# A Lightweight CNN Model for Detecting Respiratory Diseases from Lung Auscultation Sounds using EMD-CWT-based Hybrid Scalogram

Samiul Based Shuvo[1,], Shams Nafisa Ali[1,], Soham Irtiza Swapnil[1,],
Taufiq Hasan[1], *Member, IEEE* and Mohammed Imamul Hassan Bhuiyan[2], *Member, IEEE*

*Abstract*—Listening to lung sounds through auscultation is vital in examining the respiratory system for abnormalities. Automated analysis of lung auscultation sounds can be beneficial to the health systems in low-resource settings where there is a lack of skilled physicians. In this work, we propose a lightweight convolutional neural network (CNN) architecture to classify respiratory diseases using hybrid scalogram-based features of lung sounds. The hybrid scalogram features utilize the empirical mode decomposition (EMD) and continuous wavelet transform (CWT). The proposed scheme's performance is studied using a patient independent train-validation set from the publicly available ICBHI 2017 lung sound dataset. Employing the proposed framework, weighted accuracy scores of 99.20% for ternary chronic classification and 99.05% for six-class pathological classification are achieved, which outperform well-known and much larger VGG16 in terms of accuracy by 0.52% and 1.77% respectively. The proposed CNN model also outperforms other contemporary lightweight models while being computationally comparable.

*Index Terms*—Lung auscultation sound, respiratory disease detection, lightweight convolutional neural networks, empirical mode decomposition, continuous wavelet transform, scalogram.

## I. INTRODUCTION

LUNG diseases are the third largest cause of death in the world [1]. According to the World Health Organization (WHO), the five major respiratory diseases [2], namely chronic obstructive pulmonary disease (COPD), tuberculosis, acute lower respiratory tract infection (LRTI), asthma, and lung cancer, cause the death of more than 3 million people each year worldwide [3], [4]. These respiratory diseases severely affect the overall healthcare system and adversely affect the lives of the general population. Prevention, early diagnosis and treatment are considered key factors for limiting the negative impact of these deadly diseases.

Auscultation of the lung using a stethoscope is the traditional diagnostic method used by specialists and general practitioners for the initial investigation of the respiratory system. Although physicians use various other investigation strategies such as plethysmography, spirometry, and arterial

[1]Samiul Based Shuvo, Shams Nafisa Ali, Soham Irtiza Swapnil and Taufiq Hasan are with Department of Biomedical Engineering, Bangladesh University of Engineering and Technology (BUET), Dhaka-1205, Bangladesh, Email: {sbshuvo.bme.buet, snafisa.bme.buet, swapnil.buetbme}@gmail.com, taufiq@bme.buet.ac.bd.
[2]Mohammed Imamul Hassan Bhuiyan is with Department of Electrical and Electronic Engineering, Bangladesh University of Engineering and Technology, Dhaka-1205, Bangladesh, Email: imamul@eee.buet.ac.bd
These authors share first authorship on, and contributed equally to, this work.
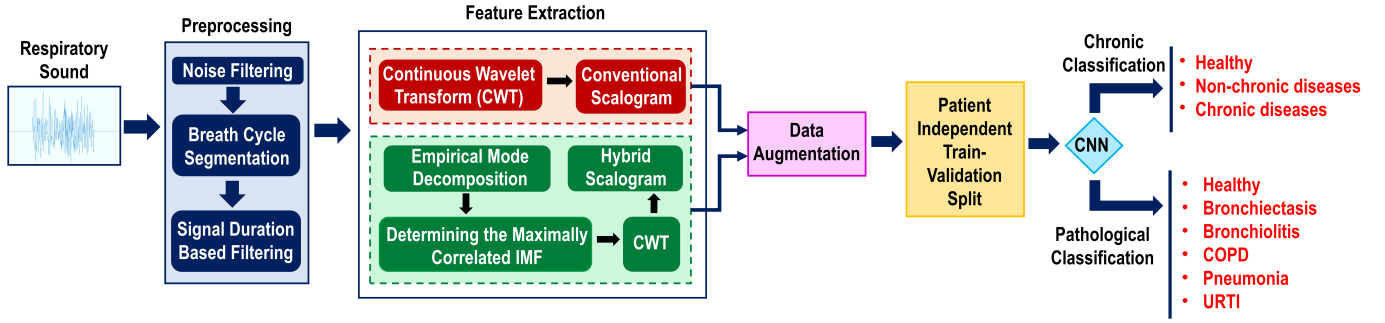Manuscript received August XX, 20XX; revised September XX, 20XX.

blood gas analysis, lung sound auscultation remains as a vital tool for physicians due to its simplicity and low-cost [5]. The primary classification of these non-periodic and non-stationary sounds consists of two groups: normal (vesicular) and abnormal (adventitious) [6]. The first group is observed when there are no respiratory diseases, while the latter group indicates complications in the lungs or airways [7]. Crackle, wheeze, rhonchus, squawk, stridor, and pleural rub are the commonly known abnormal lung sounds. These anomalies can be differentiated from the normal lung sounds on the basis of frequency, pitch, energy, intensity, timbre, and musicality [8], [9]. Therefore, lung sounds are of particular importance for recognizing specific respiratory diseases and assessing its chronic-nonchronic characteristics. However, the subtle differences between some of the adventitious lung sound classes can be a strenuous task even for a specialist and may introduce subjectivity in the diagnostic interpretation [10]. In this scenario, artificial intelligence (AI)-powered algorithms can be of benefit in automatically interpreting lung sounds, especially in underdeveloped regions of the world, with a scarcity of skilled physicians.

In the past decade, a number of research approaches have been considered and evaluated for automatic identification of respiratory anomalies from lung auscultations sounds. Numerous feature extraction techniques including statistical features [11], entropy-based features [12], wavelet coefficients [13], Mel Frequency Cepstral Coefficients (MFCC) [10], spectrograms [14], scalograms [15]etc. have been adopted in conjunction with a diverse set of machine learning (ML) algorithms [10]–[24].

