



Automatic Non-Invasive Cough Detection based on Accelerometer and Audio Signals

Madhurananda Pahar¹ · Igor Miranda² · Andreas Diacon³ · Thomas Niesler¹

Received: 31 August 2021 / Revised: 9 January 2022 / Accepted: 23 February 2022 / Published online: 19 March 2022
© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

Abstract

We present an automatic non-invasive way of detecting cough events based on both accelerometer and audio signals. The acceleration signals are captured by a smartphone firmly attached to the patient's bed, using its integrated accelerometer. The audio signals are captured simultaneously by the same smartphone using an external microphone. We have compiled a manually-annotated dataset containing such simultaneously-captured acceleration and audio signals for approximately 6000 cough and 68000 non-cough events from 14 adult male patients. Logistic regression (LR), support vector machine (SVM) and multilayer perceptron (MLP) classifiers provide a baseline and are compared with three deep architectures, convolutional neural network (CNN), long short-term memory (LSTM) network, and residual-based architecture (Resnet50) using a leave-one-out cross-validation scheme. We find that it is possible to use either acceleration or audio signals to distinguish between coughing and other activities including sneezing, throat-clearing, and movement on the bed with high accuracy. However, in all cases, the deep neural networks outperform the shallow classifiers by a clear margin and the Resnet50 offers the best performance, achieving an area under the ROC curve (AUC) exceeding 0.98 and 0.99 for acceleration and audio signals respectively. While audio-based classification consistently offers better performance than acceleration-based classification, we observe that the difference is very small for the best systems. Since the acceleration signal requires less processing power, and since the need to record audio is sidestepped and thus privacy is inherently secured, and since the recording device is attached to the bed and not worn, an accelerometer-based highly accurate non-invasive cough detector may represent a more convenient and readily accepted method in long-term cough monitoring.

Keywords Accelerometer · Audio · Cough detection · LR · SVM · MLP · CNN · LSTM · Resnet50

1 Introduction

Coughing is a common symptom of respiratory disease and the forceful expulsion of air to clear up the airway [1]. It is distinctive in nature and is an important indicator used by physicians for clinical diagnosis and health monitoring in more than 100 respiratory diseases [2], including tuberculosis (TB) [3], asthma [4], pertussis [5] and COVID-19 [6]. Machine learning algorithms can be applied to the acoustic features extracted from the cough audio for automatic cough detection and classification [7–9]. However, using an audio-based monitoring system raises privacy issues [10], especially when the audio is captured by a smartphone [11, 12] and complex filtering processes might be required to preserve the privacy for continuous monitoring [13].

Acceleration measurements can be an alternative to the audio. Due to accelerometer's much lower sampling rates, less computing and processing power is required than the

✉ Madhurananda Pahar
mpahar@sun.ac.za

✉ Thomas Niesler
trn@sun.ac.za

Igor Miranda
igordantas@ufrb.edu.br

Andreas Diacon
ahd@sun.ac.za

¹ Department of Electrical and Electronic Engineering,
Stellenbosch University, Stellenbosch 7600, Western Cape,
South Africa

² Federal University of Recôncavo da Bahia,
Cruz das Almas 44.380-000, Bahia, Brazil

³ TASK Applied Science, Cape Town, Western Cape,
South Africa

audio [14]. Automatic cough detection based on accelerometer measurements is also possible when the device is placed on the patient's body and the acceleration signals are used for feature extraction [15]. Since an accelerometer is insensitive to environmental and background noise, it can be used in conjunction with other sensors such as microphones, ECG and thermistors [16]. Body-attached accelerometers have for example proved to be useful in detecting coughs when placed in contact with a patient's throat [17, 18] or at the laryngeal prominence (Adam's apple) [15]. A cough monitoring system using a contact microphone and an accelerometer attached to the participant's suprasternal (jugular) notch was developed in [19]. The participants moved around their homes while the cough audio and vibration was recorded. A similar ambulatory cough monitoring system, using an accelerometer attached to the skin of the participant's suprasternal notch using a bioclusive transparent dressing, was developed in [20]. Here, the recorded signal is transmitted to a receiver carried in a pocket or attached to a belt. Two accelerometers, one placed on the abdomen and the second on a belt wrapped at dorsal region, have been used to measure cough rate after cross-correlation of the two sensor signals [21]. Regression analysis, carried out on both audio and accelerometer signals gathered from 50 children, was able to achieve 97.8% specificity and 98.8% sensitivity when the accelerometer was placed in the centre of the abdomen between the navel and sternal notch [22]. Finally, multiple sensors, including ECG, thermistor, chest belt, accelerometer and audio microphones were used for cough detection in [23].

However, attaching an accelerometer to the patient's body is inconvenient and intrusive. We propose the monitoring of coughing based on the signals captured by the on-board accelerometer of an inexpensive consumer smartphone firmly attached to the patient's bed, as shown in Fig. 1. This eliminates the need to wear a measuring equipment and the system uses machine learning classifiers as promising results were reported in the studies mentioned previously, making it an automatic and also non-invasive cough detection system. The work presented here extends our previous study [24] by using three additional shallow classifiers along with the deep architectures in the cough detection process and by comparing the performance between the proposed accelerometer-based classifiers and the baseline systems that classify audio signals of the same cough events. Such audio-based cough detection systems have been reported to discriminate between coughing and other sounds with areas under the ROC curve (AUCs) as high as 0.96 [7] and specificities as high as 99% [25]. Although we have found that the audio-based cough detection still outperforms the accelerometer-based detection, we demonstrate that the difference in the performance is narrow, as our best 50-layer residual architecture (Resnet50) based

