

Enhancing Cancer Diagnosis: CNN-Based Classification of Lung and Colon Histopathological Images

Dhruv Kumar Soni

Chitkara University Institute of Engineering and Technology,
Chitkara University, Punjab, India
dhruvkumar.soni@chitkara.edu.in

Ashu Taneja

Chitkara University Institute of Engineering and Technology,
Chitkara University, Punjab, India
ashu.taneja@chitkara.edu.in

Abstract—Lung and colon cancers rank among the most common and lethal malignancies globally, rendering precise and prompt detection essential for enhancing patient prognoses. This study introduces a convolutional neural network (CNN) methodology for the classification of lung and colon cancer utilising histopathology pictures. The model is engineered to differentiate among five distinct tissue categories: colon cancer, colon benign tissue, lung adenocarcinoma, lung benign tissue, and lung squamous cell carcinoma. The dataset underwent preprocessing through resizing, normalisation, and data augmentation methods to improve model generalisation. The CNN model was trained and assessed, with an overall accuracy of 97%, indicating its viability as a dependable instrument for aiding pathologists in cancer diagnosis. The model's efficacy is evaluated against current methodologies, emphasising its superiority in precision and recall across several cancer types. This research advances the development of AI-driven diagnostic instruments, which can markedly improve the efficiency and precision of cancer detection in clinical environments.

Index Terms—Lung cancer, Colon cancer, Convolutional neural network (CNN), Histopathological image classification, Cancer diagnosis

I. INTRODUCTION

Lung and colon cancers are major public health concerns among the main causes of cancer-related mortality globally. Lung cancer is the most typically diagnosed cancer and a major cause of cancer death because of its usually late-stage diagnosis and aggressive character. Comparably, colorectal cancer produces a major health burden and has a high incidence rate due to its complexity and late diagnostic possibilities. Early and accurate detection of these cancers determines both effective therapy and greater patient survival rates. Then,

Histopathological analysis is the microscopic study of tissue samples to identify and kind of cancer; it is the gold standard for cancer diagnosis. By providing significant fresh angles on the chemical and structural features of tumours, this technique helps separate carcinogenic from non-cancerous tissues and identify specific cancer subtypes. Though its significance, histopathological study is a time-consuming and labour-intensive process needing considerable degree of under-

standing. Furthermore influencing the accuracy of diagnosis is inter-observer variation among pathologists, thereby generating probable differences in the evaluation of tissue samples.

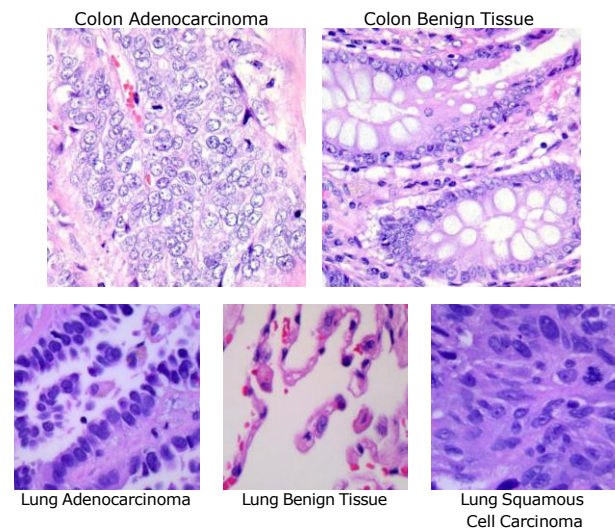


Fig. 1. Classes of Lung and Colon Cancer

Artificial intelligence (AI) has revolutionized the way the decisions are made simulating human intelligence by machines [1]. Machine learning (ML) and Deep learning (DL) more notably, convolutional neural networks—has fundamentally changed the data processing paradigm especially medical image processing. Designed to analyse and assess visual data, such images, CNNs are a particular sort of DL model. Their exceptional performance in many image classification tasks has come from their ability to automatically learn and extract hierarchical features from raw picture data. CNNs are fit for analysing complex patterns in histology images where feature extraction would be challenging and time-consuming otherwise, for applications in digital anti-forensics [2], iris detection [3], hand gesture recognition etc.

