

# Lung Cancer Detection using CNN

Mr. Abhay Chavan<sup>1</sup>, Mr. Ajay Jadhav<sup>2</sup>, Mr. Aditya Gite<sup>3</sup>, Mr. Shrinivas Bijja<sup>4</sup> Guide of: Prof. Kirti G. Walke<sup>5</sup>

Undergrad. Student, Dept. of Information Technologies SKN Sinhgad Institute of Technology & Science, Lonavala, Maharashtra

Abhay Chavan, abhaychavan.sknsits.it@gmail.com

Aditya Gite, adityagite.sknsits.it@gmail.com

*Ajay Jadhav*, <u>ajayjadhav.sknsits.it@gmail.com</u> *Shrinivas Bijja*, <u>shrinivasbijja.sknsits.it@gmail.com</u>

Abstract: In recent years, automatic detection of defects in CT scan images has become crucial due to increasing lungrelated issues caused by declining air quality. To effectively treat conditions like lung cancer, thorough examination of CT images, particularly pulmonary nodules, is essential. Machine learning plays a key role in this process, efficiently distributing the workload and processing large datasets to yield accurate results. The diagnostic approach involves three phases: CT image pre-processing, deep learning, and CNN utilization. Pre-processing converts raw data, deep learning assigns data significance, and CNNs determine lung health status, distinguishing between normal and abnormal conditions.

*Keywords:* Feature Extraction CNN, Classification, Modeltraining

#### I. INTRODUCTION

In computer vision, CNNs have greatly advanced, especially in medical imaging. Our research focuses on using CNNs to automatically detect cancer cells in lung tissue WSIs, which is vital for improving accuracy. We begin by extracting the ROI within the WSI to reduce computational load. CNN-based algorithms are then employed to classify image patches into tumor and normal tissue categories. This work is part of the ACDC@LUNGHP project, and we present initial findings here. Notably, our extensive literature review didn't reveal other studies that specifically address CNN-based assessment of lung cancer images. It's worth emphasizing that lung cancer originates in lung cells, and when other cancers metastasize to the lungs, they are treated based on their primary site of origin, such as metastatic breast cancer when breast cancer spreads to the lungs.

## HUMAN LUNG ANATOMY



Fig. 1. Computed tomography images from dataset (left image has nodule, right image has not nodule)

#### Motivation

Research into Lung Cancer Detection using CNNs is driven by the urgent healthcare need to address a significant challenge. Lung cancer is a leading global cause of cancer-related deaths, emphasizing the importance of early detection for better patient outcomes. Conventional methods are often

slow and prone to human error, while CNNs offer the potential for faster and more accurate detection of lung cancer. This can enable early intervention and improved survival rates. The primary goal of this research is to utilize CNNs to create an efficient system for detecting lung cancer in medical images, ultimately advancing medical diagnostics and enhancing patient care.

#### **II.LITERATURE SURVEY**

#### 1 Lung Cancer Detection using CNN :

Several research studies have explored the use of Computer-Aided Diagnosis (CADx) systems to enhance the precision of lung cancer detection. Hua et al. [6] introduced a deep-learning approach, employing a Deep Belief Network (DBN) architecture, resulting in a sensitivity of 73.40% and a specificity of 82.20%. They also experimented with a Convolutional Neural Network (CNN) architecture, yielding a sensitivity of 73.30% and a specificity of 78.70%. Kumar et al. [7] proposed a deep learning system using Stacked Autoencoders (SAE) with an accuracy of 75.01%. Kuruvilla and Gunavathi [8] introduced a technique based on texture features and artificial neural networks (ANN), achieving an impressive accuracy of 93.30%. Gupta and Tiwari [9] developed a technique based on ANN features, delivering an accuracy of 90%. Dandil et al. [10] utilized an ANN and texturebased approach with Principal Component Analysis (PCA), achieving an accuracy of 90.63%.

In a different study, Parveen and Kavitha [5] put forth a method relying on texture characteristics and employed Support Vector Machine (SVM) as the classifier, achieving a sensitivity of 91.38% and a specificity score of 89.56%. Nascimento et al. [11] proposed a technique based on texture properties, utilizing both SVM and Linear Discriminant Analysis (LDA) as classifiers, and achieved an overall accuracy of 92.78%. Orozco et al. [12] presented a technique based on texture characteristics with k-Nearest Neighbors (KNN) as the classifier, obtaining an accuracy rate of 82.66%. Meanwhile, Krewer et al. [13] introduced a methodology combining shape and texture characteristics, resulting in an accuracy of 90.91%.

#### **II.MATERIALS AND METHODOLOGY**

As Illustrated in Figure 1, the proposed system unfolds in three distinct steps. The initial step introduces a method to acquire CT images from the LIDC-IDRI database. In the subsequent phase, nodule segmentation is expertly performed through markings.

Finally, after comprehensive testing and evaluation, the Convolutional Neural Network (CNN) model concludes the diagnosis, distinguishing between benign and malignant tissue.



