

Lung cancer detection using Machine Learning

Dipanwita Adhikary

Machine Learning Internship

at

Feynn Lab

27.01.2024

1. Abstract

Due in large part to delaying diagnosis and restricted access to cutting-edge screening methods, lung cancer is among the world's leading causes of mortality from cancer. In medical diagnostics, machine learning (ML) has become a revolutionary technique that provides reliable and effective early lung cancer detection methods. ML models can precisely evaluate medical imaging, including CT scans and X-rays, to detect cancerous nodules and categorize them according to their severity using advanced algorithms and big datasets. While reducing mistakes by humans, methods like deep learning, feature extraction, and predictive modeling improve diagnosis accuracy. Additionally, ML-based systems optimize healthcare resource allocation by enabling risk assessment and customized screening procedures. By addressing issues like data quality, understanding, and clinical advancement, this study examines how machine learning might enhance the detection of lung cancer. It also emphasizes how it has the potential to revolutionize healthcare procedures to improve patient outcomes and obtain earlier diagnoses.

2. Problem Statement

2.1 High Mortality Rate: Mostly due to late-stage diagnosis, lung cancer continues to be one of the world's primary causes of death due to cancer.

2.2 Problems with Manual Diagnostics: Conventional diagnostic techniques, such as CT scans and X-ray interpretations, are labor-intensive, prone to human error, and significantly dependent on the knowledge of radiologists.

2.3 Missed Early Detection: When lung cancer is detected in its early stages, it frequently exhibits minimal symptoms that are easy to overlook during manual examinations. This can result in treatment delays and unfavorable patient outcomes.

2.4 Increasing Imaging Workload: Healthcare providers are under a lot of pressure to maintain consistency and accuracy in diagnosis due to the increasing amount of medical imaging data.

2.5 Need for Automation: Automation is required since there aren't many automated systems that can evaluate imaging data efficiently, precisely, and economically, especially in environments with limited resources.

2.6 Interpretability and Trust: Healthcare practitioners must accept and trust machine learning models that offer interpretable and actionable insights.

2.7 Scalability and Integration: To manage various patient groups, effective solutions need to be scalable and easily incorporated into current clinical procedures.

2.8 Data Challenges: To create reliable machine learning models, it is essential to address problems with data quality, heterogeneity, and a lack of labeled datasets.

3. Business Need Assessment

3.1 Market Demand for Early Detection

- Although lung cancer continues to be the world's largest cause of cancer-related death, early detection methods are essential.
- The demand for precise, quick, and non-invasive diagnostic technologies is growing as disease prevention is widely recognized.

3.2 Cost-Effective Healthcare Solutions

- Early detection significantly reduces treatment costs compared to advanced-stage interventions.
- Automation via machine learning minimizes the need for redundant or resource-intensive diagnostic procedures.

3.3 Efficiency in Healthcare Delivery

- Machine learning models can analyze imaging data faster than traditional manual methods, streamlining diagnostic workflows.
- Automation reduces radiologists' workload and improves efficiency, especially in regions with limited healthcare professionals.

3.4 Scalability for Broader Reach

- ML-based tools can be deployed in various settings, from high-end hospitals to low-resource areas, making lung cancer detection accessible to a larger population.

3.5 Data-Driven Insights

- Machine learning can provide predictive analytics and risk assessments, enabling personalized screening strategies and targeted therapies.

- Insights from ML models can support pharmaceutical companies in clinical trials and drug development.

By addressing these commercial requirements, machine learning for lung cancer diagnosis not only enhances patient outcomes but also stimulates innovation and operational efficiency in the medical field, providing substantial benefits to all stakeholders involved.

4. Target specifications and characterizations

4.1 Accuracy and Sensitivity

- **Target Specification:** Achieve high diagnostic accuracy (>95%) for detecting lung cancer in medical imaging datasets.
- **Characterization:** Minimize false negatives to ensure early-stage cancers are detected effectively.