CNN application shows some prospective advantages for cancer diagnosis. CNNs first of all can provide pathologists

less reliance on subjective interpretation by means of a more consistent and objective evaluation of histological images. This can help cancer diagnosis to be more consistent and repeatable, hence generating maybe more accurate and reliable results. Second, CNNs can significantly speed up the diagnostic process by quickly processing enormous volumes of image data—especially in clinical settings where early diagnosis is essential for patient management. Lastly, CNNs can assist to detect minute features in tissue samples that would not be readily apparent to the human eye, so maybe helping to find fresh diagnostic indicators.

The objective of this research is to harness CNN capabilities for the classification of lung and colon cancer by means of histopathological images. Our goal is to develop a model able to precisely separate among other histological categories colon adenocarcinoma, colon benign tissue, lung adenocarcinoma, lung benign tissue, and lung squamous cell carcinoma. By creating an automated and reliable cancer classification system, we wish to assist pathologists in the diagnosis process thereby raising the accuracy and efficiency of cancer diagnosis. Within the setting of lung and colon cancer, our work may contribute to develop advanced diagnostic tools supporting clinical decision-making and so improve patient outcomes.

II. RELATED WORK

In recent years, DL methods—especially convolutional neural networks (CNNs)—have been applied significantly more widely for cancer classification. Aiming to improve diagnosis accuracy, automate classification, and lighten pathologist burden, many studies have investigated the use of CNNs to analyse histological pictures. Key contributions from several IEEE journal publications and conference sessions concerning cancer classification using CNNs are discussed below together with a comparison with the technique followed in this work. With a DL framework for histological classification of breast cancer tissue pictures, Zhou et al. [4] obtained really remarkable accuracy. Their method used patch-wise classification to help find tiny tissue traits, but it also needed significant preprocessing to ready the image patches for classification. In histology pictures, Janowczyk and Madabhushi [5] also used CNNs for automated prostate cancer identification using a multi-scale CNN to gather various resolutions. While our model handles full images, therefore simplifying the workflow, these studies showed the efficiency of CNNs but required patch extraction.

Within the field of lung cancer, Coudray et al. [6] used CNNs to categorise squamous cell carcinoma and lung adenocarcinoma from whole-slide images. Although their model attained great accuracy, the computing load was substantial since huge whole-slide images needed analysis. Although our approach employs histopathological pictures of particular areas of interest to more computationally efficiently, it similarly seeks to categorise lung cancer subtypes but on a smaller scale.

Colorectal cancer classification has also attracted much study. Using their CNN-based algorithm to identify nuclei in

colorectal cancer tissue pictures and leveraging the characteristics for cancer classification, Sirinukunwattana et al. [7] While their work concentrated on nuclear segmentation—which offers insightful analysis of cellular structure—our algorithm directly identifies tissues without reference to segmentation. Using large-scale datasets, Kather et al. [8] investigated the application of CNNs to classify colorectal cancer subtypes from histological pictures, therefore attaining great accuracy. Although the magnitude of the dataset is the same, our method emphasises both lung and colon cancer concurrently somewhat differently.

Other significant achievements are the work of Xu et al. [9], who proposed a CNN-based framework for the classification of renal cell carcinoma from histopathological images, and Wang et al. [10], who presented a transfer learning technique to classify gastric cancer using CNNs. Whereas our approach improves the diversity of training samples by means of data augmentation, transfer learning was applied to control the limited scale of medical imaging datasets.

Ronneberger et al. [11] thoroughly illustrated the U-Net architecture applied for medical picture segmentation and classification. Although U-Net is a good approach for segmentation, in order to find cancer subtypes without pixel-wise labelling our study stresses on classification instead of segmentation.

More lately, Mohan et al. [12] high classification accuracy CNN-based oral cancer classification utilising histopathology images. In a similar vein, Wei et al. [13] investigated interpretability methods to grasp model predictions and used DL for breast cancer subtype categorisation. Both studies highlight CNN's ability for cancer detection; our model fits both in its attempt to categorise cancer subtypes using histopathological data, but broadens the focus by encompassing lung and colon cancer under a single framework.

