



## LUNG DISEASE IDENTIFICATION USING CHEST SCANS

Uday Kiran Irukula  
Dept. Of ECE VNRVJIET  
Hyderabad,India  
[udaykiranirukula@gmail.com](mailto:udaykiranirukula@gmail.com)

Shravanthi Kangula  
Dept. Of ECE VNRVJIET  
Hyderabad,India  
[shravanthikangula23@gmail.com](mailto:shravanthikangula23@gmail.com)

Shaik Khadar Sharif  
Dept. Of ECE VNRVJIET  
Hyderabad,India  
[shaikkhadarsharif@vnrvjiet.in](mailto:shaikkhadarsharif@vnrvjiet.in)

**Abstract--**Identification of lung diseases is a vital activity in the medical industry since it helps doctors diagnose and treat patients successfully. Deep learning algorithms have recently demonstrated promising outcomes in properly diagnosing various lung ailments from medical photos. In this paper, we evaluate the performance of two deep learning algorithms for identifying lung diseases: Convolutional Neural Networks (CNN), and VGG-16 (Visual Geometry Group with 16 layers depth). A dataset consisting of 13,622 chest X-ray pictures from patients with different lung conditions, such as pneumonia, tuberculosis, lung cancer and covid-19, was gathered by our team. The two algorithms were trained on the training set after dividing the dataset into training and testing sets. The trained models were subsequently put to the test on the testing set to evaluate how well they identified the pulmonary ailments. According to the findings of the project, the CNN algorithm had an accuracy of 96.87% in recognizing lung illnesses, followed by VGG-16 with an accuracy of 94.44. These findings show how deep learning algorithms, in particular CNN, are adept at correctly diagnosing lung disorders from medical images.

**Keywords--** Deep Learning, Convolutional Neural Networks, Visual Geometry Group

### I. INTRODUCTION

Lung illnesses affect the airways, lung disease tissues and other lung structures [6]. They are also referred to as respiratory disorders. Lung diseases include Pneumonia, Tuberculosis, Lung Cancer, and Coronavirus (COVID-19) illnesses. 334 million

people worldwide have asthma, and millions more die from pneumonia, lung cancer, and nearly 1.4 million people get tuberculosis every year, according to the Forum of International Respiratory Societies [9]. The COVID-19 epidemic affected every country in the world and put strain on the healthcare system by infecting millions of people. Lung conditions rank among the world's major killers and disablers. To improve long-term survival rates and boost the likelihood of recovery, early detection is essential. A subset of machine learning called "deep learning" is concerned with developing algorithms that use, are inspired by, and perform like the human brain [9]. Deep learning and other recent developments in machine learning make it possible to identify, quantify, and categorize diseases from medical images. These developments were made possible by deep learning's capacity to learn traits and features exclusively from data as opposed to manually creating features based on domain-specific expertise. Deep learning has been quickly reaching the state-of-the-art and improving productivity in a wide range of medical applications. Hence, these improvements aid clinicians in accurately identifying and classifying certain medical problems[9].

This research proposes an approach to categorize and identify lung problems as illnesses. In this study, the efficient model is used to identify the disease kind. Different models have been put out and are being compared in this research. This paper covers the detection of diseases like Covid-19, TB, Lung Cancer, and pneumonia. Also, we'll be able to tell if the CT scan image is of a typical normal person. The models that are being studied for this project are the Convolutional Neural Network (CNN), and the Visual Geometry Group with 16 layers depth (VGG16).



## II. LITERATURE SURVEY

Deep learning has improved extensively in the field of disease identification and classification. The significance can be observed from the previous years. Several studies have been done to understand the various models for the best effective solution. Due to the inadequate resources in the past, implementation of better and advances image processing techniques was not possible. After the improvement of GPU and CNN better disease detection was performed.

The work in Ref. [1] mainly focused on the disease detection and was experimented on the Covid19 CXR dataset. The dataset consisted of the normal CXR and Covid19 CXR. This method employed the ResNet50 and DenseNet which produces a validation accuracy of 86.67% and 98.33% respectively [1].

The work done in the Ref. [4] mainly focused on the need for detecting Covid-19, Pneumonia and Pneumothorax by implementing the VGG16 models. This study shows a promising result of accuracy lying between 93%-100% [4].