4.2 Specificity

- **Target Specification:** Ensure a high specificity (>90%) to reduce false positives and prevent unnecessary biopsies or treatments.
- **Characterization:** Correctly classify benign nodules to optimize clinical decision-making.

4.3 Interpretability

- **Target Specification:** Provide explainable results that allow clinicians to understand the rationale behind the predictions.
- **Characterization:** Generate heatmaps, feature importance scores, or region-specific annotations for transparency.

4.4 Integration Capability

- **Target Specification:** Seamlessly integrate with existing radiology systems (e.g., PACS) and electronic health records (EHRs).
- **Characterization:** Compatibility with DICOM standards and interoperability with healthcare IT infrastructure.

4.5 Scalability

- **Target Specification:** Scalable to handle large volumes of data from multiple healthcare facilities.
- **Characterization:** Cloud-based or edge-based deployment options for efficient performance across regions.

5. Benchmarking

To determine the efficacy, accuracy, and dependability of machine learning (ML) models and systems, benchmarking for lung cancer detection involves assessing them to pre-established performance measures, datasets, and standards. The main features and methods of benchmarking are listed below:

5.1 Datasets

- **Commonly Used Datasets:**
 - **LIDC-IDRI (Lung Image Database Consortium and Image Database Resource Initiative):** Annotated CT scans with nodule-level details for lung cancer analysis.
 - **NLST (National Lung Screening Trial):** Includes low-dose CT scan data for lung cancer screening.
 - **SPIE-AAPM Lung CT Challenge Dataset:** Focuses on lung nodule classification.
- **Benchmarking Criteria:**
 - Evaluate model performance on publicly available and widely accepted datasets.
 - Ensure consistent preprocessing (e.g., normalization, segmentation).

5.2 Model Architectures

- **State-of-the-Art (SOTA) Models:**
 - **Convolutional Neural Networks (CNNs):** Commonly used for image classification and segmentation (e.g., ResNet, EfficientNet).

- **3D CNNs:** Analyse volumetric data from CT scans.
- **Transformers:** Emerging architectures (e.g., Vision Transformers) for medical imaging.
- **Hybrid Models:** Combining traditional machine learning with deep learning for improved interpretability.
- **Benchmarking Criteria:**
 - Compare new models with baseline ML algorithms (e.g., Random Forest, SVM).
 - Measure training time, computational efficiency, and hardware requirements.

6. Applicable patents

When developing machine learning models for lung cancer detection, it's essential to be aware of existing patents to navigate the intellectual property landscape effectively.

Notable patents in this domain include:

- **Automated Lung Cancer Detection from Registered CT and PET Scans.** This patent (US20220245821A1) describes a system for the automated detection and segmentation of lung cancer by processing data from registered pairs of thoracic CT and PET scans using deep learning methods.
- **Methods and Machine Learning Systems for Predicting the Likelihood or Risk of Having Cancer** Patent WO2016094330A2 outlines non-invasive methods and diagnostic tests that measure biomarkers, combined with machine learning systems, to assess the likelihood or risk of a patient having cancer.

7. Applicable regulations

To ensure reliability, efficacy, and moral application, machine learning (ML) systems for lung cancer diagnosis in healthcare must be developed and implemented following several regulatory frameworks. The main relevant rules and guidelines are shown below:

7.1 Medical device regulations

ML-based lung cancer detection systems often qualify as medical devices, requiring adherence to specific regulations:

United States:

FDA (Food and Drug Administration):

- **21 CFR Part 820 (Quality System Regulation):** Ensures quality management in device manufacturing.
- **510(k) Premarket Notification or PMA (Premarket Approval):** For demonstrating safety and effectiveness of AI/ML-based diagnostic tools.
- **Software as a Medical Device (SaMD) guideline:** Specific to AI/ML systems in healthcare.

