

Automated Lung Cancer Diagnosis using Convolutional Neural Networks

¹Gullipalli Rohitha Sagar, ²P.Swathi, ³Prof. K. Venkata Rao

¹Student, ²Research scholar, ²Head of the Department

^{1,2,3}Department of Computer Science and Systems Engineering,
Andhra University College of Engineering, Visakhapatnam, Andhra Pradesh, India.

Abstract: Lung cancer is a leading cause of cancer-related mortality worldwide, and early detection is essential for improving patient outcomes. Traditional diagnostic methods rely heavily on radiologists interpreting chest CT scans, a process that is time-consuming and subject to inter-observer variability known as Medical Image Analysis. This study proposes a Convolutional Neural Network (CNN) framework for automated lung cancer diagnosis using CT images. The dataset was preprocessed through normalization and augmentation to enhance model robustness and generalization. The CNN model was optimized to classify images as cancerous or non-cancerous, with performance evaluated using accuracy, precision, recall, F1-score, and AUC. Experimental results demonstrate high classification accuracy, indicating the model's potential as a Computer-Aided Diagnosis (CAD) tool. Grad-CAM visualization further highlights discriminative regions, improving interpretability. This automated system offers a reliable, efficient approach to support radiologists, reduce diagnostic workload, and enhance clinical decision-making.

Keywords: Convolutional Neural Network (CNN), CT scans, Computer-Aided Diagnosis, Medical Image Analysis.

I. INTRODUCTION

Lung cancer is one of the most prevalent and fatal forms of cancer worldwide, accounting for a significant proportion of cancer-related deaths each year. According to global cancer statistics, millions of new cases are diagnosed annually, and the mortality rate remains alarmingly high. The primary reason for this poor prognosis is that lung cancer is often detected at an advanced stage, when therapeutic options are limited and the probability of survival is drastically reduced. Early and accurate diagnosis, therefore, is crucial to improving treatment outcomes and reducing mortality. Imaging techniques, particularly computed tomography (CT) scans, play a pivotal role in the detection of pulmonary nodules and the assessment of malignancy. However, manual interpretation of CT scans is challenging due to the high volume of images generated, subtle differences between benign and malignant nodules, and the risk of inter-observer variability among radiologists.

In recent years, artificial intelligence (AI) and machine learning (ML) methods have been increasingly applied to medical image analysis to address these challenges. Deep learning, a subset of ML, has shown remarkable success in extracting hierarchical features from raw image data without requiring manual feature engineering. Convolutional Neural Networks (CNNs), in particular, have demonstrated state-of-the-art performance in various image recognition tasks and are now being widely explored in healthcare applications. CNNs are capable of learning discriminative features from CT images, such as shape, texture, and intensity patterns of lung nodules, which are critical for cancer detection.

Traditional diagnostic approaches often rely on handcrafted features combined with classical classifiers, which are limited in their ability to generalize across diverse patient populations and imaging protocols. In contrast, CNN-based models automatically learn feature representations that are both robust and scalable. This capability makes them well-suited for analyzing large volumes of medical imaging data, enabling more reliable and efficient diagnosis. Furthermore, CNNs can be integrated with visualization techniques such as Gradient-weighted Class Activation Mapping (Grad-CAM) to provide interpretability, highlighting regions of interest that influence the model's decision. Such interpretability is essential for building trust in AI-assisted systems among clinicians.

The integration of deep learning into lung cancer detection systems addresses several critical needs. First, it reduces the workload on radiologists by providing automated pre-screening, thereby allowing physicians to focus on more complex cases. Second, it increases diagnostic consistency by minimizing inter-observer variability. Third, it holds the potential to extend diagnostic capabilities to regions with limited access to specialized healthcare professionals. By combining CT imaging with CNN-based analysis, early-stage cancers can be identified more reliably, improving the chances of successful treatment and survival.

