# Classification and Prediction of Lung Cancer with Histopathological Images Using VGG-19 Architecture

N. Saranya[✉], N. Kanthimathi[✉], S. Boomika[✉], S. Bavatharani[✉], and R. Karthick raja[✉]

ECE, Bannari Amman Institute of Technology, Sathyamangalam, India
{saranyan,kanthimathi,boomika.ec18,bavatharani.ec18,
karthickrajar.ec18}@bitsathy.ac.in

**Abstract.** In recent times, for the diagnosis of several diseases, many Computer-Aided Diagnosis (CAD) systems have been designed. Among many life-threatening diseases, lung cancer is one of the leading causes of cancer-related deaths in humans worldwide. It is a malignant lung tumor distinguished by the uncontrolled growth of cells in the tissues of the lungs. Diagnosis of cancer is a challenging task due to the structure of cancer cells. Predicting lung cancer at its initial stage plays a vital role in the diagnosis process and also increases the survival rate of patients. People with lung cancer have an average survival rate ranging from 14 to 49% if they are diagnosed in time. The current study focuses on lung cancer classification and prediction based on Histopathological images by using effective deep learning techniques to attain better accuracy. For the classification of lung cancer as Benign, Adenocarcinoma, or Squamous Cell Carcinoma, some pre trained deep neural networks such as VGG-19 were used. A database of 15000 histopathological images was used in which 5000 benign tissue images and 10000 malignant lung cancer-related images to train and test the classifier. The experimental results show that the VGG-19 architecture can achieve an accuracy of 95%.

**Keywords:** Lung cancer · VGG-19 · Histopathological images · Convolutional neural network

## 1 Introduction

Abnormal growth of cells in the body causes cancer. Lung cancer is a dreadful disease amongst the most widely recognized malignant tumors, also known as lung carcinoma. It affects the estimation of 2.3 million people every year all over the world. About 85% of patients affected by lung cancer are by smoking and tobacco consumption. 10–15% of cancer arises in a person who never smoke is because of their exposure to secondhand smoke, air pollution, or any other chemical exhaust. Two different types of Lung cancer are present, Small Cell Lung Cancer (SCLC) and Non-Small Cell Lung Cancer (NSCLC).

NSCLC grows out of control when there is a change in the healthy cells of the lung, which in turn form a mass called a tumor, or a nodule. These cells can grow anywhere in the lung and the tumor can be cancerous or benign. NSCLC is of three different types. Adenocarcinoma- It is more often found in former smokers or non-smokers. This type of cancer usually grows in the outer edges of the lungs, Squamous Cell carcinoma- This cancer type will develop more in smokers or former smokers. It will start growing in the middle part of the lungs near the bronchi, Large Cell carcinoma- It is very hard to treat. This tumor grows rapidly and effectively spreads to different organs. The tumor size will determine the stages of cancer in the nodes of the lung. The biopsy report is used to confirm the disease. By utilizing traditional techniques that are widely used by physicians, radiologists around the world, lung cancer can generally be detected at the mature stage and after it has spread to a great extent. A patient's chances of survival when lung cancer is discovered at that stage are very low. In the previously mentioned issue, the problem of misdiagnosis is also a cause for concern. Sometimes it is possible for doctors to diagnose benign conditions as malignant and vice versa. The life of the patients will also be put in very high-risk situations. Computer-based analysis techniques can be used as support tools for radiologists and physicians to overcome this concern. Transfer learning reuse already acquired knowledge gained from the pre-trained model in previous tasks and use it to improve generalization with new data. The idea is to use a model which is already been trained on a larger dataset for a long time and has proven to work well. By providing an input Histopathological image and possibly adding appropriate infected person metadata, to deliver a measurable result linked to lung cancer risk. According to the framework, there are two objectives to be taken into account. The first step was to minimize the inconsistency in the evaluation and observation of lung cancer by inferring with the diverse clinicians. Inevitably, Computer-based strategies have been shown to enhance a physician's ability to work across diverse medical backgrounds.

## 2  Literature Survey

Many authors and researchers were continuing to detect the early stage of lung cancer with high accuracy by using a suitable algorithm.