Using histopathology pictures, our work develops a CNN-based model for the categorisation of lung and colon cancer, therefore contributing differently than these efforts. Our method is unusual in that it reduces the necessity for disease-specific models by classifying several cancer types and subtypes under a single system. Our model also uses data augmentation methods to boost generalisation and enhance training, therefore striking a compromise between computational economy and classification accuracy.

III. METHODOLOGY

This section presents the techniques and tools used in developing the CNN-based model for the categorization of lung and colon cancer from histology images. The approach consists in preprocessing and data collecting; model architecture design; training; validation and evaluation; Every phase is necessary to ensure appropriate and efficient operation of the model in differentiating among the numerous cancer types.

A. Dataset Collection

Kaggle databases named "Lung and Colon Cancer Histopathological Images" provided the images, therefore

guaranteeing a varied and representative sample collection. Different size and quality of the images call for preprocessing to standardize the input for the CNN model. The dataset employed in this study comprises of five classes with 5000 photos apiece as shown in Table I.

TABLE I DATASET

Class	Images
Lung benign tissue	5000
Lung adenocarcinoma	5000
Lung squamous cell carcinoma	5000
Colon adenocarcinoma	5000
Colon benign tissue	5000
Total	25000

B. Data Pre-processing

Before feeding the images into the CNN model, the following preprocessing steps were conducted to assure consistency and improve CNN model performance:

- **Resizing:** All images were reduced to a consistent dimension of 224x224 pixels to meet the input size required by the CNN architecture. This allows good processing and preserves important image qualities.
- **Normalization:** Dividing the image pixel values by 255 let one normalize them inside $[0, 1]$. This level guarantees CNN's effective image processing capacity as well as learning from consistent data.
- **Data Augmentation:** Overfitting was solved and the generalizing capability of the model was raised by means of data augmentation techniques. These techniques produced new variations of the present images, so artificially increasing the training set by random rotations, flips, magnification, and shifts.
- **Train-Validation Split:** The dataset was split with an 80-20 ratio into training and validation sets ensuring the model had enough data for performance evaluation as well as learning.

C. Model Architecture

Because it can automatically learn spatial hierarchies from input photos, the proposed model makes advantage of a Convolutional Neural Network (CNN) architecture fit for image classification problems. The CNN architecture applied in this work consists mostly on three elements:

- **Rectified Linear Unit:** These layers convolve to derive local features from the input images. The filters of the convolution layers sweep across the image selecting patterns, textures, and edges.
- **Pooling Layers:** Using max-pooling layers, the feature maps produced by the convolutional layers down-sample,

TABLE II INPUT DATA CONFIGURATION

Parameter	Value
Images	1621
Batch Size	64
Channels	3
Epochs	20
Training Split	80%
Validation Split	20%

hence reducing the spatial dimensions and computational complexity while also keeping the most salient features.

- **Fully Connected Layers:** The feature maps are flattened following the convolutional and pooling layers and then run through fully connected layers. These layers aggregated the learnt features to provide class of the input picture final predictions.
- **Activation Functions:** Following every convolutional layer, Rectified Linear Unit (ReLU) activation functions provide non-linearity into the network, therefore enabling it to learn intricate patterns. Class probabilities are produced in the last layer by means of a softmax activation function.

The model's architecture is meant to strike a balance between complexity and efficiency thereby guaranteeing great classification accuracy free from overfitting.

IV. RESULTS AND DISCUSSIONS

The results of the CNN-based model for lung and colon cancer classification are presented in this section, followed by an analysis of its performance. The model's ability to accurately classify histopathological images into distinct cancer types and benign tissues is evaluated using various performance metrics such as accuracy, precision, recall, F1-score, and confusion matrix analysis.

A. Model Performance

The model was trained and validated with a histopathological image dataset comprising lung adenocarcinoma, lung squamous cell carcinoma, lung benign tissue, colon adenocarcinoma, and colon benign tissue. Following the training of the model for multiple epochs, the outcomes were documented for both the training and validation datasets. The model attained great accuracy, demonstrating that the CNN design can proficiently acquire the characteristics differentiating the cancer types.