#### A. Dataset

The LIDC-IDRI database [14] used for this research work is a collaboration between the LIDC and the IDRI, and 1018 tests of CT are included online. Out of 1018 CT tests, 185 tests were not included in the study due to two reasons. The first reason is the examinations display nodules of less than or equal to 3 mm and the second reason is incorrect markings [15],[16]. Consequently, 833 CT exams were used in the proposed methodology The report contains an XML file containing slice contour information and a variety of characteristics including calcifications, texture, and malignancy of one to five values for lung nodules of more than three mm. In two phases, four experts carried out the procedure of recording nodules in LIDCIDRI database. Each expert examined the tests individually in the first step. The results were analyzed by four experts in the second step. At this point, each expert re-examined the examinations and automatically re-examined the annotations [11]. All nodules suggested by the expert review will be considered in this study. This paper only takes account of one case per nodule in order to mitigate the effect of subjectivity on examinations. The benign or malignant classification is achieved first by measurement as shown by Jabon et al. [17], summarizing the characteristics of each nodule calculated by the four specialists. Therefore, in this paper, malignant nodules are those cases with moderately suspicious or very suspicious semantic values of malignancy and benign nodules show high or modest signs of a benign tumor. The value that includes greater boundaries between the four markers during the annotation has been used with respect to the contour. Overall, 1405 three- dimensional nodules were collected (1011 benign and 394 malignant).

I



#### **B.** Segmentation

Data for this study was sourced from an XML file containing nodule coordinates alongside analytical criteria provided by multiple experts for nodule segmentation. In this work, the segmentation approach exemplifies the nodules identified by four expert markings. It's important to note that CT images typically exhibit higher spatial resolution on the x and y-axes, while the z-axis tends to have lower resolution, as highlighted by Hua et al. [6], [18]. The dataset comprises 8296 nodules, consisting of 4329 benign and 3967 malignant cases. Each 2-D CT slice serves as a training sample. For CNN architecture training, CT images are resized to 28×28, as depicted in Figure 2. Lower-dimensional images, shown on the left side of Figure 2, are transformed into 28×28 background images. Largerdimensional nodules, displayed on the right side of Figure 2, are proportionally resized to fit within 28×28 images. Each CT image is then fed into the CNN classification model, with a focus on classifying lung nodules as malignant or benign, without involving calculations of morphological and textural characteristics.

#### Convolutional Neural Networks (CNNs):

CNNs are a fundamental methodology in image classification tasks. The code implements a custom CNN model, which consists of convolutional layers followed by fully connected layers. CNNs are specifically designed for features extraction from images and are widely use in computer vision tasks.



#### II. Proposed Framework

The classification stage is carried out using the Support Vector Machine (SVM) method. Inputs used in this stage are parameter values in the form of cancer area, contrast, energy, entropy, and homogeneity. While the output produced in the classification stage is a decision in the form of normal, benign, or malignant. There are 2 stages in the classification process using SVM, namely training and testing.

#### A. TRAINING PHASE:

The architecture of our Convolutional Neural Network (CNN) is crucial in the study, designed to assess the output of each node based on various parameters, including epoch numbers, learning rate, filter sizes, and the training and test datasets. It plays a pivotal role in distinguishing between malignant and benign lung nodules. In the architecture, C(N) represents the convolutional layer, P(S) signifies the pooling layer with a kernel size of S, CC(K) corresponds to the Fully Connected (FC) layer with Kneurons, while R and S denote the rectified linear unit (ReLU) and softmax function, respectively. The inclusion of a dropout layer (D) serves to prevent overfitting.

Figure 3 visually presents the proposed CNN architecture, which comprises an initial convolutional layer with 32 filters, each utilizing ReLU activation. This is followed by a second convolutional layer with 16 filters, also employing ReLU activation. Subsequently, a pooling layer with a kernel size of 2 is applied, followed by a dropout layer for regularization to mitigate overfitting. Further, there is an FC layer with 16 neurons, ReLU activation, and an additional dropout layer. The integration of these layers and functions forms a robust framework for our CNN, enhancing its capacity to effectively differentiate between benign and malignant lung nodules.

#### **SVM Training:**

After the creation of the .txt file is done, the next step is the training process that is carried out through the command prompt using the existing library in Open CV.

The command is used at the command prompt. The training process through this command prompt generates a file .model that contains a database or a place to store the model parameters studied by SVM-train. This file .model will be used for predictions in the testing process.

#### **B. TESTING PHASE:**

The testing process in OpenCV is accomplished using the SVM::predict library, with Figure 4 illustrating how this process functions. During testing, the input consists of parameter values generated during the feature extraction stage, specifically, parameters such as cancer area, contrast, energy, entropy, and homogeneity derived from CT-Scan images. These parameter values are compared to the parameter values obtained during the training process, which are stored in the database within a .model file.