Authors Acucena R. S. Soares, et al. [1], this method uses a 3D U net and 3D V net, which takes the 888 images which contain 1159 nodules of CT images from LIDCIDRI. After preprocessing the samples obtained are 1664. This model achieved an accuracy of 74% for 3D U-Net and 99% accuracy for the 3D V-net model accordingly. In this work, they attain the best result in 3D V-net. Aishwarya Kalra [2]  the dataset was collected from LUNA and the prediction of nodules was done in Kaggle. Next the detection and segmentation were achieved by using the LUNA dataset to provide us with cancer in the lungs using the CNN model which obtains the accuracy of AOU is 97% recall 74% precision 87% and it has the specificity of 97%. Author BardhRushiti [3], the lung cancer classification was done using Artificial Intelligence, the model used are VGG16, VGG19, ResNet 34, and ResNet 50. At last ResNet 50 obtains an accuracy of 88.93%, precision is 95.83%, and F1 score is 88.46%. For this, the dataset was get from LIDC-IDRI and the training was done in Google Colab.

Authors Siddharth Bhatia, et al. [4] used UNet and ResNet models for detection. The dataset of CT images from LIDC-IDRI, to highlight the lung regions. This model achieves an accuracy of 84%. For classifying the image the architecture used are XGBoost and Random forest. Authors Muhammad ZiarurRehman, et al. [5] has used a model CAD system for detection. Here they used SVM for classification and this model achieves an accuracy of 88.97% while the CAD system shows higher sensitivity of 89.9%. Ashnil Kumar, et al. [6] used CT images to diagnose lung cancer by using computer-aided diagnosis applications. They compared their working model to baseline techniques for the purpose of multi-modality image fusion, multi-branch techniques, and multichannel and segmentation. This detection shows that the CNN model shows a maximum accuracy of 99%.
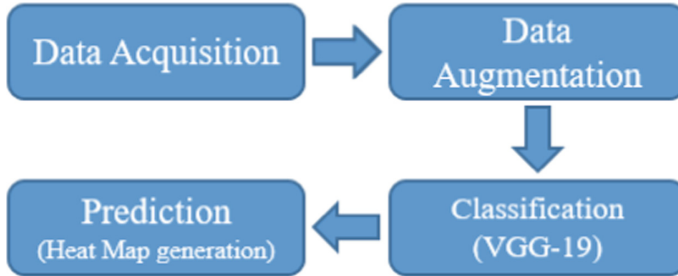
SajjaTulasi Krishna [7], used Convolution Neural Network (CNN) to analyze and classify both malignant and benign tissues from computed topology images. In this work neural network is designed on GoogleNet with a high dropout ratio to minimize the time of processing.25 epochs were used for training. GoogleNet shows 95.42% accuracy. Authors AbdarahmaneTraore, et al. [8], have used deep learning detectors such as Faster-RCNN, YOLO, SSD, RetinaNet, and EfficientNet in the critical task for the detection of the pulmonary nodule. This Faster RCNNmodel achieved a precision score of 35.73% and RetinaNet achieved a precision score of 34.15%. Wafaa Alakwaa, et al. [9], the network can be trained end to end from crude image patches. Its principle necessity is the accessibility of training database, yet otherwise no suspicions are made about the objects of interest or basic image methodology. Later on, it very well may be feasible to stretch out our present model to not just decide if the patient has disease, yet additionally decide the specific location of the cancerous knobs.

Authors Muhammad Attique Khan, et al. [10], started their work in the field of lung cancer identification with the assist of certain CT images. In this work, they approached the VGG-SegNet for segmenting the nodule in the lungs. Then they used the VGG19 algorithm with the SVM-RBF classifier for the purpose of the classification part. Finally, they acquired a maximum accuracy of about 94.83%. Mohammad Ali Abbas [11], used the Histopathological images in the diagnosis of cancer cells that are present in the lungs. In this paper, they used 3 CNN architectures and compared the results between them. VGG-19, AlexNet, ResNet-50 are the three methods used in this work. By comparing F1 score VGG-19 attains 0.997, Alex net attains 0.973 and ResNet-50 attains 0.999. Authors Imam Ali, et al. [12], diagnosed the malignant lung tissue with the assist of CT medical images. In this work, they approached the Transferrable texture CNN architecture for classification, this model mainly consists of three convolution layer networks and only one Energy layer for replacing the pooling layer. Totally this CNN model comprises nine layers for extracting the features and then goes for the classification part. Finally, this texture model acquired a higher accuracy of 96.69%.