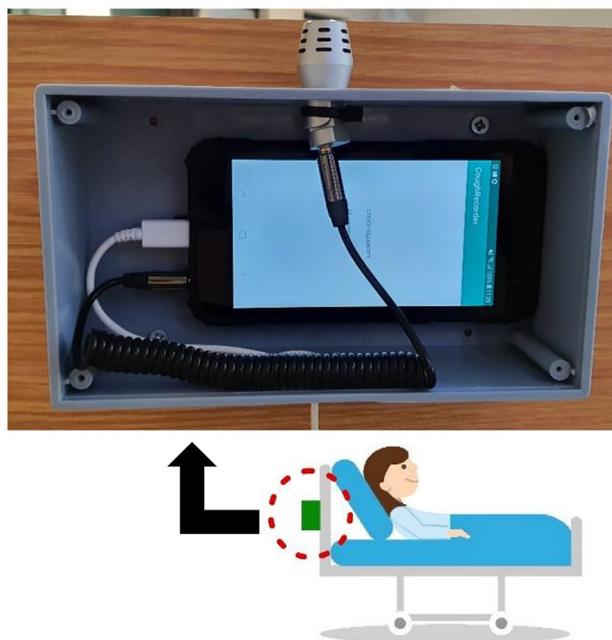


Figure 1 Recording Equipment: A plastic enclosure housing an inexpensive smartphone (Samsung Galaxy J4) running data gathering software is attached behind the headboard of each bed. The acceleration signal from the on-board accelerometer as well as the audio signal from the external microphone (BOYA BY-MM1), connected via a 3.5 mm audio jack, are monitored. Recording is triggered if activity is detected in either of these two signals.

cough detector achieves an AUC of 0.996 for audio-based and 0.989 for accelerometer-based detection respectively. Thus we are able to demonstrate that an automatic non-invasive accelerometer-based cough detection system is a viable option for long-term monitoring of patient's recovery.

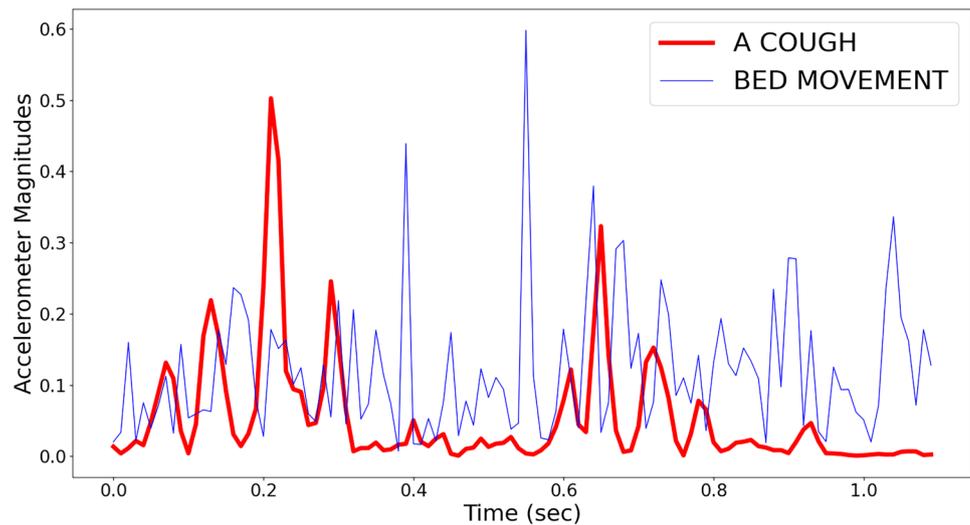
The structure of the reminder of this paper is as follows. Section 2 describes data collection while Sect. 3 details the features we extract from this data. The classifiers we use for experimentation are introduced in Sect. 4 and the classification process itself is elaborated in Sect. 5. The results are presented in Sect. 6 and discussed in Sect. 7. Finally, Sect. 8 concludes the paper.

2 Dataset Preparation

2.1 Data Collection

Data has been collected at a small 24h TB clinic near Cape Town, South Africa, which can accommodate approximately 10 staff and 30 patients. The clinic contains 8 wards and each ward has four beds, thus four patients at one time can be monitored inside a ward. The overall motivation of this

Figure 2 The accelerometer magnitudes for a cough event (red) and a non-cough event (blue). In this case, the non-cough event was the patient moving while on the bed.



study was to develop a practical method of automatic cough monitoring for the patients in this clinic, so that the recovery progress can be monitored.

Figure 1 shows the recording setup, where an enclosure housing an inexpensive consumer smartphone is firmly attached to the back of the headboard of each bed in a ward. An Android application, developed specifically for this study, monitors the accelerometer and the audio signals. The on-board smartphone accelerometer has a sampling frequency of 100Hz. Although this sensor provides tri-axial measurements, we record only the vector magnitude. A BOYA BY-MM1 external microphone was used to capture audio signals (visible in Fig. 1) at a sampling rate of 22.05 kHz. Using a simple energy detector, activity on either the acceleration or the audio channels triggers the simultaneous recording of both. This results in a dataset consisting of a sequence of non-overlapping time intervals during which both acceleration and audio have been recorded.

2.2 Data Annotation

A large volume of both the audio and accelerometer data has been captured by using this energy-threshold-based detection for both audio and acceleration signals. Ceiling-mounted cameras simultaneously recorded continuous video to assist with the data annotation process. The audio signals and the video recordings allowed the presence or absence of a cough in an event to be unambiguously confirmed during manual annotation. In the remainder of this paper, we will define an ‘event’ to be any interval of activity in either the accelerometer or the audio signals.

The non-cough events are generated mostly due to the patients getting in and out of the bed, moving while on the bed, sneezing or throat-clearing. Examples of the accelerometer

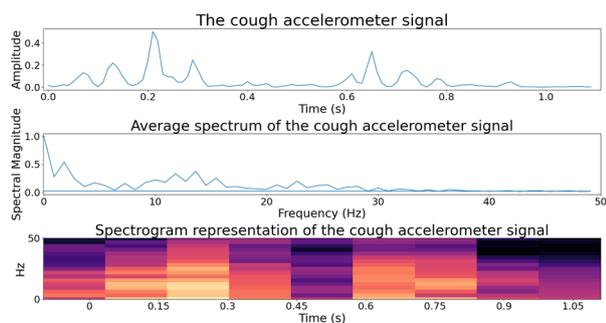
magnitude signals for a cough event and a non-cough event (in this case due to the patient moving while on the bed) are shown in Fig. 2. The spectrogram representations of these two signals are shown in Fig. 3. Manual annotation was performed using the ELAN multimedia software tool, which allowed easy consolidation of the accelerometer, audio and video signals for accurate manual labelling [26].