European Union:

- **EU MDR (Medical Device Regulation) 2017/745:** Governs CE marking for medical devices, including ML-based systems.
- **ISO 13485:** Quality management for medical devices.

7.2 Data Privacy and Security

Ensuring patient data privacy and secure handling of sensitive health data is critical:

- **HIPAA (Health Insurance Portability and Accountability Act):** U.S. regulation governing the protection of patient health information.
- **GDPR (General Data Protection Regulation):** EU regulation emphasizing data protection, including explicit consent for using patient data in ML systems.
- **ISO/IEC 27001:** International standard for information security management systems.

7.3 Standards for AI in Healthcare

- **ISO/IEC TR 24028:** Guidelines for AI trustworthiness.
- **IEC 62304:** Standards for medical device software lifecycle processes.
- **DICOM Standards:** For integrating ML models into radiology systems like PACS.

7.4 Ethical and Social Responsibility

Follow principles outlined by organizations such as:

- **World Health Organization (WHO):** Guidance on AI ethics in healthcare.
- **Bioethics Committees:** These committees are for balancing innovation and patient well-being.

7.5 Research and Training

Adherence to research ethics protocols such as:

- Institutional Review Board (IRB) approvals.
- Transparent reporting of ML model development and validation following standards like CONSORT-AI or TRIPOD.

8. Applicable constraints

Several financial, spatial, and skill-related limitations need to be taken into account while creating and implementing a Machine Learning (ML) system for lung cancer detection:

8.1 Data Storage:

- Medical imaging data, such as CT scans and X-rays, are highly voluminous and require significant storage capacity.
- Constraints include the need for high-performance storage solutions (local servers or cloud platforms) that comply with data protection regulations (e.g., HIPAA, GDPR).

8.1.1 Computational Infrastructure:

- ML model training and inference require dedicated space for high-performance computing resources (e.g., GPUs, TPUs, or specialized hardware for deep learning).
- For on-site deployment, space for housing hardware and maintaining adequate cooling systems is necessary.

8.1.2 Integration with Existing Systems:

- Space must accommodate integration with hospital IT systems, such as PACS (Picture Archiving and Communication Systems), without disrupting current workflows.

8.2 Budget Constraints

8.2.1 Development Costs:

- Training ML models requires access to high-quality labeled datasets, which may incur costs for annotation and licensing.
- Hiring or outsourcing to experienced data scientists and ML engineers adds to the budget.

8.2.2 Infrastructure Costs:

- Investment in high-performance hardware (e.g., GPUs, TPUs) for training and deployment can be significant.
- Cloud computing services (e.g., AWS, Google Cloud, Azure) offer scalability but may incur recurring costs.

8.3 Expertise Constraints

8.3.1 Machine Learning Expertise:

Developing an ML model for lung cancer detection requires skilled professionals proficient in deep learning, medical imaging, and model interpretability.

8.3.2 Domain Knowledge:

Collaboration with radiologists, oncologists, and medical experts is essential to ensure clinical relevance and proper annotation of training data.

8.3.3 Regulatory Knowledge:

Teams must have expertise navigating healthcare regulations, including medical device certifications and data privacy laws.

8.4 Additional Constraints

8.4.1 Data Access:

- Limited access to large, diverse, and high-quality datasets for training can be a bottleneck.
- Constraints may arise due to ethical concerns, privacy regulations, and costs associated with acquiring data.

8.4.2 Time Constraints:

Regulatory approval processes and clinical trials can significantly extend development timelines.

9. Business model

A well-designed business plan emphasizes the application of technology to help patients, healthcare providers, and other stakeholders while maintaining sustainability, scalability, and profitability. A business model framework is shown below:

9.1 Value Proposition

For Patients:

- Early and accurate detection of lung cancer, leading to improved survival rates and better outcomes.
- Reduced need for invasive diagnostic procedures like biopsies.