The classification of medical images is one of the primary techniques of Computer-Aided Diagnosis (CAD) systems. Traditionally, image classification is based on shape, color, and texture features, and also their combinations. These features are probably problem-specific and have ended up being integral in clinical images diagnosis. Therefore, we have a framework that is unequipped for catching undeniable level issue area ideas and has restricted capacities in summing up models. As of late, deep learning techniques have ended up being a powerful method for developing a general model for registering the name for clinical images got from raw pixels. Not-withstanding,

deep learning models are costly and have restricted layers and channels because of the resolution of the clinical images and the small dataset size [13–18].

## 3   Methodology



### 3.1   Database Collection

Artificial Intelligence relies heavily on information. Computers find it difficult to learn without data. It is the most critical factor that makes the teaching of algorithms possible. The images will be in DICOM format, which means that the added data is useful as a means of visualizing and classifying. The dataset consists of 15000 histopathological images with three different classes namely Adenocarcinoma, Squamous Cell Carcinoma, and lung normal each with 5000 images. Images are in the size of 768 * 768 pixels and are in .jpeg file format. The images are converted to the .jpeg format with the same labels for effective classification. The converted images are stacked in separate directories based on the type of cancer
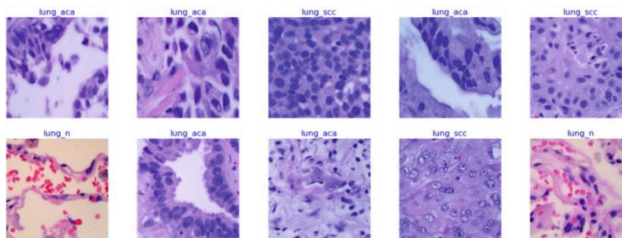


**Fig. 1.**  Histopathological images

### 3.2   Data Augmentation

As our model is based on deep learning, the number of samples will have a significant impact on its performance. Augmentation is a technique largely used to enhance the training and testing accuracy of convolutional neural networks, and further increases the performance of deep neural networks. A model can take into account what the image subject looks like in different situations by randomly adjusting the rotation, brightness, or

scaling of an input image. It helps to expose our classifier to a wider variety of situations to make our classifier more robust. It is a technique in computer vision includes adding noise, cropping, flipping, rotation, scaling, translation, brightness, contrast, color augmentation, and saturation, and for the Natural language processing EDA, Back translation, and text generation. Data augmentation algorithm increases the accuracy of machine learning models.

### 3.3   Transfer Learning

Transfer learning reuses already acquired knowledge gained from the pre-trained model in previous tasks and use it to improve generalization with new data. The idea is to use a model which is already been trained on a larger dataset for a long time and has proven to work well. Transfer learning improves learning new tasks by transferring knowledge and skills of tasks already learned solving other problems. Transfer learning is a method of the optimization training process with pre-trained models in a different task. Once the model is trained it can be reused as a base for a different tasks.

### 3.4   VGG-19 Architecture

VGG-19 is an optimized VGG model, which consists of 19 layers (16 convolution layers, three fully connected layers, five Maxpool layers, and one Softmax layer). There are different variants in VGG namely VGG11, VGG16, and so on. Unlike many other networks, the VGG-19 is a much simpler one with fewer hyper parameters. The classification is performed using 3 fully connected layers, consisting of two layers with 4096 neurons each, and one last layer with 1000 neurons.

ReLu activation functions are used in all layers except the last one, while for the last layer Softmax is used to distribute probabilities between classes. More than a million images from the ImageNet database were used to train VGG19. Using the network, images can be classified into 1000 categories, each with 19 layers. Using this network, images can be classified into 1000 object categories, each with 19 layers Fig. 1 and Fig. 2.
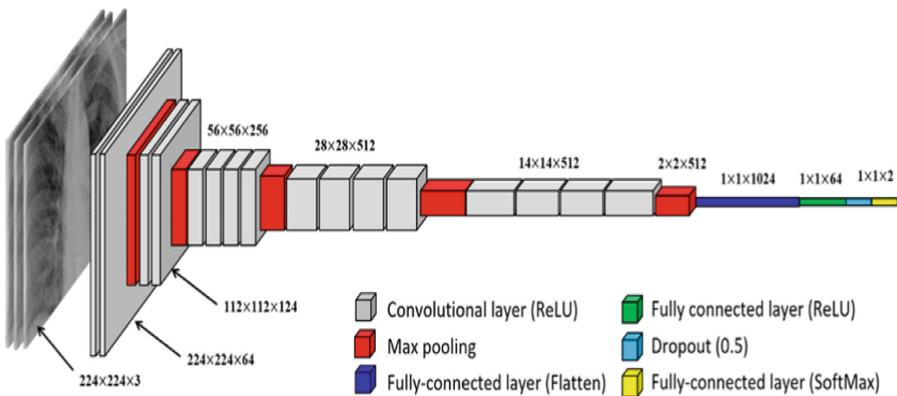