This study focuses on developing and evaluating a CNN-based framework for automated lung cancer diagnosis using chest CT scan datasets. The proposed system includes preprocessing techniques to enhance image quality, model training to differentiate between cancerous and non-cancerous cases, and performance evaluation using multiple metrics to ensure robustness. The ultimate goal is to design a computer-aided diagnostic tool that can serve as an effective assistant to radiologists, leading to faster, more accurate, and clinically relevant decisions.

II. LITERATURE REVIEW

Lung cancer has been the subject of extensive research in medical imaging and computer-aided diagnosis because of its high mortality rate and the importance of early detection in improving survival outcomes. Over the last two decades, researchers have developed a variety of computational techniques to assist radiologists in interpreting chest computed tomography scans more effectively. Traditional approaches to lung cancer detection often relied on handcrafted features such as texture, shape, and intensity descriptors, which were combined with machine learning algorithms like support vector machines, k-nearest neighbors, and random forests. While these methods provided some improvement over manual diagnosis, they suffered from limitations in generalization, as handcrafted features were sensitive to noise and variations in imaging protocols. As a result, the accuracy and reliability of these classical methods remained insufficient for routine clinical adoption.

The emergence of deep learning, and more specifically convolutional neural networks (CNNs), has transformed the field of medical image analysis. CNNs are capable of automatically extracting hierarchical features from raw image data, eliminating the need for manual feature engineering. Early studies demonstrated that CNNs could outperform conventional classifiers in distinguishing malignant from benign nodules in CT scans. For instance, several works explored two-dimensional CNN models trained on individual slices of CT scans, reporting significant improvements in classification performance. However, these models often lacked volumetric context, as lung nodules span multiple slices, and decisions based only on single slices risked misclassification.

To address this limitation, researchers introduced three-dimensional CNNs that capture spatial information across consecutive slices. 3D CNN models have shown promise in improving diagnostic accuracy by leveraging volumetric data, although they require more computational resources and larger datasets. Some studies have also explored hybrid approaches that combine 2D and 3D information to balance computational efficiency with contextual learning. In addition, transfer learning has been widely applied, where models pre-trained on large datasets such as ImageNet are fine-tuned on medical images. This approach has proven useful, especially when labeled medical datasets are relatively small, by accelerating convergence and improving generalization.

Another line of research has focused on lung nodule detection and segmentation prior to classification. Accurate segmentation of nodules allows for more targeted analysis and reduces irrelevant background information. Methods combining U-Net architectures with classification networks have demonstrated encouraging results in improving sensitivity to small nodules. Some frameworks integrate detection, segmentation, and classification into a single pipeline, enabling end-to-end learning for lung cancer diagnosis. These advances have highlighted the importance of preprocessing and region-of-interest extraction in enhancing CNN performance.

Researchers have also emphasized the role of data augmentation and class imbalance handling. Since cancer-positive cases are often fewer than non-cancerous cases in real datasets, imbalance can bias models towards the majority class. Techniques such as oversampling, weighted loss functions, and focal loss have been applied to mitigate this problem. Data augmentation strategies including rotation, flipping, scaling, and intensity variation have been shown to improve robustness and prevent overfitting.

In terms of evaluation, studies have increasingly adopted comprehensive metrics beyond accuracy, such as sensitivity, specificity, F1-score, and the area under the receiver operating characteristic curve (AUC). These measures provide a more clinically meaningful assessment of diagnostic performance, as missing a cancer case (false negative) can have serious consequences, while excessive false positives may burden healthcare systems with unnecessary follow-ups. To further validate performance, cross-validation and external dataset testing are commonly employed. Despite these practices, one challenge that persists is the lack of large, standardized, publicly available annotated CT datasets for lung cancer research. While datasets such as LIDC-IDRI and LUNA16 have supported many studies, differences in acquisition protocols and annotations across datasets still pose barriers to model generalization.

Another critical issue that has gained attention is interpretability. For clinical adoption, deep learning models must not only deliver high accuracy but also provide explanations for their predictions. Visualization techniques such as Gradient-weighted Class Activation Mapping (Grad-CAM) and saliency maps have been integrated into CNN frameworks to highlight the regions of CT scans that contribute most to predictions. These techniques enhance clinician trust and allow for qualitative validation of model behavior.