Another study was done for lung disease detection using deep learning techniques mentioned in the Ref. [6]. In this, three models namely Sequential, Functional and transfer models were trained, and a state-of-art CNN model was used. This paper proposed a comparative result of the three models showing a better accuracy of the Sequential model compared to other models.

## III. METHODOLOGY

Lung disorders affect millions of people each year. The identification of diseases like tuberculosis, COVID-19, pneumonia, lung cancer, etc., which are the most prevalent diseases affecting the world each year, requires X-ray images for disease detection. Deep learning-based models are beneficial techniques for improving performance and accuracy in a system for identifying lung disease in patients utilising X-ray pictures. Two different models are proposed in this project and are compared based on different metrics.

### A. DATASET ACQUISITION:

A dataset consisting of lung X-ray images is chosen for implementation of the CNN and VGG-16 models. The dataset consists of a total of 13,622 X-ray images. The dataset has a total of five classes which includes Covid-19, pneumonia, lung cancer, Tuberculosis and Normal. The dataset is divided into testing set, training set and validation set. The testing

dataset and training datasets were divided based on the ratio of 20:80. Here the training set contains 10,875 images and is further divided into five classes with 2552 images for Covid-19, 3990 images for pneumonia, 418 images for lung cancer, 535 images for tuberculosis and the remaining 3380 images to train the Normal class. The testing set has a total of 1421 images. This set has 121 images to test Covid-19, 704 images to test Pneumonia, 224 images to test Lung Cancer, 52 images for Tuberculosis and the rest 320 images are the testing set images for Normal class. The validation dataset consists of 1326 images.

**Table 1:** Classification of dataset

Disease Type	Testin g Images	Traini ng Images	Validati on images	Tota l
Covid-19	121	2552	136	2809
Pneumonia	704	3990	649	5343
Lung Cancer	224	418	50	692
Tuberculosis	52	535	165	752
Normal	320	3380	326	4026
Total	1421	10875	1326	13622

### A. IMAGE PRE-PROCESSING AND DATA AUGMENTATION

This phase is important as it enhances the quality of the image and helps in the reduction of noise. The enhanced or the modified images are used for the training phase. The image size is resized further to a size of 224,224. For data augmentation the ImageDataGenerator is used which is imported from the keras library. The images are normalized by rescaling the images by a factor of 255. The resizing of the images helps to train the images faster. The grayscale images are converted into the RGB format. The horizontal flip was enabled for the images. The shear range and zoom range was adjusted to a factor of 0.2 for better training process.

### B. TRAINING MODELS:

#### CNN Model:

In this study, CNN is used in the context of image classification. The "convolution" operation, which gives CNN its name, is carried out by the convolutional layer, which is regarded as one of its key layers [17]. The convolutional layer's inputs are subjected to its kernels.



A feature map is created by convolving all the convolutional layer outputs. Given that images are inherently nonlinear, the Rectified Linear Unit (ReLU) has been used in this study's activation function with a convolutional layer, which helps to augment the nonlinearity in the input image. The ReLU function is defined as,

$$ReLU(x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

Another crucial component of CNN is the pooling layer, often known as the subsampling layer. Pooling in CNN model can be of different types, these include the max, average or sum. In this study we are using the max pooling the implemented CNN model.