**Fig. 2.**  VGG-19 architecture

## 4   Results and Discussion

Below are the results of VGG-19 model for classifying the lung cancer.

| Validation Accuracy | 95% |
|---|---|
| Training Accuracy | 95.80% |
| Validation Loss | 30% |
| Training Loss | 29% |

For a deep learning model to achieve better while training the neural networks, loss function consider as a main key for adjusting the weights of the neural networks. While training the datasets during backpropagation, loss function penalizes the model assuming that there is any deviation between the label predicted by model and actual target label. Hence the use of loss function is hyper-critical to attain efficient model performance. Triplet loss is used as loss function in this work. This was first developed for recognition of face to find the similarities that occur in face. This method can be used when the images are blurred with the help of the distances between faces of similar and different identities. The triplet loss guides the model to minimize the distance between images of the similar category and increases the distance between images that belong to dissimilar categories. The use of triplet loss method resulted better accuracy in binary classification while the use of base model shows the less accuracy.

In Fig. 3, ten epochs were given to attain an expected accuracy based on training and validation datasets. The epochs indicate the number of passes that the learning algorithm will work through the entire training dataset comprised of one or more batches. VGG-19 model does not overfit. This model worked efficiently and gives the accuracy and loss by using all layers in VGG-19.



Starting training using base model VGG19 training all layers

| Epoch | Loss | Accuracy | V_loss | V_acc | LR | Next LR | Monitor | Duration |
|---|---|---|---|---|---|---|---|---|
| 1 /10 | 2.981 | 80.431 | 1.82777 | 85.667 | 0.00100 | 0.00100 | accuracy | 87.72 |
| 2 /10 | 1.193 | 90.667 | 1.11538 | 88.111 | 0.00100 | 0.00100 | val_loss | 87.09 |
| 3 /10 | 0.702 | 92.417 | 1.20456 | 54.222 | 0.00100 | 0.00050 | val_loss | 88.16 |
| 4 /10 | 0.740 | 93.306 | 1.11964 | 64.333 | 0.00050 | 0.00025 | val_loss | 87.58 |
| 5 /10 | 0.775 | 93.000 | 0.85363 | 90.333 | 0.00025 | 0.00025 | val_loss | 88.03 |
| 6 /10 | 0.642 | 93.931 | 0.71352 | 85.778 | 0.00025 | 0.00025 | val_loss | 85.67 |
| 7 /10 | 0.541 | 94.208 | 0.54909 | 93.000 | 0.00025 | 0.00025 | val_loss | 87.91 |
| 8 /10 | 0.453 | 95.083 | 0.46061 | 92.444 | 0.00025 | 0.00025 | val_loss | 85.53 |
| 9 /10 | 0.402 | 94.833 | 0.49259 | 89.778 | 0.00025 | 0.00013 | val_loss | 86.42 |
| 10 /10 | 0.398 | 95.403 | 0.50076 | 92.111 | 0.00013 | 0.00006 | val_loss | 86.14 |

**Fig. 3.**  Accuracy analysis for each epoch

## 4.1   Plots of Accuracy and Loss

Fig. 4 and 5 indicates the plots of accuracy and loss on the training and validation datasets over training epochs. The training and validation accuracy of the fine-tuned VGG-19 model is observed to have a similar convergence rate in both training and validation outcomes.
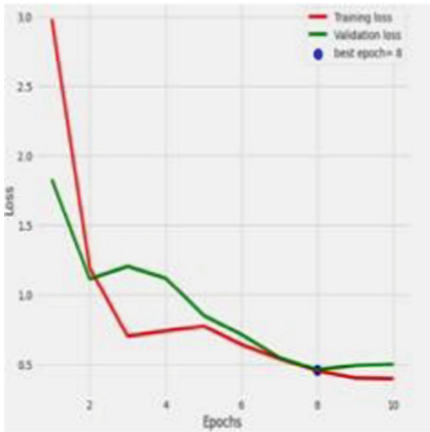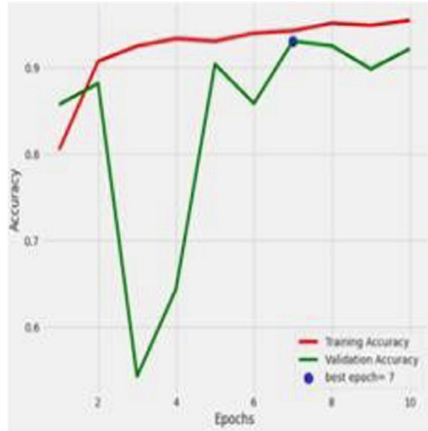