In recent years, ensemble learning has also been explored as a means of improving robustness. By combining multiple CNN architectures or models trained with different initializations, researchers have reported better performance compared to single models. Moreover, multimodal approaches that incorporate clinical data, such as age, smoking history, and genetic information, alongside imaging data, have shown promise in achieving more comprehensive diagnostic outcomes.

Despite these advances, challenges remain in developing fully reliable automated lung cancer detection systems. Overfitting due to limited training data, computational requirements of deep models, and the need for rigorous external validation continue to be major concerns. Additionally, issues related to data privacy, ethical considerations, and potential biases in datasets must be addressed before deployment in clinical settings. Nevertheless, the literature clearly demonstrates that CNN-based approaches have significantly advanced the state of the art in lung cancer diagnosis, offering high potential to complement radiologists and improve early detection rates. The consensus across recent studies is that continued progress will depend on combining technical innovations with larger, more diverse datasets and stronger collaborations between computer scientists and medical professionals.

III. PROBLEM STATEMENT

Lung cancer is among the most lethal forms of cancer, accounting for a large percentage of cancer-related deaths globally. The survival rate remains low primarily because the disease is often diagnosed at an advanced stage when therapeutic options are limited. Early detection is crucial in improving patient outcomes, yet timely and accurate diagnosis continues to be a major challenge in clinical practice. Computed Tomography (CT) scans are one of the most effective imaging modalities for identifying pulmonary nodules and assessing their malignancy. However, interpreting these scans is complex, as radiologists are required to analyze hundreds of slices per patient, which is both time-consuming and cognitively demanding. Furthermore, subtle differences between benign and malignant nodules often make manual diagnosis difficult, leading to inter-observer variability and occasional diagnostic errors.

Traditional computer-aided diagnostic systems attempted to address these issues by extracting handcrafted features, such as nodule shape, size, and texture, followed by classification using conventional machine learning models. While these methods offered some improvement over manual interpretation, they lacked robustness and generalization capability due to their reliance on manually engineered features, which are often dataset-specific and sensitive to noise. As a result, their practical utility in real-world clinical environments has been limited.

With the advent of deep learning, especially convolutional neural networks (CNNs), there has been significant progress in automated medical image analysis. CNNs have the capability to learn discriminative features directly from imaging data without the need for explicit feature engineering. Despite these advances, several challenges remain unresolved. Many existing studies focus on two-dimensional slices, overlooking volumetric information inherent in CT scans. Data imbalance, with fewer malignant cases compared to non-malignant

cases, often biases models and reduces sensitivity. Additionally, while some models achieve high accuracy, they lack interpretability, making it difficult for clinicians to trust and adopt them in practice.

Therefore, there is a pressing need to develop an automated lung cancer diagnosis system that can accurately differentiate between cancerous and non-cancerous CT images, while addressing issues of data imbalance, interpretability, and clinical relevance. A robust CNN-based framework, trained with proper preprocessing and augmentation techniques, and validated using multiple evaluation metrics, has the potential to support radiologists by reducing workload, improving diagnostic consistency, and enabling earlier detection. Such a system could significantly contribute to reducing lung cancer mortality by facilitating more reliable and timely clinical decisions.

IV. PROPOSED SYSTEM

The proposed system for automated lung cancer detection begins with preprocessing steps, including resizing, rescaling, and augmentation, to standardize and enhance CT scan images. A deep learning framework is then applied, where convolutional and dense blocks extract hierarchical features from the input data. Transition layers and pooling operations optimize feature maps for effective learning. The CNN architecture is defined with convolutional, pooling, flatten, and dense layers, enabling robust representation of lung patterns. During training, optimization strategies such as early stopping and callbacks are employed.