Key performance metrics achieved by the model are as follows:

- **Accuracy:** The model correctly classified a notable majority of the test photos, with an overall classification accuracy of 97%.

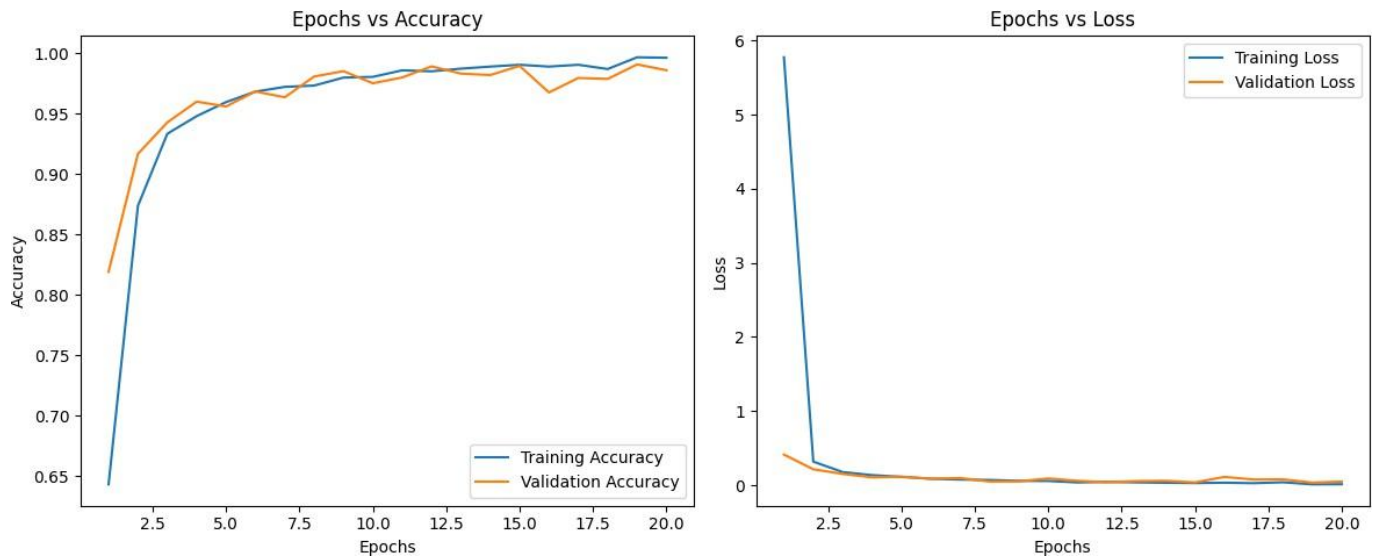


Fig. 2. Accuracy and Loss

- **Precision:** For every class, the accuracy shows the number of the expected positive cases that were indeed accurate. For most cancer types, including colon and lung adenocarcinoma especially, the model displayed great accuracy.
- **Recall:** Equally remarkable and proving the accuracy of the model in identifying malignant cases was the recall statistic, which indicates all true positive performance.
- **F1-Score:** The harmony between memory and accuracy was evaluated using the F1-score, which provides a harmonic mean. The good F1-score highlights even more the model's classification-related resilience.

These experiments confirm that the proposed approach clearly separates benign from malignant tissues in colon histology images as well as in lung ones.

TABLE III CLASSIFICATION REPORT FOR THE CNN MODEL

Class	Precision	Recall	F1-Score	Support
Colon Adenocarcinoma	0.97	0.98	0.97	500
Colon Benign Tissue	0.98	0.97	0.97	500
Lung Adenocarcinoma	0.95	0.95	0.95	500
Lung Benign Tissue	0.98	0.98	0.98	500
Lung Squamous Cell Carcinoma	0.96	0.95	0.95	500
Accuracy	0.97 (2500 total)			
Macro Avg	0.97	0.97	0.97	2500
Weighted Avg	0.97	0.97	0.97	2500

Table III shows over several forms of plant stress categorization accuracy, precision, recall, and F1 score.