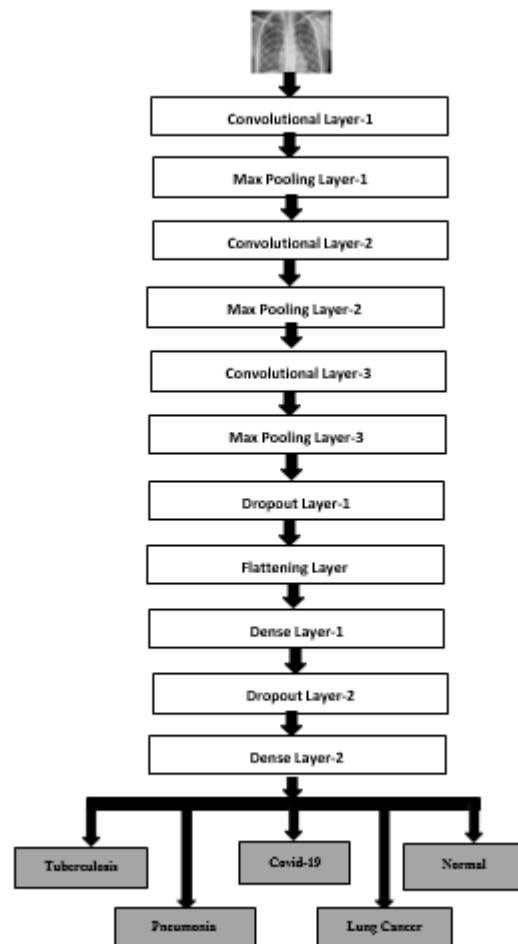
The implemented CNN model consists of three Convolutional layers (Conv2D), three max pooling layers (MaxPool2D), two dropout layers, one flattening layer and two dense layers. The input image for the CNN Model has a shape of (150,150,3). This means that the image has a dimension of 150-by-150 and is a RGB image. The model proposed is a sequential model. For a simple stack of layers with precisely one input tensor and one output tensor for each layer, a sequential approach is appropriate. For the first convolution layer 32 filters are used, for the second convolution layer 64 filters are used and for the third convolution layer 128 filters are used. In each convolution layer, the activation function 'ReLU' is used. The kernel size of 3 x 3 is used at convolution stages. After each convolution layer a max pool layer follows. In the pooling layer, a pool size of 2 x 2 is used. After the final pooling a dropout layer is used with 20% dropout rate for efficient computations. To implement the fully connected layers, the 3-dimensional convolution and pooling layers output must be converted into a single dimensional output. Hence to flatten the output a flattening layer is used followed by two dense layers with the 'softmax' activation function. The softmax function is defined as,

$$\sigma(o_i) = \frac{e^{o_i}}{\sum_{j=1}^n e^{o_j}}$$

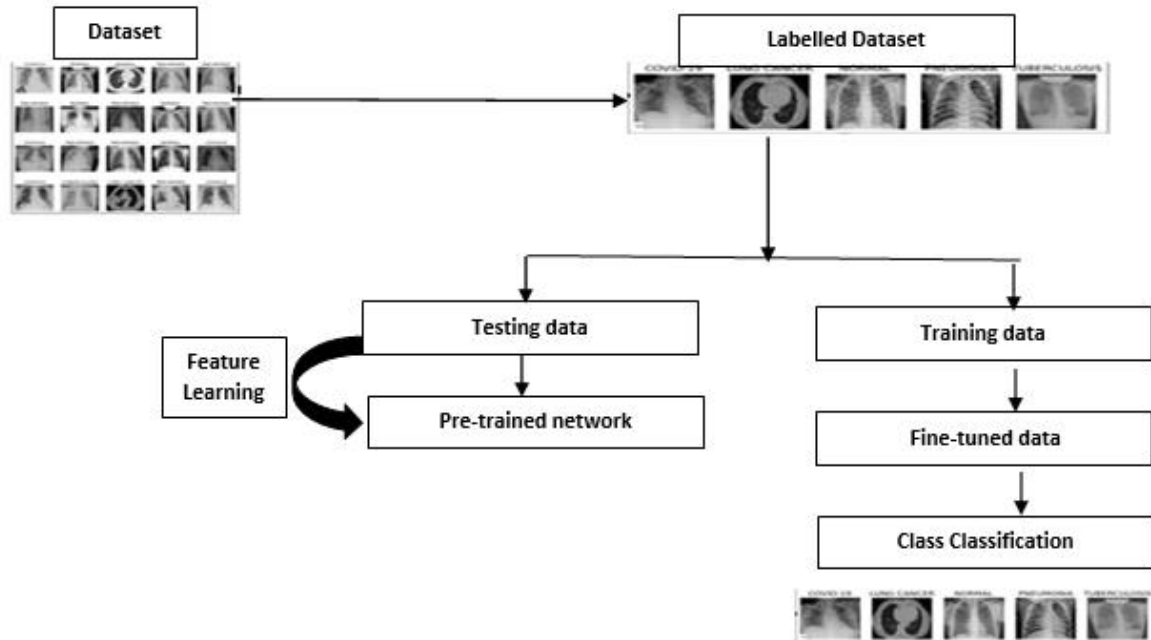
Where, i is the index and  $i=0, \dots, n-1$  and o is defined as the vector output. To compile the model, the categorical cross entropy loss function is used since it is used for multi-class classification and when there are two or more output labels. The optimizer used in this model is the 'Adam' optimizer. The Adam optimizer yields better results, require very less parameters for tuning and have faster computation time. The Adam optimizer works on the following formula,

