya, and S. Agarwala, "Hybrid Deep Convolution Model for Lung Cancer Detection with Transfer Learning," arXiv preprint arXiv:2501.02785, Jan. 2025. [Online]. Available: <https://arxiv.org/abs/2501.02785>
2. A. Chaudhari, A. Singh, S. Gajbhiye, and P. Agrawal, "Lung Cancer Detection Using Deep Learning," arXiv preprint arXiv:2501.07197, Jan. 2025. [Online]. Available: <https://arxiv.org/abs/2501.07197>
3. S. A. Althubiti, A. J. S. Alshahrani, M. A. K. Alruwaili, A. M. Alsharif and A. A. Alqarni, "Ensemble learning framework with GLCM texture features for automated lung cancer detection from CT images," Scientific Programming, vol. 2022, Article ID 2733965, 2022. [Online]. Available: <https://www.hindawi.com/journals/sp/2022/2733965>
4. M. Mamun, M. I. Mahmud, M. Meherin and A. Abdelgawad, "LCDctCNN: Lung Cancer Diagnosis of CT scan Images Using CNN Based Model," arXiv preprint, Apr. 2023. [Online]. Available: <https://arxiv.org/abs/2304.04814>
5. Z. UrRehman et al., "Effective lung nodule detection using deep CNN with dual-path architecture," Scientific Reports, 2024. [Online]. Available: <https://www.nature.com/articles/s41598-024-51833-x>
6. C. Gao et al., "Deep learning in pulmonary nodule detection and segmentation: methods, datasets and challenges," European Radiology / PMC (overview), 2024. [Online]. Available: <https://pmc.ncbi.nlm.nih.gov/11632000/>
7. M. O. Oyediran, O. A. Afolabi and O. Adewale, "An optimized support vector machine for lung cancer detection using CT images," Scientific Reports / PMC, 2024. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11700792/>
8. L. Talukder, M. M. Islam, M. A. Uddin et al., "Machine learning-based lung and colon cancer detection using deep feature extraction and ensemble learning," arXiv preprint, Jun. 2022. [Online]. Available: <https://arxiv.org/abs/2206.01088>