2.3 Final Dataset

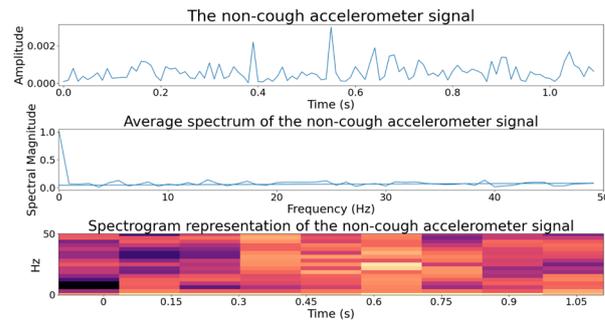
The final dataset, summarised in Table 1, contains approximately 6000 cough and 68000 non-cough events from 14 adult male patients. Cough events are on average 1.90 sec long, with a standard deviation of 0.26 sec. Non-cough events are on average 1.70 sec long, with a standard deviation of 0.24 sec. The total lengths of all cough and non-cough events are 11397.60 sec (3.16 hours) and 115928.12 sec (32.20 hours) respectively. No other information regarding patients is recorded due to the ethical constraints of the study. This dataset was used to train and evaluate six classifiers, introduced in Sect. 4 within a leave-one-out cross-validation framework, described in Sect. 5.

2.4 Dataset Balancing

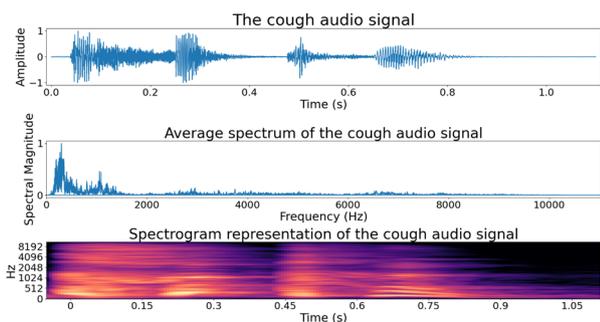
According to Table 1, cough events are outnumbered by non-cough events in our dataset. This imbalance can affect the machine learning classifiers detrimentally [27, 28]. We have applied the synthetic minority oversampling technique (SMOTE) to create additional synthetic samples of the minor class instead of for example oversampling randomly while training the classifiers [29, 30]. This addresses the class imbalance for both the accelerometer and audio events.



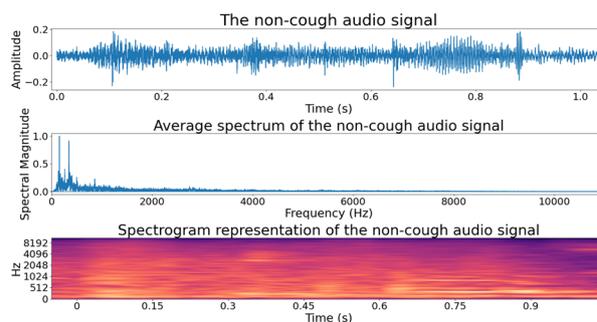
(a) Accelerometer signal of the cough event.



(b) Accelerometer signal of the non-cough (patient moving on the bed) event.



(c) Audio signal of the cough event.



(d) Audio signal of the non-cough (patient moving on the bed) event.

Figure 3 Spectrogram representation of the cough and non-cough events shown in Fig. 2: The cough event is shown in (a) and (c) and the non-cough event (patient moving on the bed) in (b) and (d). The accelerometer and audio signals are shown in (a) & (b) and (c) & (d)

SMOTE has previously been successfully applied to cough detection and classification based on audio recordings [9, 24, 31].

3 Feature Extraction

The feature extraction process is illustrated for both accelerometer and audio signal in Fig. 4.

3.1 Accelerometer Features

Power spectrum, root mean square (RMS) value, kurtosis, moving averages and crest factor are extracted from the accelerometer magnitude samples. No de-noising has been applied prior to the feature extraction process. Power spectra [32] has been used to represent sensor data for input to classifiers, including neural networks, in several studies [33–36]. RMS [37] values from the sensor data have also been found to be useful features [38, 39]. The kurtosis has

respectively. The audio signal has a higher sampling rate and thus contains more frequency and time-domain information than accelerometer measurements.

also been useful for machine learning applications as it indicates the prevalence of higher amplitudes [40]. Moving averages indicate the smoothed evolution of a signal over a time period and have been found to be useful features for sensor analysis [41]. Finally, the crest factor measures the ratio of the peak and the RMS signal amplitudes and have also been found to help machine learning prediction [42] including deep learning [43].

3.2 Audio Features

Features such as mel-frequency cepstral coefficients (MFCCs), zero crossing rate (ZCR) and kurtosis are extracted from the audio signal. MFCCs are successfully used as features in audio analysis and especially in automatic speech recognition [44, 45]. They have been successfully used to differentiate dry coughs from wet coughs [46] and also to identify coughs associated with tuberculosis [47] and COVID-19 [9, 48]. We have used the traditional MFCC extraction method considering higher resolution MFCCs along with the velocity (first-order

Table 1 Ground Truth Dataset: ‘PATIENTS’: list of the patients; ‘COUGHS’: number of confirmed cough events; ‘NON COUGHS’: number of confirmed events that are not coughs; ‘COUGH TIME’: total amount of time (in sec) for cough events; ‘NON-COUGH TIME’: total amount of time (in sec) for non-cough events.

PATIENTS	COUGHS	NON COUGHS	COUGH TIME	NON-COUGH TIME
Patient 1	88	973	169.16	1660.67
Patient 2	63	1111	117.67	1891.92
Patient 3	469	11025	893.91	18797.32
Patient 4	109	9151	204.06	15596.71
Patient 5	97	7826	188.26	13344.98
Patient 6	192	12437	360.72	21197.35
Patient 7	436	14053	825.23	23953.15
Patient 8	368	2977	702.05	5077.89
Patient 9	2816	3856	5345.27	6569.32
Patient 10	649	2579	1236.84	4400.42
Patient 11	205	527	391.42	901.38
Patient 12	213	323	402.61	547.62
Patient 13	213	712	401.61	1211.75
Patient 14	82	455	158.77	777.64
TOTAL	6000	68005	11397.6	115928.12

difference, Δ) and acceleration (second-order difference, $\Delta\Delta$) as adding these has shown classifier improvement in the past [49]. The ZCR [50] is the number of times a signal changes

its sign within a frame, and indicates the variability present in the signal. Finally, the kurtosis [51] indicates the prevalence of higher amplitudes in the samples of an audio signal. These features have been extracted by using the hyperparameters described in Table 2 for all cough and non-cough audio events.