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \frac{\partial L}{\partial w_t}, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) \frac{\partial L}{\partial w_t}$$

where, B1 and B2 terms refer to the decay rate of the average of gradients. The proposed CNN model has been trained for 20 epochs. The batch size of 64 is chosen for the proposed model. The learning rate for an effective model lies between 0.0 and 1.0. In this study we have used a learning rate of 0.001. A drop factor of 0.25 is used for implementation of the proposed model.



**Figure 1** Block Diagram of Proposed Model



**Figure 2** Architecture of Proposed Model

**B. CLASSIFICATION PHASE:**

This is the final phase of all the three models where the trained models will predict which class does a particular image belongs to. Based on our study, there are five classes our models can predict the output. These classes specify if the image belongs to a healthy person or if the image is of a person suffering from TB or Pneumonia or Covid-19 or Lung Cancer. Each of the class is given a label for easy identification.

**IV. RESULTS AND DISCUSSIONS**

We experimented upon the dataset containing the Lung X-ray images of patients with various lung diseases. During the training process, the various effects of parameters are analysed. This analysis was used to identify the optimum adjustment of these parameters to achieve better accuracy. The Python ‘Keras’ was used in the implementation of CNN and VGG models.

**V. PROPOSED CNN MODEL**

Parameter	Values
Optimizer	Adam
Learning rate	0.0001
Drop factor	0.25

Pooling method	Max pooling
Maximum Epochs	20
Batch size	32

**Table 2:** Proposed CNN model parameters

The CNN model proposed uses an Adam optimizer, with a batch size of 32. The model is trained on 13,622 images to predict the disease. In this, four diseases and five classes are predicted. These include a Normal class, TB, Pneumonia, Lung Cancer and Covid-19.

*Effects of hyper-parameters on the proposed CNNModel*

*(1) Effect of Optimizer*

An optimizer reduces loss thereby increasing the accuracy. In the proposed CNN model, an ‘Adam’ optimizer is used. This optimizer works on the Adam’s working formula. By making use of the Adam optimizer the computation time is less and provides efficient results

*(2) Effect of Learning rate*

This hyper-parameter lies in the range of 0.0 to 1.0. Learning rates affect the performance of the model and are related to drop factor. We can observe that when the learning rate is reduced to 0.0001 from 1.0,





there is a greater accuracy and requires more epochs compared to the case when the learning rate was adjusted to 1.0.

*(3) Effect of Drop factor*

The main advantage of using a dropout layer in training a CNN model is to avoid overfitting. In this a dropout rate of 0.25 is used. This means that 25% of the neurons will randomly be dropped in each epoch. By less dropping rate we get a better accurate prediction with a greater computation time.

*(4) Effect of epochs*

All the three models are trained for a total of 20 epochs. As epochs as increased, the network can become optimal from underfitting. In the proposed model, 20 epochs were considered which were required for optimal accuracy and performance in the proposed model.

*(5) Effect of batch size*

A model cannot process or train the data all at once and hence this is done in batches. In the CNN model a batch size of 32 is taken. When batch size is high fewer steps are required for the optimal solution. But we have taken a smaller batch size since the learning rate is also small. We can observe that accuracy declines with increase in batch size.