Fig. How the Testing Process Works

The outcome of this comparison is represented as 0, 1, or 2, where 0 denotes normal, 1 corresponds to malignant lung cancer, and 2 indicates benign lung cancer as described by equation (6):

yi(wxi + b) > 0 for i = 1, 2, ..., n (6)

In this equation, xi represents the input data, Yi signifies the output result, which is either +1 or +2, and w and b denote parameter values. If the data output yi = +1, the result indicates malignant lung cancer, whereas if the data output yi = +2, it signifies benign lung cancer.

#### C. Decision Making:

The last stage is decision making. The decision-making stage is carried out after obtaining the values of the area, contrast, energy, entropy, and homogeneity of the input image which are then matched with the data of the parameter values in the result database of the training process. Decisions resulting from this system can be normal, benign lung cancers, or malignant lung cancers.

#### VI.RESULTS AND DISCUSSION

The proposed work is carried out using seven thousand nodules (3500 of each, malignant and benign) for training toassess the proposed architecture and 1296 (829 benign and 467 malignant) nodules for testing. The model has been tested and trained with a learning rate of 0.01, a kernel size of five, and a kernel size of two, with a pooling layer of 30 to 30 batch training, for 30 epochs. The proposed CNN architecture obtained 97.2% accuracy, 95.6% sensitivity, and 96.1% specificity. That can be due to the inclusion of the convolution layer, which produces further maps of characteristics and the existence in the fully connected layer of a hidden layer.

 TABLE I

 COMPARISON WITH OTHER RESEARCH STUDIES [ACCURACY (ACC), SPECIFICITY (SPE) AND SENSITIVITY (SEN)].

Research Study	Dataset	No. of Samples	Acc. (%)	Spe. (%)	Sen. (%)
[10]	Private	128	90.63	89.47	92.30
[9]	Private	120	90	93.33	86.66
[6]	LIDC	2545		78.70	73.30
[13]	LIDC-IDRI	33	90.91	94.74	85.71
[7]	LIDC	4323	75.01		83.35
[8]	LIDC	110	93.30	100	91.40
[11]	LIDC	73	92.78	97.89	85.64
[12]	NBIA-ELCAP	113	-	52.17	96.15
[5]	Private	3278		89.56	91.38
Proposed Work	LIDC-IDRI	8296	97.2	95.6	96.1

The results of this study and some related studies are compared in Table I. It is important to mention that the same image database, training and testing samples, settings for classifiers along with other parameters must be used for making a consistent comparison with these previous research works. In spite of this, the present study provides better results relative to the papers that employed deep learning methods [6], [7]. The large number, that is, seven thousand nodules (3500 of each, malignant and benign) used for training is due to this fact.

### **VI.CONCLUSION**

Lung cancer remains a global health challenge, leading to a high number of cancer-related deaths and low survival rates. Timely and accurate diagnosis is crucial for patient outcomes. This study's primary focus is on utilizing deep learning, specifically CNN architecture, to address this challenge without delving into complex morphological and textural analyses. The main goal is to effectively classify lung nodules as benign or malignant. Rigorous evaluation of the LIDC-IDRI database yielded impressive results, with 97.2% accuracy, 95.6% sensitivity, and 96.1% specificity. These outcomes outperform alternative learning methods, as demonstrated through comparative assessments. The introduction of ALCDC is poised to be a transformative development in medical diagnosis research and healthcare systems, with the potential for a significant and enduring impact in the field.



#### **VI.REFERENCES**

[1] Abdul, W., Ali, Z., Ghouzali, S., and Alsulaiman, M. (2017). Security and privacy for medical images using chaotic visual cryptography. Journal of Medical Imaging and Health Informatics, 7(6), 1296-1301.

[2] Jabeen, F., Hamid, Z., Akhunzada, A., Abdul, W., and Ghouzali, S. (2018). Trust and reputation management in healthcare systems: Taxonomy, requirements and open issues. IEEE Access, 6, 17246-17263.

[3] Srichai, M. B., Naidich, D. P., Muller, N. L., and Webb, W. R. (2007). Computed tomography and magnetic resonance of the thorax. Lippincott Williams and Wilkins.

[4] Leef, J. L. and Klein, J. S. (2002). The solitary pulmonary nodule. Radiologic Clinics of North America, 40(1):123–143.

[5] Parveen, S. S. and Kavitha, C. (2014). Classification of lung cancer nodules using svm kernels. International Journal of Computer Applications, 95(25).

6] Hua, K.-L., Hsu, C.-H., Hidayati, S. C., Cheng,W.-H., and Chen, Y.-J. (2015). Computer aided classification of lung nodules on computed tomography images via deep learning technique. OncoTargets and therapy, 8:2015–2022.

[7] Kumar, D., Wong, A., and Clausi, D. A. (2015). Lung nodule classification using deep features in ct images. In Computer and Robot Vision (CRV), 2015 12<sup>th</sup> Conference on, pages 133–138. IEEE.