3.3 Extraction Process

The features are extracted in a way that preserves the information regarding the beginning and the end of an event to allow time-domain patterns in the recordings to be discovered while maintaining the fixed input dimensionality, which is expected by the deep neural architectures such as a convolutional neural network (CNN).

For accelerometer signal, the frame length (Ψ) and number of segments (C) have been used as the feature extraction hyperparameters, shown in Table 2. Power spectra have the dimension of $(C, \frac{\Psi}{2} + 1)$ and each of RMS, kurtosis, moving averages and crest factor has the dimension of $(C, 1)$. Thus, the input feature matrix for the accelerometer signal, fed to the classifiers mentioned in Sect. 4, has the dimension of $(C, \frac{\Psi}{2} + 5)$, as illustrated in Fig. 4.

For audio signal, frame length (\mathcal{F}) and number of segments (\mathcal{S}) have been used as the feature extraction hyperparameters, shown in Table 2. Each of MFCCs, MFCC velocity (Δ), MFCC acceleration ($\Delta\Delta$) has the dimension of $(\mathcal{S}, \mathcal{M})$, where \mathcal{M} is the number of MFCCs. Each of ZCR and kurtosis has the dimension of $(\mathcal{S}, 1)$. Thus, the input feature matrix for audio signals,

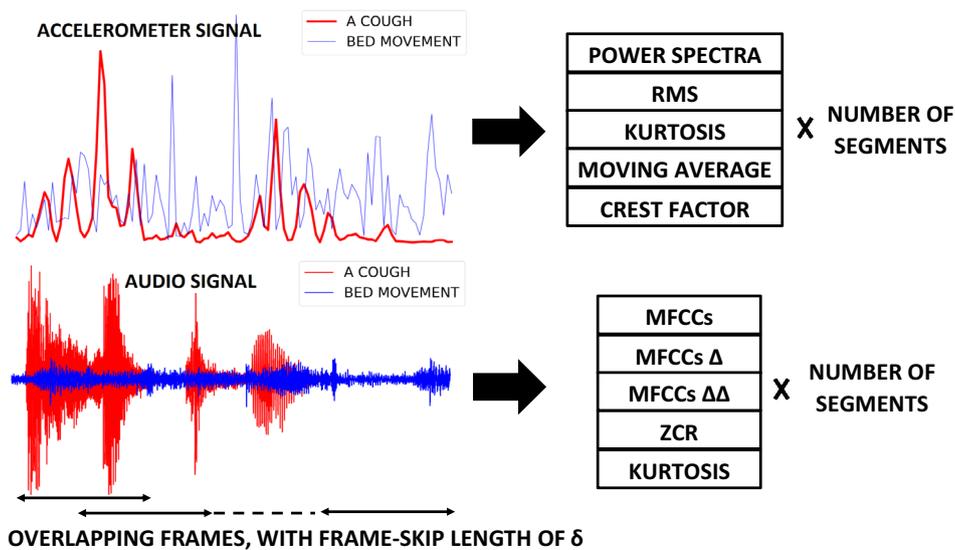


Figure 4 Feature extraction for both the accelerometer (top) and the audio (bottom) signals: Both acceleration and audio signals of the events, shown in Figs. 2 and 3, are split into a fixed number of overlapping frames. The length and number of these frames are Ψ and C for accelerometer signal & \mathcal{F} and \mathcal{S} for audio signal. For accelerometer measurements, the power spectrum, RMS, kurtosis,

moving average and crest factor of each frame are extracted. For audio signals, the MFCCs, MFCC velocity (Δ), MFCC acceleration ($\Delta\Delta$), ZCR and kurtosis are extracted. For the acceleration signal, this results in a feature matrix with dimensions $(C, \frac{\Psi}{2} + 5)$ while for the audio signal it generates a feature matrix with dimensions $(\mathcal{S}, 3\mathcal{M} + 2)$ where \mathcal{M} is the number of extracted MFCCs.

Table 2 Feature extraction hyperparameters for both accelerometer and audio signals. For accelerometer, 16, 32, 64 samples i.e. 160, 320 and 640 msec long frames overlap in such a way that the number of these frames i.e. segments (5 and 10) are the same for all events in our dataset. Similarly for audio signals, MFCCs are varied

between 13 and 65 & frames are varied between 256 samples (11.61 msec) and 4096 samples (185.76 msec) in such a way that the number of these extracted frames are varied between 50 to 150, fixed for all events in out dataset.

FEATURE EXTRACTION HYPERPARAMETERS

Accelerometer Hyperparameters		Values
Frame (Ψ)	Frame-length in samples, used to extract features	2^k where $k = 4, 5, 6$
Segments (C)	Number of frames extracted from the event	5, 10
Audio Hyperparameters		Values
MFCC (\mathcal{M})	Number of lower-order MFCCs to keep	$13 \times k$, where $k = 1, \dots, 5$
Frame (\mathcal{F})	Frame-length in samples, used to extract features	2^k where $k = 8, \dots, 12$
Segments (\mathcal{S})	Number of frames extracted from the event	$10 \times k$, where $k = 5, 7, 10, 12, 15$

fed to the classifiers mentioned in Sect. 4, has the dimension of $(\mathcal{S}, 3\mathcal{M} + 2)$, as illustrated in Fig. 4.

From every event, we extract a fixed number of features (C and \mathcal{S}) by distributing the fixed-length analysis frames (Ψ and \mathcal{F}) uniformly over the time-interval of the cough and by varying the length of frame skips, noted as δ in Fig. 4. To calculate frame skips, we divide the number of samples in an event by the number of segments and take the next positive integer. For a 1.2 sec long audio event, the length of frame skip in samples is $\left\lceil \frac{1.2 \times 22050}{100} \right\rceil = \left\lceil \frac{26460}{100} \right\rceil = 265$ samples, as the audio sampling rate is 22.05 kHz.

The frame length (Ψ) used to extract features from acceleration signal is shorter than the frame length (\mathcal{F}) used to extract features from audio in this study (Table 2) and also traditionally [52]. This is because, as already noted in Fig. 1, the accelerometer in the smartphone has a lower sampling rate of 100 Hz than the microphone 22.05 kHz and longer frames lead to deteriorated performance as the signal properties can no longer be assumed to be stationary [53]. The lower sampling rates also reduce the amount of computation needed to extract features.

In contrast with the more conventionally applied fixed, non-overlapping frame rates, this way of extracting features ensures that the entire event is captured within a fixed number of frames, allowing especially the CNN to discover more useful temporal patterns and provide better classification performance. This particular method of feature extraction has also shown promising result in classifying COVID-19 coughs, breath and speech [9, 48].