With the advent of deep learning (DL), new developments have been made in recent times, demonstrating highly promising results in diversified applications, including biomedical engineering and clinical diagnostics [25]–[30]. With the ability of automatic feature learning, deep learning (DL) approaches are more generic and can mitigate the limitations of traditional ML-based methods. In the same vein, DL-based paradigms that are employed in the recent years for the identification of respiratory anomalies and pathologies from lung auscultation data have exhibited highly promising performance [5], [31]–[42]. However, for attaining proper functionality, the deep networks require to undergo an extensive training scheme with a large training dataset that subsequently calls for a considerable amount of time and the engagement of powerful computational resources. As a result, it becomes quite chal-

Fig. 1.   **A graphical overview of the proposed framework. After several generic preprocessing steps, the lung sound signals are converted into scalograms using both conventional and hybrid approaches. The resulting images are further augmented and fed into the proposed lightweight CNN model to carry out a two-way classification of respiratory diseases: (i) Chronic and (ii) Pathological.**

lenging to incorporate the deep learning frameworks in the currently available wearable devices and mobile platforms. In order to reduce the number of parameters of these networks, various methods have been investigated, including weight quantization [36], lightweight networks [43], and low precision computation [44].

While constructing AI-assisted automated medical diagnosis frameworks, patient specificity in the train and validation dataset should be considered a salient factor to produce reliable results for unseen patient data, especially for chronic diseases [45], [46]. Due to the available medical data's sparse nature, this factor is often neglected in the existing literature. The random adoption of 80%-20% or any other percentage of the train-validation split of the dataset corrupts most of the works with intra-patient dependency and ultimately the obtained results do not stand out to be consistent and generalized in case of a new patient [36]. Although this patient-independent division requires additional time and effort, the achieved results are more generalizable and represent the real-world scenarios.

In this work, a lightweight CNN architecture is proposed to perform respiratory disease classification (ternary chronic classification and six class pathology classification) utilizing the ICBHI 2017 scientific challenge respiratory sound database [47] while maintaining patient independent train-validation dataset splitting strategy. A hybrid approach for obtaining scalograms from respiratory sound signals is presented wherein continuous wavelet transform (CWT) is performed only on the maximally correlated intrinsic mode function (IMF) of the empirically decomposed (EMD) respiratory sound signals. The class discrimination capability of a hybrid scalogram is evaluated with respect to the CWT-based conventional scalogram. Subsequently, along with the proposed CNN model, complex CNN models such as VGG16 [48], AlexNet [49] and several contemporary lightweight architectures including MobileNet V2 [50], NASNet [51] and ShuffleNet V2 [52] are used for classifying the scalogram images to detect respiratory diseases in different categories. A comparative study among the proposed CNN model and the others is presented in terms of detection performance and being a lightweight network.

The rest of the paper is organized as follows. Previous studies related to lung sound classification using different ML-based approaches are discussed in Section II. Section III describes the dataset, feature extraction process, and the proposed lightweight CNN model. The experimental setup and results are discussed in Section IV. The performance of the proposed method is compared with other works in Section V. Finally, the concluding remarks are provided in Section VI.

## II. RELATED WORKS

Many research works employing machine learning, and deep learning have been reported on developing automated systems for respiratory sound classification. However, the majority of the works have focused on respiratory anomaly prediction, basically classifying the lung sounds as wheeze, crackles [10]–[24], [31]–[36] rather than directly predicting respiratory diseases from lung auscultation recordings. The few works geared towards pathology classification are very recent and mostly involve elaborate processing or dedicated CNN and RNN frameworks due to the inherent complexity of the signal [37]–[42]. However, at pathology-level, so far, the classification task has been investigated at three different resolutions; the binary classification (healthy, pathological) [37], [38], the ternary chronic classification (healthy, chronic disease, non-chronic disease) [38], [42] and multi-class distinct disease classification [39], [42]. Among the diseases, Upper and Lower Respiratory Tract Infection (URTI and LRTI), bronchiolitis and pneumonia have been included in the non-chronic disease class while COPD, asthma and bronchiectasis have been combined to form the chronic class [38].

In [37], a novel CNN based ternary classification approach has been implemented and performed considerably well with 82% accuracy and 88% ICBHI score. Later, the same authors proposed a Mel-Frequency Cepstral Coefficient (MFCC) and Long Short-term Memory (LSTM) based framework capable of conducting both binary and ternary classification of respiratory diseases [38] which demonstrated excellent performance with 99% and 98% accuracy, respectively. Another work involving complex RNN architecture and extensive preprocessing has reported accuracy of 95.67%0.77% in predicting six class pathology-driven diseases [39]. However, by employing a CRNN network with a CNN-Mixture-of-Experts (MoE) baseline to learn both spatial and time-sequential features

from the spectrograms, recent work has achieved a specificity of 83% and a sensitivity of 96% in ternary respiratory disease classification [40]. For binary classification, the same work has reported specificity and sensitivity of 83% and 99%, respectively. As an extension of [40], a separate study involving the robust Teacher-Student learning schemes with knowledge distillation has been conducted, which resulted in a substantially reduced specificity while maintaining the sensitivity [41].

Since the existing heavily imbalanced datasets of lung auscultations further exacerbate the task of respiratory disease classification, a contemporary study has dealt with this issue by experimenting with several data augmentation techniques, such as SMOTE, Adaptive Synthetic Sampling Method (ADASYN) and Variational autoencoder (VAE) [42]. Among the methods, the VAE-based Mel-spectrogram augmentation strategy, in conjunction with a CNN model, has achieved the best results with 98.5% sensitivity and 99.0% specificity in ternary chronic classification. The strategy has also exhibited an equally sophisticated performance with 98.8% sensitivity and 98.6% specificity in the case of six class respiratory disease classification [42].

Although the scope of DL-based frameworks with a spectrogram-based feature extraction strategy has been investigated in several works for direct classification of respiratory diseases from lung auscultations [40]–[42], to the best of the knowledge of the authors, scalogram based approaches have not been explored in this domain. Additionally, no dedicated lightweight, efficient CNN framework has been developed and investigated for the respiratory disease classification task. Furthermore, none of the studies consider the issue of intra-patient dependency in the train-validation split. Inspired by all of these factors, a scalogram based approach in conjunction with a lightweight CNN is proposed in this paper for the prediction of respiratory diseases from lung auscultations, maintaining patient independence. The proposed framework is schematically represented in Fig. 1.