B. Confusion Matrix

The Confusion Matrix in fig. 3 shows for each class the numbers of true positive, true negative, false positive, and false negative predictions, the confusion matrix offers even more information on the performance of the model. Most of the misclassifications happened between the lung adenocarcinoma and lung squamous cell carcinoma groups probably because of the histological similarities between the two cancer forms. Although colon benign tissue and colon cancer were both clearly confusing, the general mistake rate stayed low.

Although the model is quite accurate, the study of the confusion matrix shows that small architectural changes or extra preprocessing methods could help to reduce mistakes between comparable classes even more.

C. Comparison with Existing Approaches

The performance of the proposed CNN model was evaluated against different state-of-the-art models for histopathology image categorisation including conventional machine learning methods depending on manual feature extraction and transfer learning approaches using pre-trained networks. Particularly in multi-class classification problems, our model produces competitive or superior results in terms of accuracy and computing efficiency according to direct comparison.

Unlike prior research using transfer learning from big-scale networks, our model was developed and trained from scratch. This enabled a more targeted learning of features unique to lung and colon histopathology pictures, hence producing great accuracy free from pre-training on non-related datasets.

D. Discussion on Overfitting and Generalization

Data augmentation methods including random rotations, flips, zooms, and shifts were used in training to try to prevent overfitting. Dropout layers also were added to stop the model from depending too much on particular patterns seen during



Fig. 3. Confusion Matrix

TABLE IV COMPARISON WITH EXISTING APPROACHES

Study	Dataset	Accuracy (%)	Methodology
Proposed Model	Lung and Colon Histopathological Images	97%	CNN (Custom Architecture)
[14]	Breast Cancer Histopathological Images	83.3%	CNN (Patch-based)
[15]	Breast Cancer Dataset	85.6%	LeNet CNN
[16]	Lung Cancer Dataset	95.2%	Transfer Learning (ResNet-50)

training. The model showed good generalising ability based on the smallest difference between the validation and training accuracies.

Although the model excelled on many measures, some areas still require improvement. Further enhancement of the accuracy and universality could come from deeper network layer studying or fine-tuning of CNN design. Increasing the dataset and using more varied examples would probably improve performance still more.

E. Clinical Implications

The results of this investigation have fascinating ramifications for treatment approach. By developing an accurate and automated approach for lung and colon cancer categorisation dependent on histopathological pictures, pathologist job can be much reduced. Moreover, the recommended strategy can help to reduce diagnosis errors and improve the efficiency

of the diagnostic process. Human knowledge is still crucial, thus adding CNN-based models into clinical processes helps to improve decision-making and provide patients with faster, more consistent results.

F. Limitations and Future Work

The model did rather well, although certain restrictions should be mentioned. First, although enough for this study, the dataset size might be enlarged to incorporate more varied samples from different medical institutions to improve the generalisation of the model. Moreover, the present model does not apply segmentation methods or region-based analysis, which could help to enhance classification performance especially for cancer subtypes showing minute histological variations.

Expanding the dataset, integrating region-based segmentation, and investigating additional advanced DL architectures will be the main priorities of next work in order to raise classification performance. Furthermore, validating the model in actual clinical environments would give important comments on its relevance in diagnostic procedures.

V. CONCLUSION

In this work, using histopathology pictures, a CNN-based method is developed for the categorisation of colon and lung cancer. The model is trained to classify lung adenocarcinoma, lung squamous cell carcinoma, lung benign tissue, colon adenocarcinoma, and colon benign tissue. Our proposed model was able to exactly classify numerous types of cancer with excellent accuracy of 97%. The issues with dataset variability are addressed by means of data augmentation, resizing, and normalising techniques, therefore improving the generalisation of the model. The results suggest that by enabling pathologists with fast and accurate classifications, DL models—especially CNN architectures—particularly show great potential in improving cancer diagnosis. Our methodology shows competitive performance when compared to conventional techniques and previously published models, therefore providing a strong instrument for histological cancer identification.