4 Classifier Training

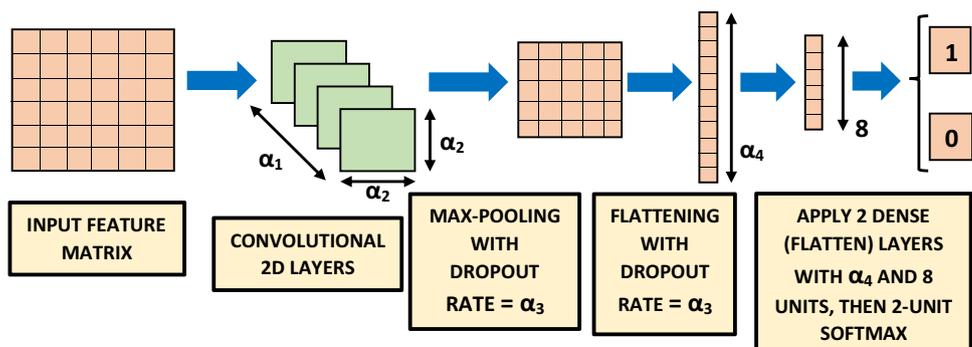
We have trained and evaluated six machine learning classifiers on both audio and accelerometer signals. Table 3 lists the classifier hyperparameters that were optimised during leave-one-out cross-validation.

First, we establish the baseline results by training and evaluating three shallow classifiers: logistic regression (LR), support vector machine (SVM) and multilayer perceptron (MLP). Then, we improve the cough detection performance

Table 3 Classifier hyperparameters, optimised using the leave-one-patient-out cross-validation.

Hyperparameters	Classifier	Range
Regularisation strength (γ_1)	LR, SVM	10^i where, $i = -7, -6, \dots, 6, 7$
$l1$ penalty (γ_2)	LR	0 to 1 in steps of 0.05
$l2$ penalty (γ_3)	LR, MLP	0 to 1 in steps of 0.05
Kernel coefficient (γ_4)	SVM	10^i where, $i = -7, -6, \dots, 6, 7$
No. of neurons (γ_5)	MLP	10 to 100 in steps of 10
Batch size (ξ_1)	CNN & LSTM	2^k where $k = 6, 7, 8$
No. of epochs (ξ_2)	CNN & LSTM	10 to 250 in steps of 20
No. of convolutional filters (α_1)	CNN	3×2^k where $k = 3, 4, 5$
kernel size (α_2)	CNN	2 and 3
Dropout rate (α_3)	CNN & LSTM	0.1 to 0.5 in steps of 0.2
Units in the dense layer (α_4)	CNN & LSTM	2^k where $k = 4, 5$
LSTM units (β_1)	LSTM	2^k where $k = 6, 7, 8$
Learning rate (β_2)	LSTM	10^k where $k = -2, -3, -4$

Figure 5 CNN Classifier, trained and evaluated using leave-one-out cross-validation [64] on 14 patients. The results are shown in Tables 4, 5 and 6 for feature extraction hyperparameters mentioned in Table 2.



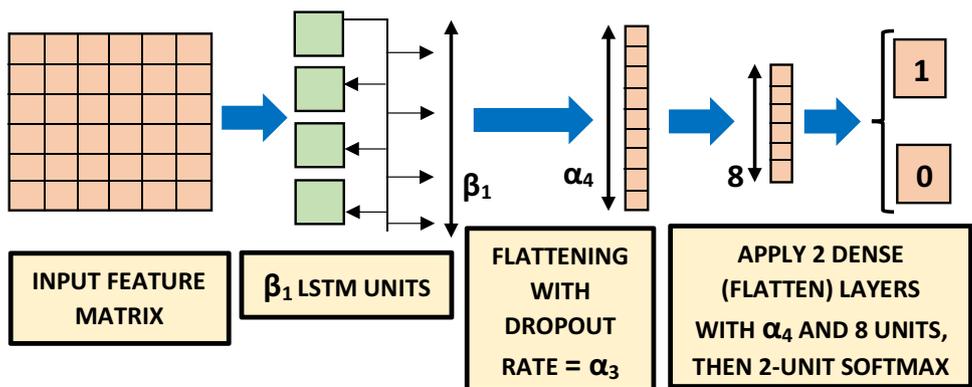
by implementing three deep neural network (DNN) classifiers: CNN, long short-term memory (LSTM) and Resnet50.

LR models have outperformed other more complex classifiers such as classification trees, random forests, SVM in several other clinical prediction tasks [3, 54, 55]. The gradient descent weight regularisation as well as lasso (l_1 penalty) and ridge (l_2 penalty) estimators [56, 57] were the hyperparameters, listed in Table 3, optimised inside the nested cross-validation during training. SVM classifiers have also performed well in both detecting [58, 59] and classifying [60] cough events in the past. The independent term in kernel functions is the hyperparameter optimised for the SVM classifier. An MLP, consisting multiple layers of neurons [61], is capable of learning non-linear relationships. It has produced promising results in discriminating influenza coughs from other coughs [62] in the past. MLP has also been applied to classify TB coughs [47, 59] and detect coughs in general [25, 63]. The penalty ratios, along with the number of neurons are used as the hyperparameters, optimised during leave-one-out cross-validation (Fig. 7 and Sect. 5).