For Healthcare Providers:

- Enhanced diagnostic efficiency through automated analysis of medical imaging.
- Reduced workload for radiologists, allowing them to focus on critical cases.
- Improved diagnostic accuracy, minimizing errors and misdiagnoses.

9.2 Target Customers

- **Primary Customers:**
 - Hospitals and healthcare providers.
 - Diagnostic imaging centers.
 - Cancer research organizations and institutes.
- **Secondary Customers:**
 - Pharmaceutical companies (for clinical trials and drug development).
 - Insurance companies interested in reducing costs associated with late-stage cancer treatments.
- **Government and Non-Profits:**
 - Public health programs focused on cancer prevention and early detection.

9.3 Revenue Streams

- **Software-as-a-Service (SaaS):**
 - Subscription-based pricing for access to the ML platform, charged per user, per scan, or month.
- **Licensing:**
 - Licensing the ML technology to medical device manufacturers or healthcare organizations for integration into their systems.

9.4 Scalability and Growth

- **Market Expansion:**
 - Scale the solution to international markets, adapting to local regulations and healthcare needs.
- **AI-Powered Insights:**

Other features include risk prediction, treatment monitoring, and outcome analysis.

10. Concept generation

Concept development is a key initial phase when developing a Machine Learning (ML) system for diagnosing lung cancer. It involves identifying creative solutions

while concentrating on several areas, including data collection, model creation, deployment, and user integration. The main ideas are as follows:

10.1 Data Acquisition and Preparation

- **High-Quality Medical Imaging Data:**
 - Use publicly available datasets like LIDC-IDRI, NLST, and SPIE-AAPM.
 - Collaborate with hospitals and imaging centers to collect anonymized CT scans and X-rays.
- **Data Annotation:**
 - Engage radiologists to annotate datasets with labels for lung nodules, stages, and malignancy.
 - Leverage semi-automated tools to preprocess and augment the data (e.g., resizing, noise removal, contrast enhancement).

10.2 Machine Learning Model Development

- **Model Architecture:**
 - Implement convolutional neural networks (CNNs) for feature extraction from 2D X-rays and CT slices.
 - Use 3D CNNs or transformers for volumetric CT data analysis to capture spatial features.
- **Multi-Modal Analysis:**
 - Combine imaging data with patient demographics, genetic information, and biomarkers for a holistic analysis.
- **Ensemble Learning:**
 - Develop ensemble models to aggregate predictions from multiple ML architectures for improved accuracy.

10.3 Diagnostic Features

- **Nodule Detection:**
 - Automate the detection of lung nodules using region-based CNNs or YOLO models.
- **Nodule Classification:**
 - Classify nodules as benign or malignant using supervised learning with labeled datasets.
- **Cancer Staging:**

- Predict cancer stages using machine learning models trained on multi-class datasets.
- **Risk Assessment:**
 - Predict the likelihood of cancer development based on imaging and patient history.
 - Predict cancer stages using machine learning models trained on multi-class datasets.

10.4 User Experience and Interaction

- **User-Friendly Interfaces:**
 - Develop intuitive dashboards for radiologists to upload scans, view results, and access insights.
- **Decision Support System:**
 - Provide actionable recommendations (e.g., "Refer to oncology," "Schedule follow-up scan") alongside predictions.
- **Alerts and Notifications:**
 - Generate alerts for high-risk cases to prioritize radiologist attention.
- **Collaboration Tools:**
 - Enable team-based review and annotations for consensus building in diagnostics.

11. Final product prototype

This app aims to promote early cancer detection, and it is very user-friendly. The key components of this app are listed below:

1. User Categories

- **Radiologists/Doctors:** Upload medical images (CT scans, X-rays) and view diagnostic results.
- **Technicians:** Assist in data preparation, input, and workflow management.
- **Administrators:** Monitor system usage, manage user accounts, and access logs.