## III. MATERIALS AND METHODS

### A. ICBHI 2017 Dataset

ICBHI (International Conference on Biomedical Health Informatics) 2017 database is a publicly available benchmark dataset of lung auscultations [47]. It is collected by two independent research teams of Portugal and Greece. The dataset contains 5.5 hours of audio recordings sampled at different frequencies (4 kHz, 10 kHz, and 44.1 kHz), ranging from 10s to 90s, in 920 audio samples of 126 subjects from different anatomical positions with heterogeneous equipment [53].

The samples are professionally annotated considering two schemes: 1. according to the corresponding patients pathological condition, i.e. healthy and seven distinct disease classes, namely Pneumonia, Bronchiectasis, COPD, URTI, LRTI, Bronchiolitis, Asthma and 2. according to the presence of respiratory anomalies, i.e. crackles and wheezes in each respiratory cycle. Further details about the dataset and data collection methods can be found in [53].

### B. Data Prepossessing

*1) Noise filtering:* Since 50 Hz to 2500 Hz is the acknowledged frequency range of the lung auscultation signals [7], the recorded audio signals are filtered with a 6th order Butterworth bandpass filter, thus retaining 50 Hz to 2500 Hz frequency components. Subsequently, all the sample signals are resampled to 22050 Hz for ensuring consistency and normalized to the range [-1,1] for attaining device homogeneity.

*2) Segmentation of the sound data:* Each of the audio recordings is segmented according to the annotated respiratory cycle timing with a 6s duration each. Samples with a minimum respiratory cycle duration of 3s are taken into account to obtain useful respiratory sound information [40]. Post performing this procedure, 2 of the disease classes, namely Asthma and LRTI, are found to have inadequate segmented samples for meaningful feature extraction and therefore, these two classes are not considered in our study. After these procedures, lung auscultation sounds from 87 out of 120 independent patients are usable. Table I represents data distribution at several levels of processing corresponding to the disease classes considered for this study.

TABLE I
DISTRIBUTION OF DATA AT DIFFERENT PROCESSING LEVELS
CORRESPONDING TO THE DISEASE CLASSES

| Disease Name | No. of unsegmented sound file | No. of Segmented and Filtered Sample | No. of Unique Patient | No. of Generalized Augmented Image |
|---|---|---|---|---|
| Pneumonia | 37 | 41 | 3 | 164 |
| Bronchiectasis | 16 | 55 | 6 | 220 |
| COPD | 793 | 1,963 | 51 | 1963 |
| Healthy | 35 | 42 | 13 | 168 |
| URTI | 23 | 21 | 8 | 84 |
| Bronchiolitis | 13 | 65 | 6 | 260 |
| Total | 917 | 2187 | 87 | 2859 |

### C. Feature extraction

*1) Empirical Mode Decomposition (EMD):* EMD is a powerful self-adaptive signal decomposition method especially in the time scale and energy distribution aspects and highly suitable for analysis and processing of non-linear and non-stationary signals such as lung sounds and heart sounds [54]. It decomposes a given signal *x(t)* into a finite set (N) of intrinsic mode functions, $IMF_1(t)$, $IMF_2(t)$, . . . . . , $IMF_N(t)$, depending on the local characteristic time scale of the signal, with a view to expressing the original signal as the sum of all its IMF plus a final trend either monotonic or constant called residue, *r(t)* a: $x(t) = \sum_{i=1}^{N} IMF_i(t) + r(t)$ [55]. An IMF is a simple oscillatory function with the equal number of extrema and zero crossings and its envelopes must be symmetrical with respect to zero. Thus, the EMD detrends a signal and elicits underlying spectral patterns [54].

*2) Continuous Wavelet Transform (CWT):* Wavelet transform is defined as a signal processing method that can decompose a signal into an orthonormal wavelet basis or into a set of independent frequency channels [15], [29]. Using a basis function, i.e. the mother wavelet *g(t)*, and its scaled and

dilated versions, the Continuous Wavelet Transform (CWT) can be used to decompose a finite-energy signal, $x(t)$ as [30]:

$$Z(a,b) = \int x(t) * g(t)(\frac{t-a}{b})  \qquad (1)$$

where b denotes time location and a is scale factor. Larger scale values reveal low-frequency information while the smaller scale values reveal high-frequency information [29]. The squared-modulus of the CWT coefficients Z is known as the scalogram [15].

### D. Scalogram Representations

*1) Conventional Scalogram:* Scalogram is defined as the time-frequency representation of a signal that depicts the obtained energy density using CWT [5], [56]. The segmented and filtered lung sound samples are decomposed into corresponding wavelet coefficients in MATLAB 2020a by using Morse analytic wavelet. Scalogram plots are generated with a resolution of 224 224 using these coefficients. Fig. 2 shows the scalograms of lung sounds in different disease categories.

*2) Hybrid Approach for Scalogram:* For each segmented and filtered sample under each pathological class, 9 IMFs are generated using the EMD function in MATLAB 2020a. Based on the cross-correlation between the source signal and the IMFs, the most physically significant IMF output with the highest correlation coefficient is determined [57], [58]. Subsequently, the squared-modulus of the CWT of the corresponding IMF is calculated to obtain the scalogram.

The diverse frequency bands varying from the maximum to the minimum range give the IMFs the capability to extract the temporal and spectral information [55] effectively. Hence, when this IMF based scheme is combined with CWT-oriented scalogram representation, the newly formed hybrid scalograms can demonstrate more discriminative and significant features. Thus, it has the potential to provide better classification performance by a CNN model. The box plots of the scalograms of lung sounds for various respiratory diseases are shown in Fig. 3. The distinction among the plots is more evident when using the hybrid approach than those of the conventional scalograms obtained using only CWT.