A CNN is a popular deep neural network architecture, primarily used in image classification [65], such as face recognition [66]. It has also performed well in classifying COVID-19 coughs, breath and speech [9, 48]. The CNN architecture [67, 68], shown in Fig. 5, contains α_1 2D convolutional layers with kernel size α_2 and rectified linear units as activation functions. A dropout rate α_3 has been applied

along with max-pooling, followed by two dense (flatten) layers containing α_4 and 8 units (dimensionality of the output space) respectively with rectified linear units as activation functions. An LSTM model is a type of recurrent neural network which remembers previously-seen inputs when making its classification decision [69]. It has been successfully used in automatic cough detection [7, 24], and also in other types of acoustic event detection [70, 71] including COVID-19 coughs etc. [9, 48]. The hyperparameters optimised for the LSTM classifier [72] are mentioned in Table 3 and visually explained in Fig. 6. The LSTM classifier, shown in Fig. 6, contains β_1 LSTM units i.e. the dimensionality of the output space is β_1 with rectified linear units as activation functions and a dropout rate α_3 . Then two dense (flatten) layers containing α_4 and 8 units (dimensionality of the output space) respectively have been applied with rectified linear units as activation functions. For both CNN and LSTM classifiers, a final softmax function produces one output for a cough event (i.e. 1) and the other for a non-cough event (i.e. 0), shown in Figs. 5 and 6. Features are fed into these two classifiers using a batch size of ξ_1 for ξ_2 number of epochs. The 50-layer residual network (Resnet50) architecture (Table 1 of [73]) we trained and evaluated has a very deep architecture that contains skip layers and has performed even better than existing deep architectures such as VGGNet on image classification tasks on the dataset such as ILSVRC, the CIFAR10 dataset and the COCO object detection dataset

Figure 6 LSTM classifier, trained and evaluated using leave-one-out cross-validation [64] on 14 patients. The results are shown in Tables 4, 5 and 6 for feature extraction hyperparameters mentioned in Table 2.



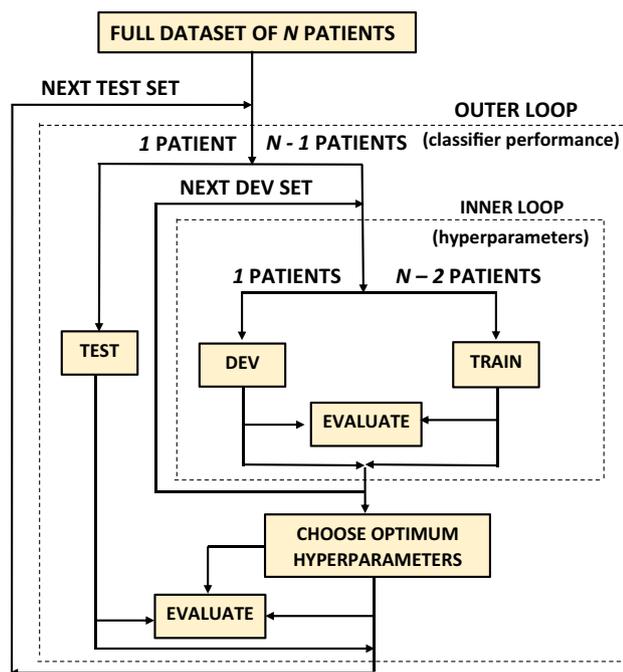


Figure 7 Leave-one-out cross-validation used to train and evaluate all six classifiers. Here, $N = 14$ (Table 1). The development set (DEV) consisting 1 patient has been used to optimise the hyperparameters while training on the TRAIN set, consisted of 12 patients. The final evaluation of the classifiers in terms of the AUC occurs on the TEST set, consisting 1 patient.

[74]. This architecture has also performed the best in detecting COVID-19 signatures in coughs, breaths and speech [9, 48]. Due to extreme computation load, we have used the default Resnet50 structure mentioned in Table 1 of [73].

5 Classification Process

5.1 Hyperparameter Optimisation

Hyperparameters for both the classifiers and feature extraction are optimised inside the leave-one-out cross-validation process and are listed in Tables 2 and 3. Different phases of an event might carry important information and our way of feature extraction preserves the time-domain information. By varying the frame lengths and number of frames to extract, this information was varied. The spectral resolution was also varied by varying the number of lower-order MFCCs to keep from the audio signal.

5.2 Cross-Validation

All six classifiers have been trained and evaluated by using a leave-one-patient-out cross-validation scheme [64], as

explained in Fig. 7. Our dataset contains only 14 patients and by using this cross-validation scheme we make the best use of our dataset, as a patient's weight, coughing intensity and distance from the microphone can affect the accelerometer and audio signals and we were not allowed to collect that vital information due to the ethical constraints.

The Fig. 7 shows that one patient is left out from 14 patients to be used for later independent testing. Then another patient is removed from the remaining 13 patients to be used as the development set where the hyperparameters, listed in Table 3, are to be optimised. AUC has always been the optimisation criterion in this cross-validation. This entire procedure is repeated until all patients are used as an independent test set in the outer loop. The final performance is evaluated by calculating and averaging AUC over these outer loops. The hyperparameters producing the highest AUC over these outer test sets are noted as the 'best hyperparameters' in Tables 4, 5 and 6. Performances produced by each classifier for each set of hyperparameters are noted by 'ID' in these tables.

6 Results

6.1 Accelerometer-Based Cough Detection

Table 4 lists the performance achieved by the shallow classifiers in systems C1 to C18 and Table 5 lists the performance achieved by the deep architectures in systems C19 to C36. These results are the averages over the 14 leave-one-patient-out testing partitions in the outer loop of the nested cross-validation.

The shallow classifiers have provided the baseline classification performance. Table 4 shows that an LR classifier has achieved the best performance of an AUC of 0.8135 along with σ_{AUC} of 0.003, specificity of 81.42%, sensitivity of 81.28% and an accuracy of 81.35% (system C4). The SVM has produced an AUC of 0.8252 with σ_{AUC} of 0.003 and 80.91% specificity, 84.11% sensitivity and 82.51% accuracy while using ten 32 samples long frames (system C10) as its best performance. However, the AUC of 0.8587, accuracy of 85.67%, specificity of 84.47% and sensitivity of 86.89% have been achieved from the MLP classifier with 40 neurons and l_2 penalty ratio of 0.7 using five 64 sample long frames (system C17) and this is the highest AUC achieved by a shallow classifier.

For the DNN classifiers, the lowest AUC of 0.9243 has been achieved from a CNN classifier in system C19 in Table 5. Table 5 also shows that the best-performing CNN uses ten 64 samples (640 msec) long frames to achieve an AUC of 0.9499, accuracy of 85.82%, specificity of 80.91% and sensitivity of 90.73% (system C24). The optimal LSTM classifier achieves the slightly higher AUC of 0.9572 when

Table 4 Accelerometer-based cough detection results for the shallow classifiers. The values are averaged over 14 cross-validation folds. The highest AUC of 0.8587 has been achieved from an MLP classifier.