**VI. CONFUSION MATRIX AND PERFORMANCE METRICS**

In this study, we have a 5\*5 confusion matrix since there are 5 classes. The various performance metrics determined based on the confusion matrix are,

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$F1\ Score = 2 * \frac{(Recall * Precision)}{(Recall + Precision)}$$

where, TP = True Positive, TN = True Negative, FP = False Positive, FN = False Negative

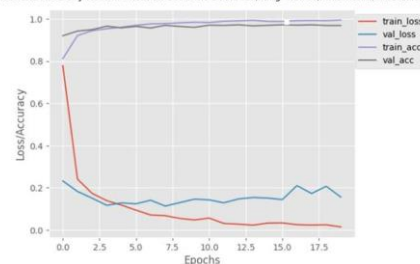
**Table 3: Confusion matrix of CNN Model**

	Covid-19	Lung Cancer	Normal	Pneumonia	TB
Covid-19	515	0	4	3	11
Lung Cancer	0	137	1	0	0
Normal	2	1	730	20	0
Pneumonia	3	0	24	1024	0
TB	11	1	1	0	137

**Table 4: Performance metrics using CNN Model**

Class	Precision(%)	Recall(%)	F1 score(%)
Covid-19	97.0	96.6	96.8
Lung Cancer	98.6	99.3	98.9
Normal	96.1	96.9	96.5
Pneumonia	97.8	97.4	97.6
TB	92.6	91.3	91.9

Training Loss and Accuracy for Classification between COVID-19,Lung Cancer,Pneumonia,Tuberculosis and Normal



**Figure 3: Graphical results of CNN model**

**V. DEPLOYING THE PROPOSED MODEL ON THE WEB**

The effective model chosen after the comparison is deployed on the web using the Flash application. The Flash framework is one such application in Python for deployment of Machine learning models on the web. In this application, a simple X-ray image can be uploaded by the user and the prediction is displayed on the



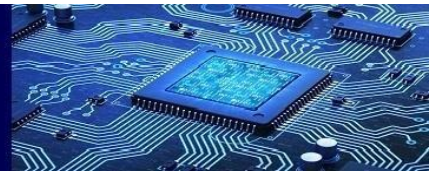
screen referring to the class or disease.

## VI. CONCLUSION

In this work two different models were proposed namely CNN, VGG16. Each of these models was trained on the same dataset containing 13,622 images. The proposed model was CNN since it outperformed other models and was finally deployed on the web due to the better accuracy of 96.87% compared to VGG16 model that produced an accuracy of 94.44%. Further adjustments to the optimizers, learning rate, and addition of more augmented data may result in significant advancements in the suggested CNN model's classification accuracy. Early stopping methods will probably give researchers more information on identifying lung conditions that can be passed down to prevent overfitting.