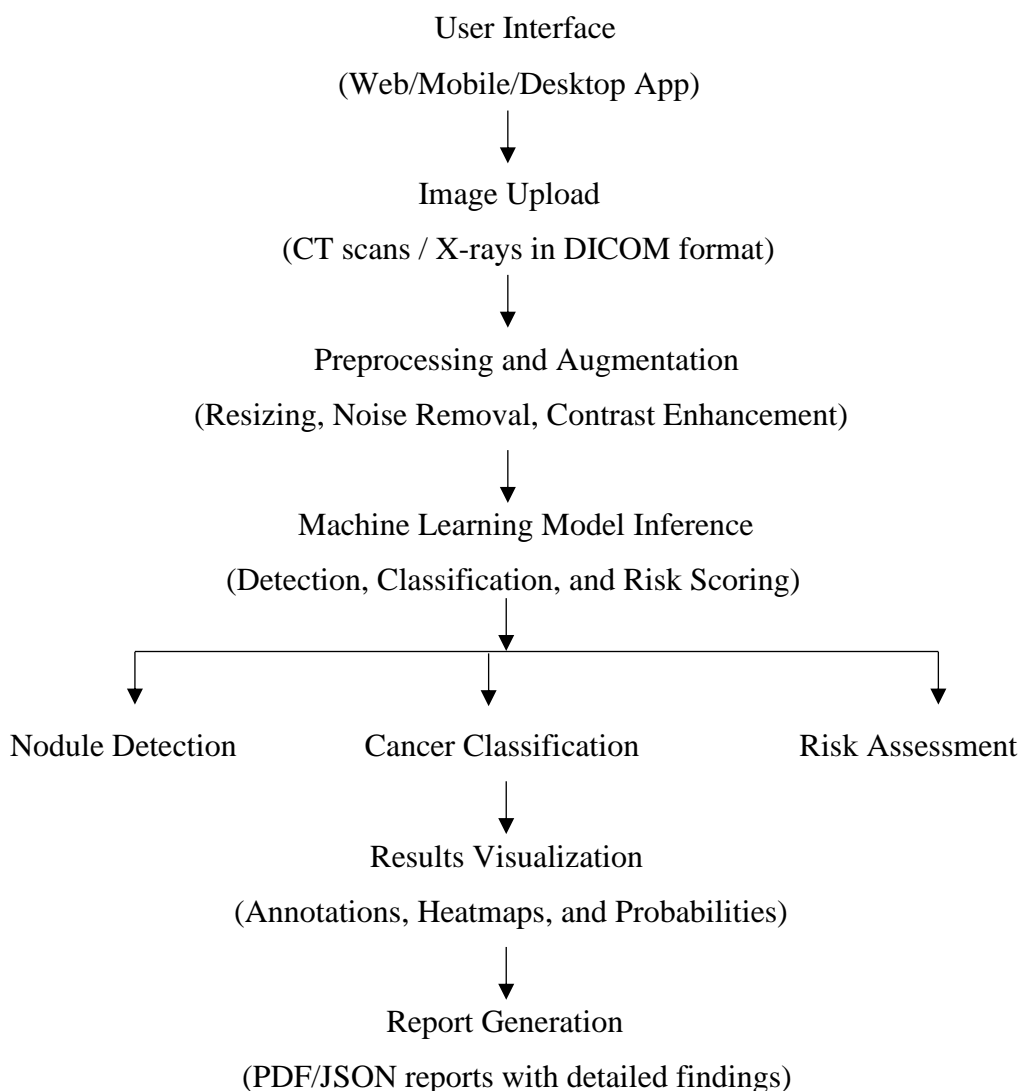
2. Key Functionalities

- **Image Upload:**
 - Upload CT scans or X-ray images in DICOM format.
- **Preprocessing:**
 - Automatically preprocess images (e.g., resizing, denoising).
- **Lung Nodule Detection:**
 - Detect and highlight nodules with bounding boxes or segmentation.