Classifier	ID	Best Feature Hyperparameters	Best Classifier Hyperparameters (Optimised inside nested cross-validation)	Performance				
				Spec	Sens	Acc	AUC	σ_{AUC}
LR	C1	$\Psi = 16, C = 5$	$\gamma_1 = 10^{-4}, \gamma_2 = 0.35, \gamma_3 = 0.65$	80.41%	80.28%	80.35%	0.8055	0.003
	C2	$\Psi = 16, C = 10$	$\gamma_1 = 10^{-2}, \gamma_2 = 0.55, \gamma_3 = 0.45$	80.25%	80.08%	80.16%	0.8058	0.003
	C3	$\Psi = 32, C = 5$	$\gamma_1 = 10^2, \gamma_2 = 0.2, \gamma_3 = 0.8$	80.39%	80.55%	80.47%	0.8072	0.003
	C4	$\Psi = 32, C = 10$	$\gamma_1 = 10^{-3}, \gamma_2 = 0.4, \gamma_3 = 0.6$	81.42%	81.28%	81.35%	0.8135	0.003
	C5	$\Psi = 64, C = 5$	$\gamma_1 = 10^{-1}, \gamma_2 = 0.25, \gamma_3 = 0.75$	80.22%	80.41%	80.31%	0.8119	0.003
	C6	$\Psi = 64, C = 10$	$\gamma_1 = 10^{-2}, \gamma_2 = 0.75, \gamma_3 = 0.25$	80.16%	80.32%	80.24%	0.8124	0.003
SVM	C7	$\Psi = 16, C = 5$	$\gamma_1 = 10^3, \gamma_4 = 10^{-3}$	80.71%	82.91%	81.81%	0.8202	0.003
	C8	$\Psi = 16, C = 10$	$\gamma_1 = 10^{-2}, \gamma_4 = 10^2$	80.22%	82.94%	81.58%	0.8248	0.003
	C9	$\Psi = 32, C = 5$	$\gamma_1 = 10^3, \gamma_4 = 10^{-2}$	80.41%	82.97%	81.69%	0.8212	0.003
	C10	$\Psi = 32, C = 10$	$\gamma_1 = 10^{-1}, \gamma_4 = 10^{-1}$	80.91%	84.11%	82.51%	0.8252	0.003
	C11	$\Psi = 64, C = 5$	$\gamma_1 = 10^{-4}, \gamma_4 = 10^{-3}$	80.28%	84.35%	82.31%	0.8245	0.003
	C12	$\Psi = 64, C = 10$	$\gamma_1 = 10^2, \gamma_4 = 10^{-4}$	80.55%	82.78%	81.68%	0.8251	0.003
MLP	C13	$\Psi = 16, C = 5$	$\gamma_3 = 0.55, \gamma_5 = 30$	82.37%	86.95%	84.68%	0.8507	0.003
	C14	$\Psi = 16, C = 10$	$\gamma_3 = 0.45, \gamma_5 = 50$	83.24%	87.08%	85.16%	0.8558	0.003
	C15	$\Psi = 32, C = 5$	$\gamma_3 = 0.35, \gamma_5 = 70$	83.55%	87.41%	85.47%	0.8552	0.003
	C16	$\Psi = 32, C = 10$	$\gamma_3 = 0.4, \gamma_5 = 20$	82.18%	86.05%	84.12%	0.8424	0.003
	C17	$\Psi = 64, C = 5$	$\gamma_3 = 0.7, \gamma_5 = 40$	84.47%	86.89%	85.67%	0.8587	0.003
	C18	$\Psi = 64, C = 10$	$\gamma_3 = 0.35, \gamma_5 = 30$	83.45%	86.84%	84.64%	0.8499	0.003

features were extracted using ten 32 samples (320 msec) long frames (system C28). However, the best performance

is achieved by the Resnet50 architecture, with an AUC of 0.9888 after 50 epochs from ten 32 samples (320 msec) long

Table 5 Accelerometer-based cough detection results for the DNN classifiers. The values are averaged over 14 cross-validation folds. DNN classifiers have outperformed the shallow classifiers by a

wide margin and a Resnet50 produces the highest AUC of 0.9888 in detecting cough events.

Classifier	ID	Best Feature Hyperparameters	Best Classifier Hyperparameters (Optimised inside nested cross-validation)	Performance				
				Spec	Sens	Acc	AUC	σ_{AUC}
CNN	C19	$\Psi = 16, C = 5$	$\alpha_1=48, \alpha_2=2, \alpha_3=0.1, \alpha_4=32, \xi_1=128, \xi_2=210$	83.47%	85.62%	84.55%	0.9243	0.002
	C20	$\Psi = 16, C = 10$	$\alpha_1=24, \alpha_2=2, \alpha_3=0.3, \alpha_4=32, \xi_1=256, \xi_2=110$	83.76%	87.56%	85.66%	0.9358	0.002
	C21	$\Psi = 32, C = 5$	$\alpha_1=96, \alpha_2=2, \alpha_3=0.3, \alpha_4=32, \xi_1=128, \xi_2=150$	76.98%	91.96%	84.47%	0.9272	0.002
	C22	$\Psi = 32, C = 10$	$\alpha_1=48, \alpha_2=2, \alpha_3=0.3, \alpha_4=16, \xi_1=256, \xi_2=110$	84.09%	86.41%	85.25%	0.9324	0.002
	C23	$\Psi = 64, C = 5$	$\alpha_1=48, \alpha_2=2, \alpha_3=0.5, \alpha_4=32, \xi_1=256, \xi_2=230$	85.47%	87.15%	86.31%	0.9339	0.002
	C24	$\Psi = 64, C = 10$	$\alpha_1=96, \alpha_2=2, \alpha_3=0.3, \alpha_4=32, \xi_1=128, \xi_2=170$	80.91%	90.73%	85.82%	0.9499	0.002
LSTM	C25	$\Psi = 16, C = 5$	$\beta_1 = 128, \beta_2 = 0.0001, \alpha_3 = 0.3, \alpha_4 = 32, \xi_1 = 256, \xi_2 = 210$	84.34%	90.82%	87.58%	0.9444	0.002
	C26	$\Psi = 16, C = 10$	$\beta_1 = 128, \beta_2 = 0.01, \alpha_3 = 0.1, \alpha_4 = 32, \xi_1 = 128, \xi_2 = 110$	85.37%	91.27%	88.32%	0.9504	0.002
	C27	$\Psi = 32, C = 5$	$\beta_1 = 256, \beta_2 = 0.001, \alpha_3 = 0.3, \alpha_4 = 16, \xi_1 = 128, \xi_2 = 130$	79.92%	94.31%	87.11%	0.9457	0.002
	C28	$\Psi = 32, C = 10$	$\beta_1 = 128, \beta_2 = 0.001, \alpha_3 = 0.1, \alpha_4 = 32, \xi_1 = 256, \xi_2 = 150$	86.41%	92.05%	89.21%	0.9572	0.002
	C29	$\Psi = 64, C = 5$	$\beta_1 = 256, \beta_2 = 0.001, \alpha_3 = 0.3, \alpha_4 = 16, \xi_1 = 128, \xi_2 = 190$	84.57%	92.79%	88.68%	0.954	0.002
	C30	$\Psi = 64, C = 10$	$\beta_1 = 128, \beta_2 = 0.01, \alpha_3 = 0.5, \alpha_4 = 32, \xi_1 = 256, \xi_2 = 230$	86.21%	89.13%	87.66%	0.9489	0.002
Resnet50	C31	$\Psi = 16, C = 5$	Default Resnet50 (Table 1 in [73])	93.81%	97.09%	95.43%	0.9802	0.002
	C32	$\Psi = 16, C = 10$	"	94.12%	98.58%	96.35%	0.9812	0.002
	C33	$\Psi = 32, C = 5$	"	94.29%	98.79%	96.54%	0.9810	0.002
	C34	$\Psi = 32, C = 10$	"	94.09%	99.33%	96.71%	0.9888	0.002
	C35	$\Psi = 64, C = 5$	"	94.71%	98.23%	96.35%	0.9854	0.002
	C36	$\Psi = 64, C = 10$	"	95.07%	97.89%	96.46%	0.9884	0.002