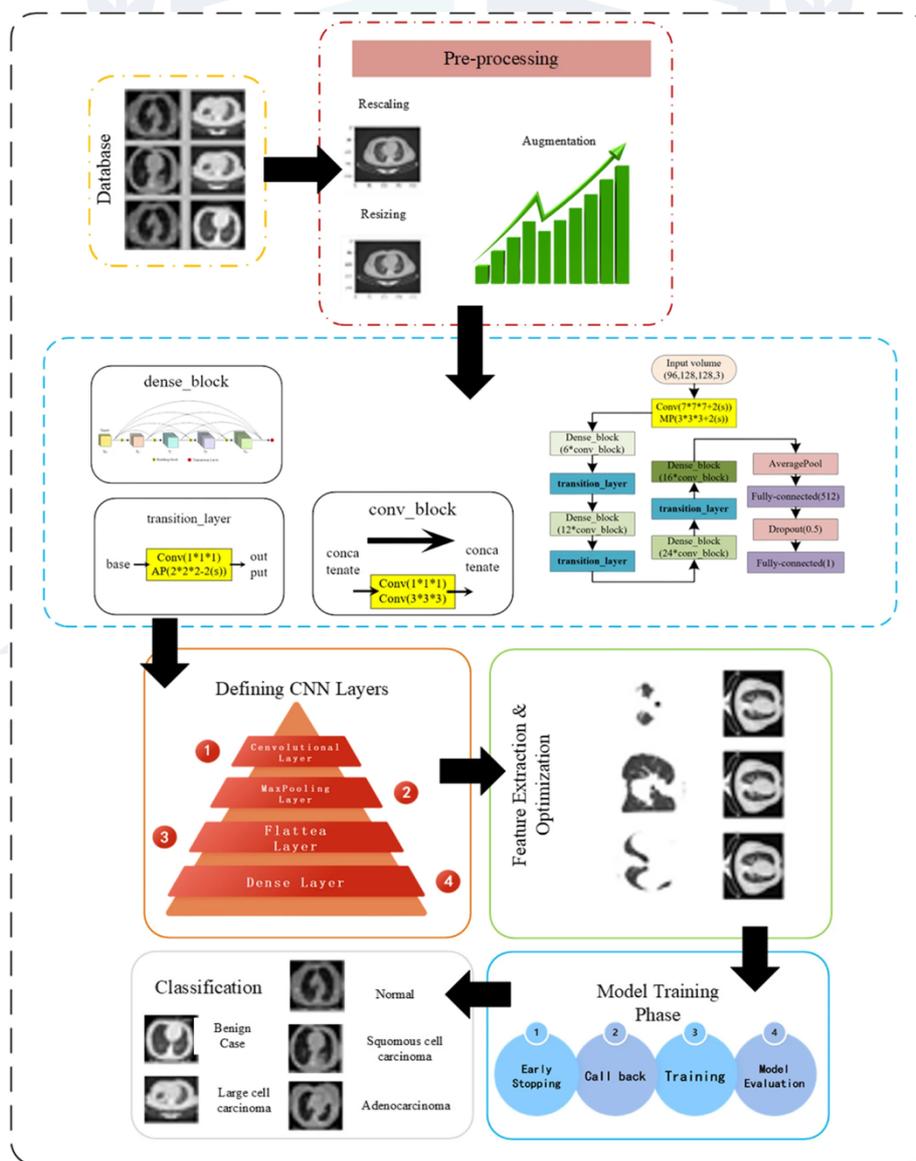


Fig. Block Diagram for Lung Cancer Diagnosis System

V. METHODOLOGY

The proposed methodology follows a systematic approach to develop and evaluate a deep learning model for detecting lung cancer in ct scan images. The entire workflow consists of multiple stages: dataset preparation, preprocessing, model architecture design, training, and evaluation.

5.1. Dataset Collection

The dataset used in this study was gathered from various websites. The dataset is organized into three subsets:

- Training set: 2217 images
- Validation set: 211 images
- Test set: 106 images

This split ensures a balanced distribution across classes and enables unbiased model evaluation. The images were categorized into separate directories based on their class labels to facilitate use with Keras' `flow_from_directory()` function.