**Fig. 4.**  Model loss                    **Fig. 5.**  Model accuracy

## 4.2   Heat Map Generation Using Vgg-19

Heat map generation is the process of analyzing and reviewing heat map data to gather insights about the performance of the model during the training process. In Fig. 4.4, The Heat Map Generation method represents the accuracy of on testing process. Prediction is based on the following parameters such as,

TRUE POSITIVE (TP) - As a test result, it correctly indicates the presence of a condition.

TRUE NEGATIVE (TN) - As a test result, it correctly indicates the absence of a condition.

FALSE POSITIVE (FP) - As a test result, it wrongly indicates that a certain condition is present.

FALSE NEGATIVE (FN) - As a test result, it wrongly indicates that a certain condition is absent.

The classification report in Table 1 shows the evaluation parameters such as Precision, Recall, F1-Score for our proposed model. These values are calculated based on the below mathematical expressions,

$$\text{Accuracy(ACC)} ACC = (TP + TN)/(TP + TN + FP + FN)$$

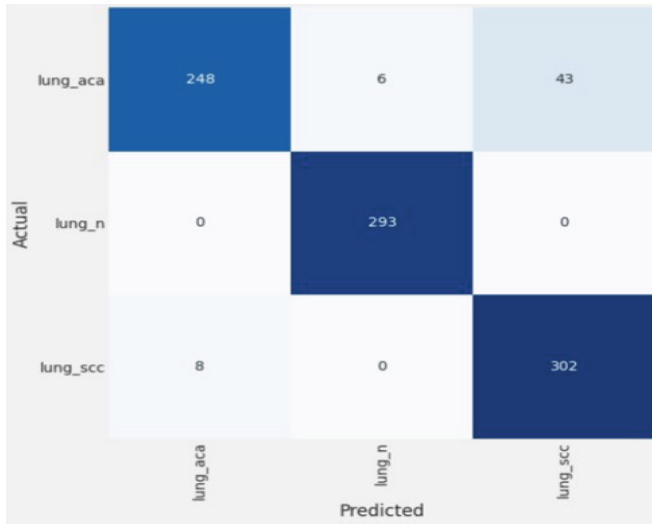$$\text{F1 Score } F1 = 2TP/(2TP + FP + FN)$$

**Fig. 6.** Heat map generation

Precision, recall, and F1 score are more helpful than accuracy performance metrics for predicting outcomes. As both indicates the accuracy of the model. These are some of the important metrics for evaluating models.

**Table 1.**   Evaluation parameters

|  | PRECISION | RECALL | F1 - SCORE | SUPPORT |
|---|---|---|---|---|
| Lung_aca | 0.97 | 0.84 | 0.90 | 297 |
| Lung_n | 0.98 | 1.00 | 0.99 | 293 |
| Lung_scc | 0.88 |  | 0.92 | 310 |
| Accuracy |  |  | 0.94 | 900 |
| Macro avg | 0.94 | 0.94 | 0.94 | 900 |
| Weighted avg | 0.94 | 0.94 | 0.94 | 900 |

## 5   Conclusion

In this project, the deep learning approach is used for classifying and predicting lung cancer using histopathological images. Early detection of lung cancer surely increases the possible chances of survival rate but it is more difficult to identify the lung cancer in the beginning stage itself. The proposed technique is based on Convolution neural networks model and transfer learning like VGG19. Here, the transfer learning method

is used for the purpose of increasing the validation accuracy. By using the VGG19 algorithm, an accuracy of 95.8% in the training phase and an accuracy of 95% in the validation phase were achieved. The Confusion matrix is the method used to evaluate parameters based on the performance of classification by VGG19. By using heat map generation, the prediction of the presence and absence of lung cancer was done. The future work is to collect the real-time lung affected dataset images to classify and detect lung cancer by using a suitable algorithm that shows526570 efficient accuracy. The proposed method VGG19 can also be applied to other cancers such as breast cancer, skin cancer, brain cancer, etc.