## V. REFERENCES:

- [1]. Anitha, J., Kalaiarasu, M., Kumar, N. S., & Sundar, G. R. (2022). Detection and classification of lung diseases using deep learning. *AIP Conference Proceedings*, 2519(1), 030001. doi:10.1063/5.0109980
- [2]. Goyal, S., & Singh, R. (2021). Detection and classification of lung diseases for pneumonia and Covid-19 using machine and deep learning techniques. *Journal of Ambient Intelligence and Humanized Computing*, 1-21. doi: 10.1007/s12652-021-03464-7
- [3]. Bharati, S., Podder, P., & Mondal, M. R. H. (2020). Hybrid deep learning for detecting lung diseases from X-ray images. *Informatics in Medicine Unlocked*, 20,100391. <https://doi.org/10.1016/j.imu.2020.100391>
- [4]. A. Chatchaiwatkul, P. Phonsuphee, Y. Mangalmurti, & N. Wattanapongsakorn. (2021). *Lung Disease Detection and Classification with Deep Learning Approach*. doi:10.1109/ITC-CSCC52171.2021.9501445
- [5]. Romalho, G. L. B., Rebouças Filho, P. P., Medeiros, F. N. S. de, & Cortez, P. C. (2014). Lung disease detection using feature extraction and extreme learning machine. *Revista Brasileira de Engenharia Biomédica*, 30. <https://doi.org/10.1590/rbeb.2014.019>
- [6]. M. Jasmine Pemeena Priyadarsini, Ketan Kotecha, G. K. Rajini, K. Hariharan, K. Utkarsh Raj, K. Bhargav Ram, V. Indragandhi, V. Subramaniaswamy, Sharnil Pandya, "Lung Disease Detection Using Various Deep Learning Algorithms", *Journal of Healthcare Engineering*, vol. 2023, Article ID 3563696, 13 pages, 2023. <https://doi.org/10.1155/2023/3563696>
- [7]. Batra, N., Goyal, S., & Chhabra, K. (2023). Lung Disease Detection Using Machine Learning Approach. In D. Gupta, A. Khanna, S. Bhattacharyya, A. E. Hassanein, S. Anand, & A. Jaiswal (Eds.), *International Conference on Innovative Computing and Communications* (pp. 251-260). Singapore: Springer Nature Singapore. [https://doi.org/10.1007/978-981-19-2821-5\\_21](https://doi.org/10.1007/978-981-19-2821-5_21)
- [8]. Tripathi, S., Shetty, S., Jain, S., & Sharma, V (2021). Lung disease detection using deep learning. *International Journal of Innovative Technology and Exploring Engineering*, 10(8), 8. DOI: 10.35940/ijitee.H9259.0610821
- [9]. Kieu STH, Bade A, Hijazi MHA, Kolivand H. A Survey of Deep Learning for Lung Disease Detection on Medical Images: State-of-the-Art, Taxonomy, Issues and Future Directions. *Journal of Imaging*. 2020; 6(12):131. <https://doi.org/10.3390/jimaging6120131>
- [10]. Gunasinghe, A. D., Aponso, A. C., & Thirimanna, H. (2019). *Early Prediction of Lung Diseases. 2019 IEEE 5<sup>th</sup> International Conference for Convergence in Technology (I2CT)*. doi:10.1109/i2ct45611.2019.9033668
- [11]. M. R. V, S. J, S. Koshy and N. G M, " A Survey on Lung Disease Diagnosis using Machine Learning Techniques," *2022 2<sup>nd</sup> International Conference on Advance Computing and Innovative Technologies in Engineering(ICACITE)*, Greater Noida, India, 2022, pp. 01-04, doi:10.1109/ICACITE53722.2022.9823787
- [12]. Diwan, M., Patel, B., & Shah, J. (2021). Classification of Lungs Diseases Using Machine Learning Technique. *International Research Journal of Engineering and Technology (IRJET)*, 9. e-ISSN: 2395-0056
- [13]. Rashmi, S. (2022). Lung Disease Identification Based on Chest X-ray and Lung Sounds Using Machine Learning and Deep Learning Techniques. *International Journal of Open Information Techniques*, 10(7), 94-100 e-ISSN: 2395-0056
- [14]. W. Ausawalaithong, A. Thirach, S. Marukatat and T. Wilaiprasitporn, "Automatic Lung Cancer Prediction from Chest X-ray Images Using the Deep Learning Approach," *2018 11<sup>th</sup> Biomedical Engineering International Conference (BMEiCON)*. Chiang Mai, Thailand, 2018, pp.1-5, doi:10.1109/BMEiCON.2018.8609997
- [15]. Zeilar, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional



- networks. In *Computer Vision-ECCV 2014: 13<sup>th</sup> European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13* (pp. 818-833). SpringerInternational Publishing.
- [16]. C. Szegedy et al., "Going deeper with convolutions," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Boston, MA, USA, 2015, pp. 1-9, doi:10.1109/CVPR.2015.7298594.
- [17]. Reshi, A. A., Rustam, F., Mehmood, A., Alhossan, A., Alrabiah, Z., Ahmad, A., ... & Choi, G. S. (2021). An efficient CNN model for COVID-19 disease detection based on X- ray image classification. *Complexity*, 2021, 1-12. <https://doi.org/10.1155/2021/6621607>
- [18]. L. V. R. Kumari, P. Shreya, M. Begum, T. P. Krishna and M. Prathibha, "Machine Learning based Diabetes Detection," 2021 6th International Conference on Communication and Electronics Systems (ICCES), Coimbatre, India, 2021, pp. 1-5, doi: 10.1109/ICCES51350.2021.9489058.
- [19]. L. V. R. Kumari, P. Bokkolla, S. F. Syeeda, S. P. Gudala and M. Dasandla, "Detection of pneumonia using Deep Learning," 2021 5th International Conference on Trends in Electronics and Informatics (ICOEI), Tirunelveli, India, 2021, pp. 1272-1279, doi: 10.1109/ICOEI51242.2021.9452830.
- [20]. Machiraju, Gayatri & Kakumani, Aruna & Sharif, Shaikh. (2021). Object Detection and Tracking for Community Surveillance using Transfer Learning. 1035-1042. 10.1109/ICICT50816.2021.9358698.