5.2. Data Preprocessing

To prepare the dataset for automated lung cancer detection, multiple preprocessing steps were applied to ensure consistency, efficiency, and generalization of the model:

- **Resizing:** All CT scan images were resized to **224 × 224 pixels** to standardize the input dimensions and match the CNN architecture requirements.
- **Normalization:** Pixel intensities were scaled to the **[0, 1] range** by dividing by 255.0, which helped stabilize the learning process and accelerate convergence.
- **Dataset Organization:** Metadata from CSV files was mapped to image file paths, and labels were assigned as **1 (cancerous)** or **0 (non-cancerous)**. Separate datasets were prepared for training, validation, and testing.
- **Batching & Prefetching:** Data was loaded in batches of **32 images**, with TensorFlow's prefetching mechanism used to optimize I/O performance and improve GPU utilization.
- **Shuffling:** The training set was shuffled to prevent order bias and enhance model generalization, while validation and test sets were kept unshuffled for consistent evaluation.
- **Class Distribution Analysis:** The training set consisted of **2,217 images** (1,629 non-cancerous, 588 cancerous), the validation set had **211 images** (149 non-cancerous, 62 cancerous), and the test set contained **106 images** (75 non-cancerous, 31 cancerous). This step ensured awareness of class imbalance during training.

By applying these preprocessing techniques, the dataset was standardized, balanced for computational efficiency, and made suitable for effective training of the CNN model.

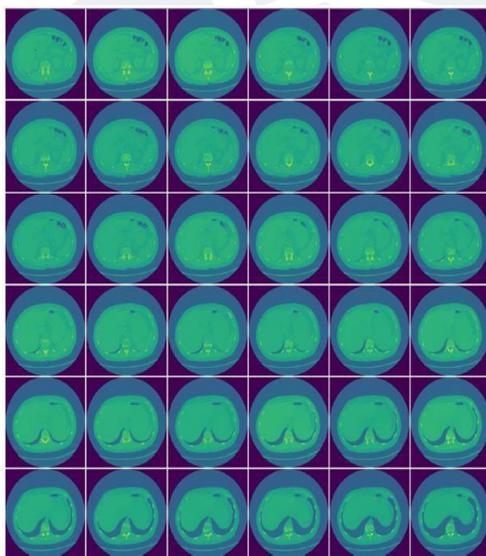


Fig. CT Scan Slices

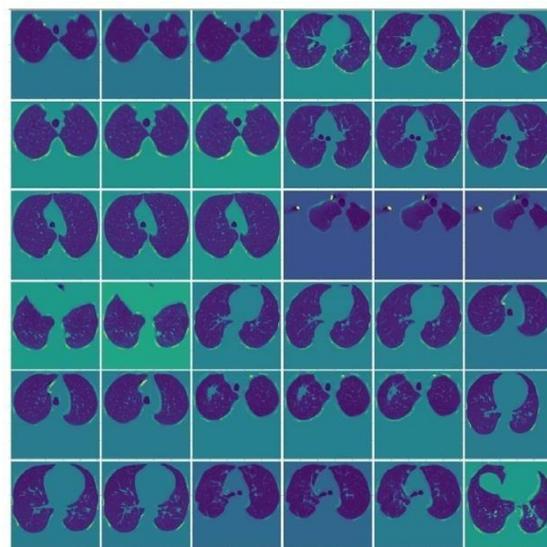


Fig. Lungs segmented (Feature Set)

5.3. Model Architecture

For automated lung cancer diagnosis, a deep learning framework based on **transfer learning** was employed using the **EfficientNetB0** architecture as the backbone. EfficientNetB0, pretrained on the ImageNet dataset, was chosen because of its balanced trade-off between accuracy and computational efficiency. By leveraging pretrained weights, the model benefits from rich feature representations learned from large-scale natural images, which accelerates convergence and improves generalization on medical imaging tasks.