# References

Soares, A.R.S., Lima, T.J.B., Rabelo, R.D.A.L., Rodrigues, J.J.P.C., Araujo, F.H.D.: Automatic segmentation of lung nodules in CT images using deep learning. In: IEEE International Conference on E-health Networking, Application & Services (HEALTHCOM), pp. 1–6 (2020)

Kalra, A., Singh, B., Chauhan, H.: An Approach for lung cancer detection using deep learning. Int. Res. J. Eng. Technol. (IRJET) **7**(9) (2020)

Rushiti, B.: Automatic lung cancer detection using artificial intelligence university of business and technology Kosovo (2019)

Bhatia, S., Sinha, Y., Goel, L.: Lung Cancer Detection: A Deep Learning Approach. In: Bansal, J., Das, K., Nagar, A., Deep, K., Ojha, A. (eds) Soft Computing for Problem Solving. Advances in Intelligent Systems and Computing, vol 817, Springer, Singapore. (2019). https://doi.org/10.1007/978-981-13-1595-4_55

ur Rehman, M.Z., Javaid, M., Shah, S.I.A., Gilani, S.O., Jamil, M., Butt, S.I.: An appraisal of nodules detection techniques for lung cancer in CT images. Biomed. Sig. Process. Cont. **41**, 140–151 (2018)

Kumar, A. Fulham, M., Feng, D., Kim, J.: Co-learning feature fusion maps from PET-CT images of lung cancer. IEEE Trans. Med. Imag. **39**(1) 204–217 (2020)

Krishna, T.K., Devarapalli, Hemantha, R.M., Kalluri, H.K.: Lung cancer detection based on CT scan images by using deep transfer learning. Traitement du Sig. **36** 339–344 (2019)

Traore, A., Ly, A.O., Akhloufi, M.A.: Evaluating Deep Learning Algorithms in Pulmonary Nodule Detection, New Brunswick Health Research Foundation (NBHRF) (2020)

Alakwaa, W., Nassef, M., Badr, A.: Lung cancer detection and classification with 3D convolutional neural network (3D-CNN). Int. J. Adv. Comput. Sci. Appl. **08** (2017)

Khan, M.A.: VGG19 network assisted joint segmentation and classification of lung nodules in CT images. Diagnostics **11**(12), 2208 (2021)

Abbas, M.A., Bukhari, S.U.K., Syed, A., Shah, S.S.H.: The histopathological diagnosis of adenocarcinoma & squamous cells carcinoma of lungs by artificial intelligence: a comparative study of convolutional neural networks (2020)

Salaken, S.M., Khosravi, A., Khatami, A., Nahavandi, S., Hosen, M.A.: Lung cancer classification using deep learned features on low population dataset. In: IEEE 30th Canadian Conference on Electrical and Computer Engineering (CCECE) (2017)

Song, Q.Z., Zhao, L., Luo, X.K., Dou, X.C.: Using deep learning for classification of lung nodules on computed tomography images. J. Healthcare Eng. **04** (2017)

Lakshmanaprabu, S.K., Mohanty, S.N., Shankar, K., Arunkumar N., González, G.R.: Optimal deep learning model for classification of lung cancer on CT images. Future Gen. Comput. Syst. **92** 374–382 (2019)

Ali, I., Muzammil, M., Haq, I.U., Khaliq, A.A., Abdullah, S.: Efficient lung nodule classification using transferable texture convolutional neural network. IEEE Access **8** 175859–175870 (2020)

Traoré, A.  Ly, A.O.,  Akhloufi, M.A.: Evaluating deep learning algorithms in pulmonary nodule detection. In: 2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC) (2020)

Kumar, A., Fulham, M., Feng, D., Kim, J.: Co-learning feature fusion maps from PET-CT images of lung cancer.  IEEE Trans. Med. Imag. **39**(1), 204–217 (2020)

Singh, G., Gupta, G.K.: Performance analysis of various machine learning-based approaches for detection and classification of lung cancer in humans. Neural Comput. Appl. **31**, 6863–6877 (2019)