It should be mentioned that the proposed scalogram is distinctly different from that of [5], [58] in that the CWT modulus is computed from the maximally correlated IMFs and thus, providing a better representation of the underlying information. Note that the works of [5] and [58] are on detecting respiratory anomalies such as crackle and wheeze,
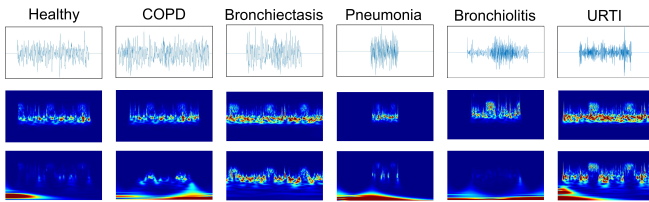
Fig. 3. **Box plot. (a) Scalogram using the conventional CWT approach; (b) Scalogram using the hybrid approach.**

and analysis and segmentation of heart sounds, respectively, whereas our objective is to detect respiratory diseases from the lung auscultations.

### E. Augmentation

The ICBHI 2017 dataset is highly imbalanced, with around 86% of the data belonging to COPD. Image augmentation using different color mapping schemes is employed to oversample the less represented classes and address the data imbalance issue [59]. Colormaps are three-column arrays containing RGB triplets where each row defines a distinct color. Scalogram representation using different color maps helps generalize the produced images.

From each of the audio samples of the less represented data classes, four scalograms are generated for each segmented sample using four different color mapping schemes: Parula, HSV, Jet, and Hot, which are available in MATLAB 2020a while for the most represented class, COPD, only one image is produced from each audio sample. Nevertheless, for ensuring generalization and homogeneity of the augmented data, all four-color mapping schemes are randomly utilized for COPD. A summary of segmented audio files and final augmented scalogram images with corresponding diseases classes are presented in Table I.

## IV. PROPOSED LIGHTWEIGHT CNN ARCHITECTURE

CNN has become a popular approach for classifying image data, and recently there have been several works using CNN on classifying images produced from sounds [5], [31], [33]. However, due to memory constraints, a regular deep CNN model is computationally expensive with its large number of learnable parameters and arithmetic operations. Thus, it is not suitable for embedded devices as they cannot afford the processing complexity and storage space for parameters and weight values of filters [36]. Cloud computing methodology requires a higher RAM for this computationally intensive training and hence are outsourced [60]. For this reason, Lightweight CNN models are gaining popularity among researchers for their faster performance and compact size without compromising the much-needed accuracy performance compared to the well-known deep learning networks [61].

The architecture of the proposed CNN model consists of an input layer corresponding to the 3-channel input of 224224 images. The architecture of the proposed model is illustrated in Fig. 4.

Fig. 2. **Scalograms of the lung auscultation sounds for 6 disease classes; lung sound recordings (1st row), conventional scalogram (2nd row) and scalogram using the proposed hybrid approach (3rd row).**
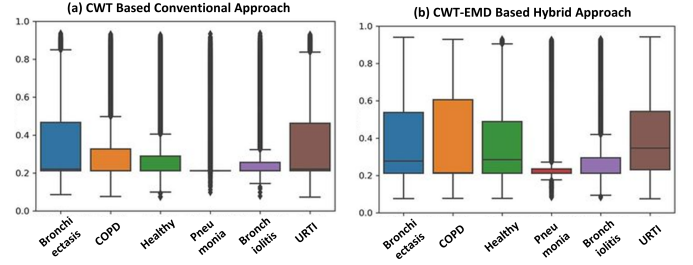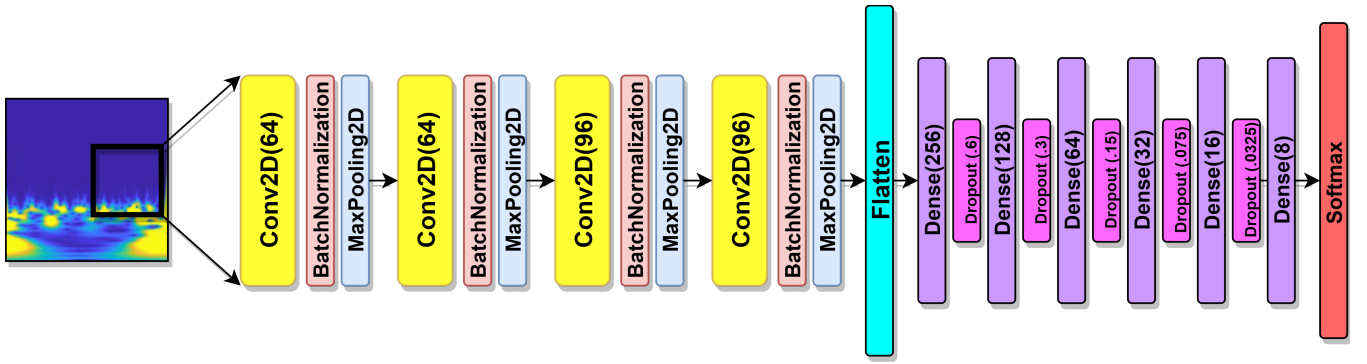
Fig. 4. **The detailed architecture of the proposed lightweight CNN model.**

The 1st convolutional layer uses 64 output filters with a 55-pixel size kernel followed by a 22-pixel max-pooling layer. Three additional convolutional layers are stacked over the first layer, each having a 33-pixel size kernel with 64, 96, and 96 filters sequentially and corresponding batch-normalization and max-pooling layers with 22 pooling window. Outputs from all these layers are flattened and connected with five pairs of FC and dropout layers, followed by a SoftMax output layer with probability nodes for each class. ReLU activation layer is applied within convolution calculation and fully connected layers and employed to introduce nonlinearity within the calculation and reduce the time for convergence. It does not get activated for any negative value. Max pooling is used after the ReLU activation; it reduces the spatial dimensionality of the extracted feature maps, extracts the most important features, and is unaffected from locational bias [35]. To overcome the problem of more diverse data variance, the Batch Normalization layer with every convolution layer normalizes the extracted feature. It gives the network a representative power with a small number of parameters and faster training capability by reducing the variance.

## V. EXPERIMENTAL RESULTS

### A. Evaluation Criteria

The augmented image sets are divided into 80% training and 20% validating parts for training and fine-tuning the model hyperparameters. Patient uniqueness, a critical aspect in the real-world applications, is maintained while dividing into training and validation parts as speaker dependency results in biased accuracy [46].

The classifier models' performance is evaluated based on the well-known evaluation matrices, namely, accuracy, recall (sensitivity), precision, and F1-score. Additionally, specificity and ICBHI-score [38], [53], a dedicated metric involving both sensitivity and specificity to assess the performance of the frameworks using the ICBHI dataset, is used to evaluate the performance of our method.