- **Cancer Classification:**
Classify nodules as benign or malignant with a probability score.
- **Risk Assessment:**
Provide a patient-specific risk score based on imaging and other clinical data.
- **Results Visualization:**
Heatmaps and annotations for detected regions to aid decision-making.
- **Report Generation:**
Auto-generate detailed diagnostic reports, including predictions and risk analysis.
- **Real-Time Notifications:**
Alert doctors for high-risk cases requiring urgent attention.

Schematic Diagram

Below is a high-level schematic diagram illustrating the architecture and workflow of the app:





Integration with PACS/EHR System
(Seamless hospital/clinic workflow integration)

3. Prototype design simulation

The app could consist of the following UI components:

Dashboard

- Overview of recently analyzed cases, pending scans, and alerts for high-risk cases.

Image Upload Page

- Drag-and-drop interface to upload CT scans or X-ray images.
- Options to preprocess images before analysis.

Analysis Results

- Visual display of the scan with detected nodules, heatmaps, and classification labels.
- Probability scores for malignancy and stage of cancer.

Patient Report

- A summary of findings with visual annotations.
- Risk assessment and suggested next steps.

Admin Panel

- Manage users, track usage statistics, and monitor model performance.

12. Product details

12.1 How does it work?

(a) User interaction

Log In:

- Users (e.g., radiologists, technicians, administrators) log in to the app with their credentials.
- Secure access is ensured via two-factor authentication (optional).

Image Upload:

- The user uploads medical images (e.g., CT scans or X-rays) in standard formats like DICOM or JPEG.
- Images can be uploaded in batches to process multiple cases simultaneously.

(b) Image preprocessing

Automated Preprocessing:

- Uploaded images are preprocessed to enhance quality:
 - Noise Reduction: Removes unwanted image noise.
 - Resizing and Cropping: Ensures images meet the required input size for the machine learning model.

- Contrast Enhancement: Improves visibility of nodules and other critical features.

(c) Machine Learning analysis

Nodule Detection:

- The system uses a trained machine learning model (e.g., CNNs, 3D CNNs) to identify and localize potential nodules in the lung regions.
- Detected nodules are highlighted with bounding boxes or segmentation masks.

Classification:

- The app classifies each detected nodule as benign or malignant based on the trained model's output.
- Probability scores (e.g., "90% malignant") are generated for each prediction.

Cancer Staging:

- If applicable, the app predicts the stage of lung cancer (e.g., early-stage, advanced-stage) using additional patient data and advanced algorithms.

Risk Assessment:

- The app calculates a comprehensive risk score for the patient, considering imaging results, demographic data (e.g., age, smoking history), and clinical history.

(d) Result visualization

- Heatmaps showing regions of high attention by the ML model.

(e) Report generation

- Detection results (e.g., number and location of nodules).
- Classification and probability scores.
- Reports are generated in PDF format for easy sharing.

12.2 Data Sources

- LIDC-IDRI (Lung Image Database Consortium and Image Database Resource Initiative)
- NLST (National Lung Screening Trial)
- SPIE-AAPM Lung CT Challenge Dataset
- TCGA-LUAD and TCGA-LUSC

12.3 Algorithm

Convolutional Neural Networks (CNNs):

- Used for feature extraction and classification from 2D and 3D medical images (e.g., CT scans, X-rays).

- **Popular Architectures:**

ResNet, EfficientNet, DenseNet, Inception.

- 3D CNNs:**

Ideal for volumetric CT data to capture spatial features.

- Object Detection Algorithms:**

- **For detecting lung nodules in CT scans:**

YOLO (You Only Look Once), Faster R-CNN, Mask R-CNN.

- Segmentation Algorithms:**

- **For delineating lung regions or nodules:**

UNet, SegNet, 3D UNet, V-Net.

- Classification Algorithms:**

- **For classifying nodules as benign or malignant:**

Random Forest, Support Vector Machines (SVM), Gradient Boosting (XGBoost, LightGBM).