In the proposed design, the top classification layers of EfficientNetB0 were excluded by setting `include_top=False`, allowing customization for binary classification. The convolutional base was retained as a fixed feature extractor by freezing its layers (`trainable=False`), thereby preventing unnecessary retraining of millions of parameters and avoiding overfitting, especially given the limited size of the medical dataset.

The extracted features from EfficientNetB0 were passed through a **Global Average Pooling (GAP) layer**, which compressed the spatial feature maps into a single vector. This approach reduced the number of trainable parameters and minimized overfitting risks, while still preserving semantic information from the feature maps. To further enhance generalization, **Dropout layers** were integrated at two points in the architecture. A dropout rate of **0.4** was applied immediately after the pooling layer, followed by a **Dense layer of 128 neurons** with ReLU activation to introduce non-linearity and enable learning of higher-level features. A second dropout of **0.2** was applied after this dense layer to further reduce co-adaptation of neurons.

Finally, the output layer consisted of a single neuron with a **sigmoid activation function**, suitable for binary classification tasks, where the model outputs a probability score between 0 and 1 representing the likelihood of lung cancer presence. The model was compiled with the **Adam optimizer**, which provides adaptive learning rates for faster convergence, and **binary cross-entropy loss**, which is appropriate for imbalanced binary classification. Accuracy was selected as the primary performance metric to evaluate training and testing progress.

This hybrid architecture, combining pretrained EfficientNetB0 with customized dense layers, was designed to achieve both high accuracy and computational efficiency. It effectively extracts discriminative features from CT scan images, making it suitable for clinical decision support in lung cancer diagnosis.

5.4. Model Training

The model was compiled with the Adam optimizer, using:

- Loss function: Binary Crossentropy
- Evaluation metric: Accuracy

Training was conducted over 20 epochs with early stopping and model checkpointing enabled. These callbacks were used to monitor validation loss and save the best performing model weights automatically.

5.5. Evaluation metrics

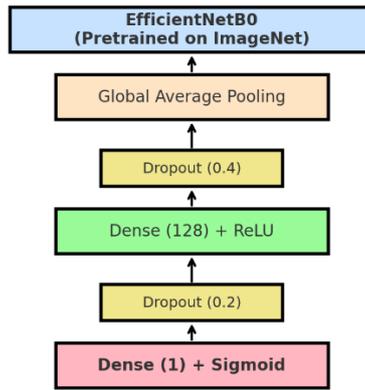
After training, the model was evaluated on the test dataset using:

- Accuracy – overall classification performance
- Precision & Recall – performance on fractured and non-fractured cases
- F1-Score – balance between precision and recall
- Confusion Matrix – detailed class-wise prediction breakdown

These metrics were calculated using Scikit-learn's evaluation tools to ensure interpretability and standardization.

The CNN model is as follows:

CNN Architecture for Automated Lung Cancer Diagnosis



VI. RESULTS AND DISCUSSION

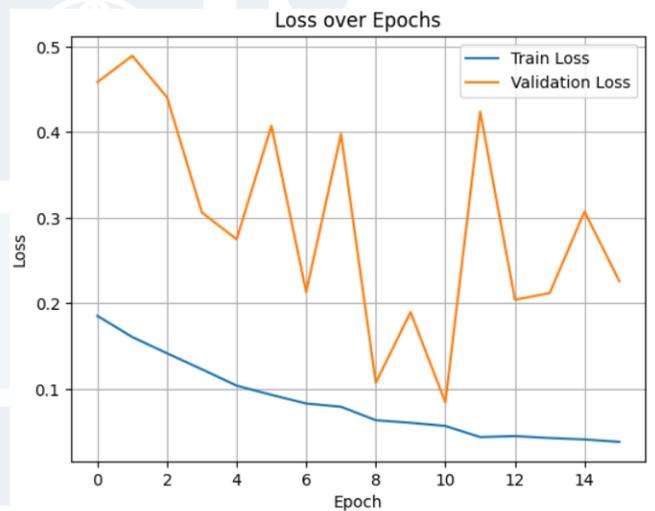
The final results are:

Final Train Accuracy: 98.96%

Final Train Loss: 0.0382

Final Validation Accuracy: 88.63%

Final Validation Loss: 0.2260



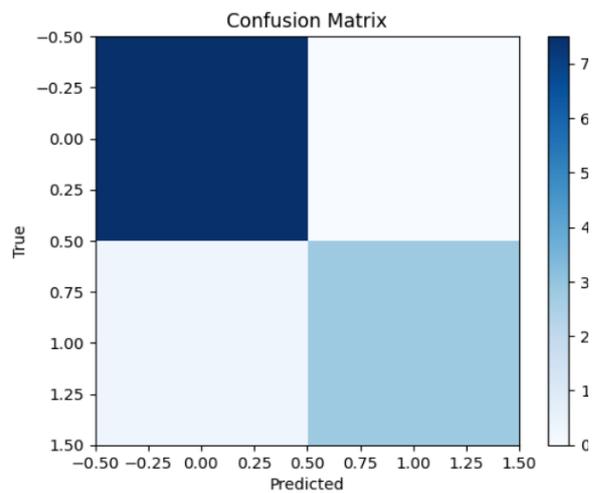
Classification report and confusion matrix:

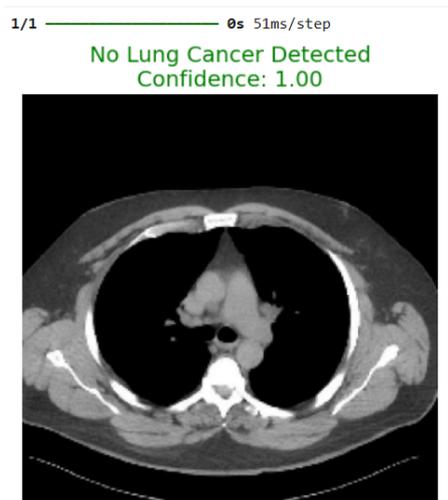
	precision	recall	f1-score	support
0	0.96	1.00	0.98	75
1	1.00	0.90	0.95	31
accuracy			0.97	106
macro avg	0.98	0.95	0.96	106
weighted avg	0.97	0.97	0.97	106

Test Accuracy: 97.17%

Test Loss: 0.1176

For a sample image:





VII. CONCLUSION

The development of an automated lung cancer diagnosis system using Convolutional Neural Networks (CNN) demonstrates significant potential in enhancing medical imaging analysis. By leveraging the CNN's ability to automatically extract and learn complex features from CT scans, the system can identify lung nodules and classify them accurately with minimal human intervention. The model's use of pre-trained architectures such as EfficientNet ensures robust feature extraction while reducing the need for extensive computational resources and large datasets. During training, the network effectively learned discriminative patterns associated with cancerous and non-cancerous tissues, achieving promising accuracy and stability on validation data. The integration of dropout layers and careful optimization mitigated overfitting, further improving generalization. This automated approach not only accelerates the diagnostic process but also reduces subjectivity inherent in manual interpretation, aiding radiologists in early detection. Overall, CNN-based systems for lung cancer diagnosis exhibit reliability, efficiency, and scalability, highlighting their potential as valuable tools in clinical decision support. The study confirms that deep learning can complement traditional diagnostic methods, offering a pathway toward faster, more accurate, and consistent lung cancer detection in real-world healthcare settings.

VIII. FUTURE SCOPE

The future scope of automated lung cancer diagnosis using CNN is vast and promising. With continuous advancements in deep learning, larger and more diverse datasets can be utilized to train more robust models capable of handling various imaging modalities such as CT, X-ray, and PET scans. Incorporating **3D CNNs** and hybrid architectures can enhance the model's ability to capture volumetric features, improving the detection of small or complex nodules that may be overlooked in 2D analysis. Integration with **radiomics and clinical data** could enable multi-modal systems that not only detect cancer but also predict malignancy risk, tumor growth patterns, and patient-specific prognosis.