### B. Experimental Setup

The proposed CNN model is constructed using Keras and TensorFlow backend, and trained using NVidia K80 GPUs provided by Kaggle notebooks. The mini-batch training

scheme is employed while feeding the image data into a model for tackling the class imbalance issue. This technique performs by oversampling the scarce classes while randomly undersampling a majority class. This strategy ensures that the CNN model takes an equal number of samples from each class during each of the training epochs and thereby forms a balanced training set [26].

The adaptive learning rate optimizer (Adam) with the learning rate of 0.00001 is used for compiling the model. The batch size needs to be a multipler of 6 since an equal number of samples from each of the 3 and 6 data classes are taken in each training and validation batch [26]. In this study, batch size 6 has been taken for training and validation of both the classification schemes.

As stated earlier, both the ternary chronic classification (chronic, non-chronic, healthy) and six class (Bronchiectasis, Bronchiolitis, COPD, Healthy, Pneumonia, and URTI) pathological classification are carried out in this work. The classification performance of the proposed CNN model is compared with that of VGG16, a well-known CNN architecture for image classification [48] in both of the classification schemes. It should be noted that the experiments are performed using both the convention CWT-based scalogram and hybrid scalogram images. In addition, the performance of our proposed CNN model is compared with a number of well-known DL architectures such as VGG16 and AlexNet and several lightweight networks in terms of computational complexity and accuracy.

### C. Classification Performance of the Proposed Framework

*1) Chronic Classification:* From Table II, it can be seen that using the hybrid scalogram method in conjunction with the proposed CNN model classifier shows the best accuracy, 99.21%. However, the corresponding accuracy obtained by using VGG16 is quite close (98.89%). Despite being a heavy model, the comparatively lower accuracy of VGG16 can be attributed to the over-fitting issue due to the limited number of images in different classes. When comparing conventional CWT scalogram to the proposed hybrid scalogram, considerable improvement in accuracy is evident for the latter using VGG16 and our proposed CNN model (9.5%-11.4%). The corresponding confusion matrices for both the models' best results are illustrated in Fig. 5. The results depict that the

TABLE II
SUMMARY OF THE CLASSIFICATION PERFORMANCE
THE RED AND BLUE MARKED VALUES REPRESENT THE HIGHEST ACCURACY OBTAINED WITH OUR PROPOSED CNN MODEL AND VGG16

| Network | Chronic Classification | | | | | | | | Pathological Classification | | | | | | | |
| | Scalogram using CWT | | | | Scalogram using EMD and CWT | | | | Scalogram using CWT | | | | Scalogram using EMD and CWT | | | |
| | Prec. | Recall | Acc. | F1 | Prec. | Recall | Acc. | F1 | Prec. | Recall | Acc. | F1 | Prec. | Recall | Acc. | F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Proposed model | 95.00 | 91.00 | 90.58 | 92.00 | 99.25 | 99.20 | 99.20 | 99.22 | 87.00 | 86.00 | 86.31 | 86.00 | 99.12 | 99.05 | 99.05 | 98.96 |
| VGG16 | 92.00 | 89.00 | 88.58 | 90.00 | 99.03 | 98.69 | 98.69 | 99.11 | 85.00 | 85.00 | 85.80 | 85.00 | 97.80 | 97.32 | 97.32 | 97.01 |

TABLE III
COMPARISON OF THE PROPOSED FRAMEWORK WITH EXISTING WORKS USING THE ICBHI 2017 DATASET

| Processing | Type of Training Network | Number of Prediction Classes | Acc. | Spec. | Sen. | ICBHI Score |
|---|---|---|---|---|---|---|
| Gammatone Spectrogram [40], [41] | C-RNN | 3 (Healthy, Chronic, Non-chronic) | - | 0.57 | 0.94 | 0.76 |
| | CNN-MoE | | - | 0.86 | 0.96 | 0.91 |
| | Ensemble | | - | 0.71 | 0.95 | 0.83 |
| MFCC [38] | CNN | 3 (Healthy, Chronic, Non-chronic) | 0.82 | 0.76 | 0.89 | 0.83 |
| | LSTM | | 0.98 | 0.82 | 0.98 | 0.90 |
| MFCC [39] | RNN | 6 classes (excluding Asthma, LRTI) | 0.9567 | - | 0.9567 | - |
| Mel-spectrogram- VAE [42] | CNN | 3 (Healthy, Chronic, Non-chronic) | 0.99 | 0.990 | 0.985 | 0.988 |
| | | 6 classes (excluding Asthma, LRTI) | 0.99 | 0.986 | 0.988 | 0.987 |
| **Hybrid Scalogram (Proposed)** | Lightweight CNN | 3 (Healthy, Chronic, Non-chronic) | **0.99** | 1.00 | 0.99 | 0.995 |
| | | 6 classes (excluding Asthma, LRTI) | **0.99** | 1.00 | 0.99 | 0.995 |

proposed method is better in ternary chronic classification than the VGG16.

*2) Pathological Classification:* For six-class Pathological classification, the proposed method involving hybrid scalogram and proposed CNN model classifier yields the best accuracy, 99.05%, as seen in Table II. Similar to the case in the ternary chronic classification scheme, the accuracy of VGG16 is slightly lower. However, since the dataset gets more segregated being divided into six different disease classes, the accuracy drop is more here. The proposed hybrid scalogram outperforms the conventional CWT scalogram with a larger margin (13.4%-14.7%) for both VGG16 and our proposed model. In general, the proposed method gives a better performance, which is apparent from the best confusion matrices shown in Fig. 6.

### D. Comparison with Other Works

*1) Respiratory Disease Classification:* As discussed in section II, none of the existing works for respiratory disease classification explore the domain of patient-specific prediction. Some of the studies address the issues regarding class imbalance [39], [42]. Nevertheless, the extensive preprocessing, coupled with the ambiguous undersampling of the COPD disease class while oversampling all other disease classes, can complicate the reproducibility of [39]. Furthermore, in [42], FFT is applied to the entire respiratory sound signals, whereas our work focuses on segmented breath sounds. In our work, complete patient independence has been maintained in the train and validation set, which is not possible while using the entire lung auscultation signal due to the low number of samples. Therefore, our work aims to overcome all the drawbacks present in the existing methods. A comparison among the various methods, including the Proposed method, is provided in Table III.