- Transformers:**

Vision Transformers (ViT) for high-performance feature extraction and classification in medical imaging.

Frameworks

Deep Learning Frameworks

1. **TensorFlow:**

- End-to-end open-source platform for training and deploying machine learning models.
- Includes TensorFlow Extended (TFX) for production pipelines.

2. **PyTorch:**

- Highly flexible and user-friendly framework, widely used for research and production in deep learning.

3. **Keras:**

- High-level API built on TensorFlow, ideal for rapid prototyping.

Software Tools

Data Management and Annotation

1. **Label Studio:**

- Open-source tool for labeling medical images (bounding boxes, segmentation masks).

2. **DICOM Viewers:**

- Tools like OsiriX, Horos, or 3D Slicer for viewing and annotating DICOM files.

3. PACS Integration:

- Open-source PACS like Orthanc for managing medical image data.

Development Environments

1. Jupyter Notebook:

- Interactive environment for developing and debugging machine learning models.

13. Team required to develop

a. Machine Learning Engineers

- Develop, train, and fine-tune machine learning models for lung cancer detection.
- Implement advanced algorithms (e.g., CNNs, 3D CNNs, Transformers).

b. Data Scientists

- Analyse and preprocess medical imaging data (CT scans, X-rays).
- Perform exploratory data analysis (EDA) and feature engineering.

c. Software Engineers

- Develop the app interface (web, mobile, or desktop).
- Ensure seamless integration of the ML models with the app.

d. DevOps Engineers

- Set up cloud infrastructure for model training and deployment.

14. What does it cost?

Its costs vary depending on the workload.

15. Conclusion

A revolutionary step toward improving early diagnosis, increasing patient outcomes, and accelerating healthcare procedures is the creation of a machine learning-based lung cancer detection system. Such a system can address important issues including late-stage diagnosis, diagnostic mistakes, and radiologist workload by utilizing sophisticated algorithms, extensive medical imaging datasets, and clinical expertise.

Key outcomes of this project include:

1. **Enhanced Diagnostic Accuracy:** Machine learning models provide consistent and reliable analysis, detecting subtle patterns that may be missed by human interpretation.
2. **Early Detection:** Automated systems enable early identification of malignant nodules, significantly increasing the chances of successful treatment.

3. **Improved Efficiency:** The system reduces the time and effort required for manual image review, freeing radiologists to focus on complex cases.
4. **Scalable and Accessible Solutions:** Cloud-based deployment and edge AI ensure the system can be scaled globally, including in resource-constrained settings.

In conclusion, by combining state-of-the-art technology with medical knowledge, this initiative has the potential to completely transform the detection of lung cancer, ultimately saving lives, lowering healthcare expenses, and advancing customized treatment. To make sure the system provides value throughout the healthcare ecosystem, the following steps include clinical trials, regulatory approvals, and iterative enhancements based on feedback from the real world.

References

1. Kadir, T., & Gleeson, F. (2018). Lung cancer prediction using machine learning and advanced imaging techniques. *Translational lung cancer research*, 7(3), 304.
2. Raoof, S. S., Jabbar, M. A., & Fathima, S. A. (2020, March). Lung Cancer prediction using machine learning: A comprehensive approach. In *2020 2nd International conference on innovative mechanisms for industry applications (ICIMIA)* (pp. 108-115). IEEE.
3. Willaime, J. M. Y., Pickup, L., Boukerroui, D., Talwar, A., Gooding, M., Gleeson, F. V., & Kadir, T. (2016, March). Impact of segmentation techniques on the performance of a CT texture-based lung nodule classification system. European Congress of Radiology-ECR 2016.
4. Gollapudi, S. K. S., Bathula, M. K., Muthuluru, M., Sathi, H. S. V., & Pravallika, P. (2024, January). Investigation of machine learning algorithms for the pre-estimating medical diagnosis. In *AIP Conference Proceedings* (Vol. 2512, No. 1). AIP Publishing.