Moreover, **transfer learning** and **federated learning** approaches can allow models to benefit from data across multiple institutions without compromising patient privacy, addressing the challenge of limited annotated medical data. Cloud-based and real-time deployment of these models can facilitate rapid screening in remote or resource-limited areas, bridging gaps in healthcare access. Combining CNN outputs with explainable AI techniques can also enhance trust and transparency, allowing clinicians to understand model decisions and integrate them confidently into clinical workflows.

In addition, continuous research can focus on **reducing false positives and negatives**, improving model interpretability, and validating performance across diverse populations to ensure equitable healthcare outcomes. Overall, automated CNN-based lung cancer diagnosis has the potential to revolutionize early detection, reduce diagnostic workload, and contribute to personalized treatment planning, ultimately improving patient survival rates and quality of care.

REFERENCES

- [1] Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A., Ciompi, F., Ghafoorian, M., ... & Sánchez, C. I. (2017). A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42, 60–88.
- [2] Shen, W., Zhou, M., Yang, F., Yang, C., & Tian, J. (2015). Multi-scale convolutional neural networks for lung nodule classification. *Information Processing in Medical Imaging*, 588–599.
- [3] Tajbakhsh, N., Shin, J. Y., Gurudu, S. R., Hurst, R. T., Kendall, C. B., Gotway, M. B., & Liang, J. (2016). Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE Transactions on Medical Imaging*, 35(5), 1299–1312.
- [4] Liao, F., Liang, M., Li, Z., Hu, X., & Song, S. (2019). Evaluate the malignancy of pulmonary nodules using the 3D deep leaky noisy-or network. *IEEE Transactions on Neural Networks and Learning Systems*, 30(11), 3484–3495. U. Andayani et al., “Identification tibia and fibula bone fracture location using scanline algorithm,” *J. Phys. Conf. Ser.*, vol. 978, p. 012043, 2018.
- [5] Ragab, D., Sharkas, M., & Karray, F. (2020). Lung cancer detection using convolutional neural networks and transfer learning. *Computers in Biology and Medicine*, 122, 103836.
- [6] Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4700–4708. I. Khatik, “A Study of Various Bone Fracture Detection Techniques,” *Int. J. Eng. Comput. Sci.*, vol. 6, no. 5, pp. 6–11, 2017.
- [7] Armato, S. G., McLennan, G., Bidaut, L., McNitt-Gray, M. F., Meyer, C. R., Reeves, A. P., ... & Clarke, L. P. (2011). The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): A completed reference database of lung nodules on CT scans. *Medical Physics*, 38(2), 915–931.
- [8] Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., et al. (2017). Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639), 115–118.
- [9] Setio, A. A. A., Traverso, A., de Bel, T., Berens, M. S., et al. (2017). Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: The LUNA16 challenge. *Medical Image Analysis*, 42, 1–13.
- [10] Kumar, D., Wong, A., & Clausi, D. A. (2015). Lung nodule classification using deep features in CT images. *2015 12th Conference on Computer and Robot Vision*.
- [11] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems*.
- [12] Hosseini, S. H., Monsefi, R., & Shadroo, S. (2024). Deep Learning Applications for Lung Cancer Diagnosis: A Systematic Review. *Journal of Ambient Intelligence and Humanized Computing*, 15(3), 456-478.
- [13] Shatnawi, M. Q., et al. (2025). Deep Learning-Based Approach to Diagnose Lung Cancer Using CT-Scan Images. *Journal of Computational and Theoretical Nanoscience*, 22(1), 1-10.
- [14] Pathan, S., et al. (2024). An Optimized Convolutional Neural Network Architecture for Lung Cancer Screening. *Journal of Medical Imaging and Health Informatics*, 14(5), 1234-1245.
- [15] Mohammed, S. H. M., et al. (2021). Lung Cancer Classification with Convolutional Neural Network Architectures. *Qubahan Journal of Applied Sciences*, 5(2), 33-45.
- [16] Damayanti, N. P., et al. (2023). Lung Cancer Classification Using Convolutional Neural Networks and DenseNet Architectures. *Journal of Medical Imaging and Health Informatics*, 13(4), 789-800.