It is observed that our proposed CNN model with the hybrid scalogram can perform on par with the existing state-of-the-art CNN and RNN models for both cases of classification while maintaining a patient independent train-validation scheme.

*2) Computational Performance as a Lightweight Network:* A detailed comparison is presented in Table IV among VGG16 [48], our proposed CNN model, AlexNet [49] and the existing state-of-the-art lightweight models such as MobileNetV2 [50], NASNet [51], ShuffleNetV2 [52] in terms of size, trainable parameters, the number of operations measured by multiply-add (MAdd) and accuracy on both chronic and pathological approach. In terms of accuracy, our proposed CNN model shows better results than VGG16 while requiring only 3% of the parameters. The proposed CNN model also outperforms the contemporary lightweight models, ShuffleNet V2, MobileNet V2, and NASNet relatively by 0.16%, 0.32%, and 0.80%, respectively, while obtaining better trade-off between the number of parameters, requiring significantly lower storage space and computational power. This makes our proposed lightweight model more suitable for real-time wearable
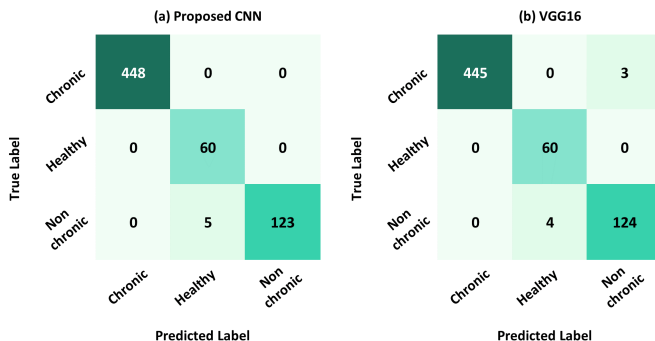


Fig. 5. **Confusion matrices for the best results obtained in ternary chronic classification. (a) Proposed CNN model with batch size 6; (b) VGG16 with batch size 6.**

devices with faster and less resource-intensive training.

We have calculated the time required for the end-to-end classification of an auscultation sound using our framework. For this experiment, we only performed the preprocessing and inference step using all of our test data and calculated the mean and standard deviation of the required CPU time. We found that the preprocessing time for EMD+CWT is $8s \pm 0.5s.$, and only CWT is $7.2s \pm 0.5s$. These processes are run on a Core i7 7500 processor with a 2.70-2.90GHz speed. Time required for the classification of a scalogram using the proposed network is $0.07s \pm 0.01s$, while the MobileNetV2 takes $0.085s \pm 0.01s$. Thus, the proposed CNN is faster in classifying a sound image as compared to MobileNetV2.

TABLE IV
COMPARISONS AMONG SEVERAL MODELS FROM LIGHTWEIGHT PERSPECTIVE

| Parameter | Network | | | | | |
|---|---|---|---|---|---|---|
| | VGG16 | AlexNet | **Proposed model** | Mobile-Net (v2) | Shuffle-Net (v2) | NASNet |
| Size(after training) | 1.5GB | 294MB | 44.85MB | 49MB | 46.9MB | 64MB |
| Trainable parameters | 138M | 25.704M | 3.7674M | 4.2M | 5.4M | 4.2M |
| MAdd | 154.7G | 725M | 371.93M | 575M | 564M | 567M |
| Accuracy (6 classes) | 97.60% | 98.237% | 99.05% | 98.89% | 98.27% | 98.73% |
| Accuracy (3 classes) | 97.60% | 99.519% | 99.21% | 98.72% | 99.06% | 98.42% |

## VI. CONCLUSION

In this work, we have proposed a lightweight CNN model to classify respiratory diseases using scalogram images of lung sounds. A hybrid approach employing both EMD and CWT is presented to generate the scalogram images. The publicly available ICBHI 2017 challenge dataset has been used for the Chronic and Pathological classification of respiratory diseases. The proposed method has provided a considerable accuracy of 99.21% for ternary chronic classification. In pathological classification among six disease classes, an accuracy of 99.05% is achieved. The obtained accuracies are higher than VGG16, which is a much larger network. In addition, for both cases



Fig. 6. **Confusion matrices for the best results obtained in six-class Pathological classification. (a) Proposed CNN model with batch size 6; (b) VGG16 with batch size 6.**

of classifications, the proposed framework provides better or a comparable performance with respect to the existing state-of-the-art methods in terms of Precision, Recall, F1-score, Sensitivity, Specificity and ICBHI score. It is worthwhile to mention that unlike most of these methods, the classification performance of the proposed technique has been assessed, keeping the training and testing data-independent in terms of patients. The proposed classifier's computational complexity has also been compared with a number of well-known CNN models and state-of-the-art lightweight networks. It has been shown to achieve high accuracy in classification while being a lightweight deep architecture. We believe that these attributes can enable the development of the automatic classification of respiratory diseases from lung auscultations in real-world clinical applications.

## REFERENCES

[1] "The Global Impact of Respiratory Disease Second Edition — CHEST Physician," 2017, [Online]. Available: https://www.mdedge.com/chestphysician/article/140055/society-news/global-impact-respiratory-disease-second-edition.

[2] A. A. Cruz, *Global surveillance, prevention and control of chronic respiratory diseases: a comprehensive approach*. World Health Organization, 2007.

[3] C. D. Mathers and D. Loncar, "Projections of global mortality and burden of disease from 2002 to 2030," *PLoS medicine*, vol. 3, no. 11, p. e442, 2006.

[4] "WHO — Global tuberculosis report 2019," 2019, [Online]. Available: https://www.who.int/tb/publications/global_report/en/.

[5] S. Jayalakshmy and G. F. Sudha, "Scalogram based prediction model for respiratory disorders using optimized convolutional neural networks," *Artificial Intelligence in Medicine*, vol. 103, p. 101809, 2020.

[6] A. Abbas and A. Fahim, "An automated computerized auscultation and diagnostic system for pulmonary diseases," *Journal of Medical Systems*, vol. 34, no. 6, pp. 1149–1155, 2010.