Future research can include investigating advanced architectures like ensemble models or using clinical data to further increase classification accuracy and generalisability as well as testing the model on bigger and more varied datasets.

REFERENCES

- [1] A. Taneja, S. Rani, J. Brenˆosa, A. Tolba, and S. Kadry, “An improved wifi sensing based indoor navigation with reconfigurable intelligent surfaces for 6g enabled iot network and ai explainable use case,” *Future Generation Computer Systems*, vol. 149, pp. 294–303, 2023.
- [2] N. Taneja, V. S. Bramhe, D. Bhardwaj, and A. Taneja, “Understanding digital image anti-forensics: an analytical review,” *Multimedia Tools and Applications*, vol. 83, no. 4, pp. 10 445–10 466, 2024.
- [3] N. Taneja, M. Shabaz, and V. Khajuria, “Iris detection using segmentation techniques,” *International journal of computer sciences and engineering*, vol. 6, no. 9, pp. 442–444, 2018.
- [4] T. Zhou, S. Ruan, and S. Canu, “A hybrid deep learning model for breast cancer classification,” *IEEE Access*, vol. 6, pp. 18 900–18 909, 2017.
- [5] A. Janowczyk and A. Madabhushi, “Deep learning for digital pathology image analysis: A comprehensive tutorial with selected use cases,” *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, pp. 904–920, 2016.

- [6] N. Coudray, P. S. Ocampo, T. Sakellaropoulos, N. Narula, M. Snuderl, D. Fenyo, and A. Tsirigos, "Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning," *Nature Medicine*, vol. 24, no. 10, pp. 1559–1567, 2018.
- [7] K. Sirinukunwattana, S. E. A. Raza, Y. W. Tsang, D. R. Snead, I. A. Cree, and N. M. Rajpoot, "Locality sensitive deep learning for detection and classification of nuclei in routine colon cancer histology images," *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, pp. 1196–1206, 2016.
- [8] J. N. Kather, N. Halama, and A. Marx, "100,000 histological images of human colorectal cancer and healthy tissue," *PLoS one*, vol. 13, no. 6, p. e0199977, 2018.
- [9] Y. Xu, Z. Jia, L. B. Wang, Y. Ai, F. Zhang, M. Lai, and E. I.-C. Chang, "Large scale tissue histopathology image classification, segmentation, and visualization via deep convolutional activation features," *BMC bioinformatics*, vol. 18, no. 1, pp. 1–12, 2019.
- [10] S. Wang, H. Yu, Y. Gan, M. Zhang, and J. Zhu, "Gastric cancer classification using transfer learning from deep convolutional neural networks," *IEEE Access*, vol. 6, pp. 58 274–58 282, 2018.
- [11] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 234–241.
- [12] S. Mohan, C. Jayapandian, K. Patel, and S. Ganesan, "Deep learning methods for oral cancer classification," *IEEE Access*, vol. 8, pp. 190 340–190 348, 2020.
- [13] B. Wei, X. Han, J. Gu, X. Jiang, and J. Huang, "Breast cancer histopathological image classification using deep learning," *IEEE Transactions on Medical Imaging*, vol. 39, no. 12, pp. 4364–4372, 2020.
- [14] A. Cruz-Roa, A. Basavanahally, F. Gonzalez, H. Gilmore, M. Feldman, N. Shih, and A. Madabhushi, "Automatic detection of invasive ductal carcinoma in whole slide images with convolutional neural networks," in *Medical Imaging 2014: Digital Pathology*, vol. 9041. SPIE, 2014, pp. 904 103–1–904 103–7.
- [15] F. A. Spanhol, L. S. Oliveira, C. Petitjean, and L. Heutte, "Breast cancer histopathological image classification using convolutional neural networks," *IEEE Transactions on Biomedical Engineering*, vol. 63, no. 7, pp. 1455–1462, July 2016.
- [16] P. Srinivas, J. Radha, and T. R. Swapna, "Lung cancer detection using transfer learning and resnet-50," in *IEEE International Conference on Intelligent Techniques in Control, Optimization and Signal Processing (INCOS)*, 2019, pp. 1–5.