[7] S. Reichert, R. Gass, C. Brandt, and E. Andrès, "Analysis of respiratory sounds: state of the art," *Clinical Medicine. Circulatory, Respiratory and Pulmonary Medicine*, vol. 2, pp. CCRPM–S530, 2008.

[8] M. Sarkar, I. Madabhavi, N. Niranjan, and M. Dogra, "Auscultation of the respiratory system," *Annals of Thoracic Medicine*, vol. 10, no. 3, p. 158, 2015.

[9] A. Bohadana, G. Izbicki, and S. S. Kraman, "Fundamentals of lung auscultation," *New England Journal of Medicine*, vol. 370, no. 8, pp. 744–751, 2014.

[10] M. Bahoura and C. Pelletier, "New parameters for respiratory sound classification," in *Canadian Conference on Electrical and Computer Engineering*, vol. 3. IEEE, 2003, pp. 1457–1460.

[11] R. Palaniappan, K. Sundaraj, and N. U. Ahamed, "Machine learning in lung sound analysis: a systematic review," *Biocybernetics and Biomedical Engineering*, vol. 33, no. 3, pp. 129–135, 2013.

[12] J. Zhang, W. Ser, J. Yu, and T. Zhang, "A novel wheeze detection method for wearable monitoring systems," in *2009 International Symposium on Intelligent Ubiquitous Computing and Education*. IEEE, 2009, pp. 331–334.

[13] M. Bahoura, "Pattern recognition methods applied to respiratory sounds classification into normal and wheeze classes," *Computers in Biology and Medicine*, vol. 39, no. 9, pp. 824–843, 2009.

[14] J. Acharya, A. Basu, and W. Ser, "Feature extraction techniques for low-power ambulatory wheeze detection wearables," in *2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2017, pp. 4574–4577.

[15] N. Gautam and S. B. Pokle, "Wavelet scalogram analysis of phonopulmonographic signals," *International Journal of Medical Engineering and Informatics*, vol. 5, no. 3, pp. 245–252, 2013.

[16] G. Serbes, C. O. Sakar, Y. P. Kahya, and N. Aydin, "Pulmonary crackle detection using time–frequency and time–scale analysis," *Digital Signal Processing*, vol. 23, no. 3, pp. 1012–1021, 2013.

[17] S. İçer and Ş. Gengeç, "Classification and analysis of non-stationary characteristics of crackle and rhonchus lung adventitious sounds," *Digital Signal Processing*, vol. 28, pp. 18–27, 2014.
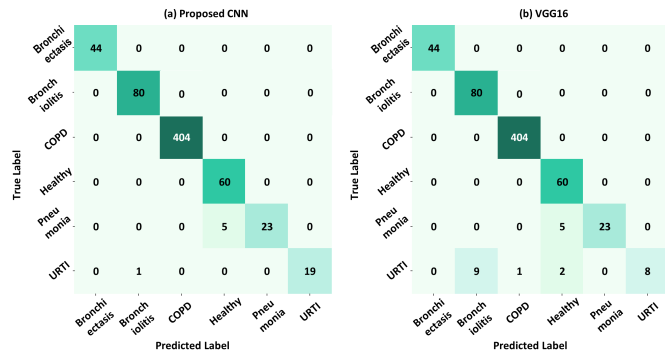
[18] F. Jin, F. Sattar, and D. Y. Goh, "New approaches for spectro-temporal feature extraction with applications to respiratory sound classification," *Neurocomputing*, vol. 123, pp. 362–371, 2014.

[19] P. Bokov, B. Mahut, P. Flaud, and C. Delclaux, "Wheezing recognition algorithm using recordings of respiratory sounds at the mouth in a pediatric population," *Computers in Biology and Medicine*, vol. 70, pp. 40–50, 2016.

[20] P. Mayorga, C. Druzgalski, R. Morelos, O. Gonzalez, and J. Vidales, "Acoustics based assessment of respiratory diseases using gmm classification," in *2010 Annual International Conference of the IEEE Engineering in Medicine and Biology*. IEEE, 2010, pp. 6312–6316.

[21] T. R. Fenton, H. Pasterkamp, A. Tal, and V. Chernick, "Automated spectral characterization of wheezing in asthmatic children," *IEEE Transactions on Biomedical Engineering*, vol. 32, no. 1, pp. 50–55, 1985.

[22] H. Pasterkamp, S. S. Kraman, and G. R. Wodicka, "Respiratory sounds: advances beyond the stethoscope," *American Journal of Respiratory and Critical Care Medicine*, vol. 156, no. 3, pp. 974–987, 1997.

[23] Z. Dokur, "Respiratory sound classification by using an incremental supervised neural network," *Pattern Analysis and Applications*, vol. 12, no. 4, p. 309, 2009.

[24] S. Rietveld, M. Oud, and E. H. Dooijes, "Classification of asthmatic breath sounds: preliminary results of the classifying capacity of human examiners versus artificial neural networks," *Computers and Biomedical Research*, vol. 32, no. 5, pp. 440–448, 1999.

[25] B. Bozkurt, I. Germanakis, and Y. Stylianou, "A study of time-frequency features for cnn-based automatic heart sound classification for pathology detection," *Computers in Biology and Medicine*, vol. 100, pp. 132–143, 2018.

[26] A. I. Humayun, S. Ghaffarzadegan, M. I. Ansari, Z. Feng, and T. Hasan, "Towards domain invariant heart sound abnormality detection using learnable filterbanks," *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 8, pp. 2189–2198, 2020.

[27] U. R. Acharya, S. L. Oh, Y. Hagiwara, J. H. Tan, and H. Adeli, "Deep convolutional neural network for the automated detection and diagnosis of seizure using eeg signals," *Computers in Biology and Medicine*, vol. 100, pp. 270–278, 2018.

[28] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, *et al.*, "Cnn architectures for large-scale audio classification," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 131–135.

[29] S. Debbal and F. Bereksi-Reguig, "Analysis of the second heart sound using continuous wavelet transform," *Journal of Medical Engineering & Technology*, vol. 28, no. 4, pp. 151–156, 2004.

[30] A. Meintjes, A. Lowe, and M. Legget, "Fundamental heart sound classification using the continuous wavelet transform and convolutional neural networks," in *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2018, pp. 409–412.

[31] K. Minami, H. Lu, H. Kim, S. Mabu, Y. Hirano, and S. Kido, "Automatic classification of large-scale respiratory sound dataset based on convolutional neural network," in *2019 19th International Conference on Control, Automation and Systems (ICCAS)*. IEEE, 2019, pp. 804–807.

[32] M. Aykanat, Ö. Kılıç, B. Kurt, and S. Saryal, "Classification of lung sounds using convolutional neural networks," *EURASIP Journal on Image and Video Processing*, vol. 2017, no. 1, p. 65, 2017.

[33] F. Demir, A. Sengur, and V. Bajaj, "Convolutional neural networks based efficient approach for classification of lung diseases," *Health Information Science and Systems*, vol. 8, no. 1, p. 4, 2020.

[34] R. Liu, S. Cai, K. Zhang, and N. Hu, "Detection of adventitious respiratory sounds based on convolutional neural network," in *2019 International Conference on Intelligent Informatics and Biomedical Sciences (ICIIBMS)*. IEEE, 2019, pp. 298–303.

[35] D. Bardou, K. Zhang, and S. M. Ahmad, "Lung sounds classification using convolutional neural networks," *Artificial Intelligence in Medicine*, vol. 88, pp. 58–69, 2018.

[36] J. Acharya and A. Basu, "Deep neural network for respiratory sound classification in wearable devices enabled by patient specific model tuning," *IEEE Transactions on Biomedical Circuits and Systems*, vol. 14, no. 3, pp. 535–544, 2020.

[37] D. Perna, "Convolutional neural networks learning from respiratory data," in *2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE, 2018, pp. 2109–2113.

[38] D. Perna and A. Tagarelli, "Deep auscultation: Predicting respiratory anomalies and diseases via recurrent neural networks," in *2019 IEEE 32nd International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE, 2019, pp. 50–55.

[39] V. Basu and S. Rana, "Respiratory diseases recognition through respiratory sound with the help of deep neural network," in *2020 4th International Conference on Computational Intelligence and Networks (CINE)*. IEEE, 2020, pp. 1–6.

[40] L. Pham, I. McLoughlin, H. Phan, M. Tran, T. Nguyen, and R. Palaniappan, "Robust deep learning framework for predicting respiratory anomalies and diseases," *arXiv preprint arXiv:2002.03894*, 2020.

[41] L. Pham, "Predicting respiratory anomalies and diseases using deep learning models," *arXiv preprint arXiv:2004.04072*, 2020.

[42] M. T. García-Ordás, J. A. Benítez-Andrades, I. García-Rodríguez, C. Benavides, and H. Alaiz-Moretón, "Detecting respiratory pathologies using convolutional neural networks and variational autoencoders for unbalancing data," *Sensors*, vol. 20, no. 4, p. 1214, 2020.

[43] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[44] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6869–6898, 2017.

[45] S. Kiranyaz, T. Ince, R. Hamila, and M. Gabbouj, "Convolutional neural networks for patient-specific ecg classification," in *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2015, pp. 2608–2611.

[46] N. U. Maheswari, A. Kabilan, and R. Venkatesh, "Speaker independent speech recognition system based on phoneme identification," in *2008 International Conference on Computing, Communication and Networking*. IEEE, 2008, pp. 1–6.

[47] "ICBHI 2017 Challenge," 2017, [Online]. Available: https://bhichallenge.med.auth.gr/.

[48] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[49] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.

[50] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

[51] X. Qin and Z. Wang, "Nasnet: A neuron attention stage-by-stage net for single image deraining," *arXiv preprint arXiv:1912.03151*, 2019.

[52] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "Shufflenet v2: Practical guidelines for efficient cnn architecture design," in *European conference on Computer Vision (ECCV)*, 2018, pp. 116–131.

[53] B. Rocha, D. Filos, L. Mendes, I. Vogiatzis, E. Perantoni, E. Kaimakamis, P. Natsiavas, A. Oliveira, C. Jácome, A. Marques, *et al.*, "A respiratory sound database for the development of automated classification," in *International Conference on Biomedical and Health Informatics*. Springer, 2017, pp. 33–37.

[54] N. Ibtehaz, M. S. Rahman, and M. S. Rahman, "Vfpred: A fusion of signal processing and machine learning techniques in detecting ventricular fibrillation from ecg signals," *Biomedical Signal Processing and Control*, vol. 49, pp. 349–359, 2019.

[55] M. Altuve, L. Suárez, and J. Ardila, "Fundamental heart sounds analysis using improved complete ensemble emd with adaptive noise," *Biocybernetics and Biomedical Engineering*, vol. 40, no. 1, pp. 426–439, 2020.

[56] Z. Ren, K. Qian, Y. Wang, Z. Zhang, V. Pandit, A. Baird, and B. Schuller, "Deep scalogram representations for acoustic scene classification," *IEEE/CAA Journal of Automatica Sinica*, vol. 5, no. 3, pp. 662–669, 2018.

[57] R. Fontugne, J. Ortiz, D. Culler, and H. Esaki, "Empirical mode decomposition for intrinsic-relationship extraction in large sensor deployments," in *Workshop on Internet of Things Applications, IoT-App*, vol. 12, 2012.

[58] D. Boutana, M. Benidir, and B. Barkat, "Segmentation and time-frequency analysis of pathological heart sound signals using the emd method," in *2014 22nd European Signal Processing Conference (EUSIPCO)*. IEEE, 2014, pp. 1437–1441.

[59] F. Y. Shih and H. Patel, "Deep learning classification on optical coherence tomography retina images," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 34, no. 08, p. 2052002, 2019.

[60] S. Y. Nikouei, Y. Chen, S. Song, R. Xu, B.-Y. Choi, and T. R. Faughnan, "Real-time human detection as an edge service enabled by a lightweight cnn," in *2018 IEEE International Conference on Edge Computing (EDGE)*. IEEE, 2018, pp. 125–129.

[61] B. Lim, B. Yang, and H. Kim, "Real-time lightweight cnn for detecting road object of various size," in *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*. IEEE, 2018, pp. 202–203.