Table 6 Audio-based cough detection results. The values are averaged over 14 cross-validation folds and the best-three performances of each classifier are shown. All classifiers have performed well in

detecting coughs but DNN classifiers have performed particularly well and their performances are very close to each other.

Classifier	ID	Best Feature Hyperparameters	Best Classifier Hyperparameters (Optimised inside nested cross-validation)	Performance				
				Spec	Sens	Acc	AUC	σ_{AUC}
LR	D1	$\mathcal{M} = 26, \mathcal{F} = 512, \mathcal{S} = 100$	$\gamma_1 = 10^{-3}, \gamma_2 = 0.25, \gamma_3 = 0.75$	87.52%	87.71%	87.61%	0.9129	0.003
	D2	$\mathcal{M} = 39, \mathcal{F} = 1024, \mathcal{S} = 70$	$\gamma_1 = 10^2, \gamma_2 = 0.4, \gamma_3 = 0.6$	87.31%	87.41%	87.36%	0.9358	0.003
	D3	$\mathcal{M} = 26, \mathcal{F} = 1024, \mathcal{S} = 100$	$\gamma_1 = 10^{-4}, \gamma_2 = 0.55, \gamma_3 = 0.45$	87.14%	87.28%	87.21%	0.9272	0.003
SVM	D4	$\mathcal{M} = 26, \mathcal{F} = 1024, \mathcal{S} = 120$	$\gamma_1 = 10^{-2}, \gamma_4 = 10^3$	86.75%	86.91%	86.83%	0.9066	0.003
	D5	$\mathcal{M} = 26, \mathcal{F} = 512, \mathcal{S} = 100$	$\gamma_1 = 10^2, \gamma_4 = 10^{-3}$	86.61%	86.68%	86.64%	0.9058	0.003
	D6	$\mathcal{M} = 39, \mathcal{F} = 1024, \mathcal{S} = 100$	$\gamma_1 = 10^3, \gamma_4 = 10^{-3}$	86.40%	86.54%	86.47%	0.9017	0.003
MLP	D7	$\mathcal{M} = 39, \mathcal{F} = 2048, \mathcal{S} = 120$	$\gamma_3 = 0.35, \gamma_5 = 30$	89.47%	90.10%	89.78%	0.9254	0.002
	D8	$\mathcal{M} = 26, \mathcal{F} = 1024, \mathcal{S} = 70$	$\gamma_3 = 0.4, \gamma_5 = 50$	89.55%	89.76%	89.66%	0.9214	0.003
	D9	$\mathcal{M} = 39, \mathcal{F} = 1024, \mathcal{S} = 100$	$\gamma_3 = 0.6, \gamma_5 = 40$	88.78%	89.04%	88.91%	0.9205	0.003
LSTM	D10	$\mathcal{M} = 26, \mathcal{F} = 1024, \mathcal{S} = 70$	$\beta_1 = 128, \beta_2 = 0.001, \alpha_3 = 0.3, \alpha_4 = 32, \xi_1 = 256, \xi_2 = 210$	94.57%	96.59%	95.58%	0.9932	0.002
	D11	$\mathcal{M} = 39, \mathcal{F} = 1024, \mathcal{S} = 100$	$\beta_1 = 128, \beta_2 = 0.001, \alpha_3 = 0.3, \alpha_4 = 16, \xi_1 = 256, \xi_2 = 130$	94.21%	94.21%	96.43%	0.9904	0.002
	D12	$\mathcal{M} = 26, \mathcal{F} = 2048, \mathcal{S} = 120$	$\beta_1 = 128, \beta_2 = 0.01, \alpha_3 = 0.3, \alpha_4 = 32, \xi_1 = 128, \xi_2 = 170$	93.95%	96.25%	95.10%	0.9857	0.002
CNN	D13	$\mathcal{M} = 26, \mathcal{F} = 1024, \mathcal{S} = 100$	$\alpha_1 = 48, \alpha_2 = 2, \alpha_3 = 0.3, \alpha_4 = 32, \xi_1 = 256, \xi_2 = 90$	93.24%	97.88%	95.56%	0.9944	0.002
	D14	$\mathcal{M} = 13, \mathcal{F} = 512, \mathcal{S} = 70$	$\alpha_1 = 24, \alpha_2 = 2, \alpha_3 = 0.3, \alpha_4 = 16, \xi_1 = 256, \xi_2 = 170$	92.18%	98.74%	95.46%	0.9891	0.002
	D15	$\mathcal{M} = 39, \mathcal{F} = 2048, \mathcal{S} = 120$	$\alpha_1 = 48, \alpha_2 = 2, \alpha_3 = 0.1, \alpha_4 = 32, \xi_1 = 256, \xi_2 = 130$	92.78%	97.56%	95.17%	0.9872	0.002
Resnet50	D16	$\mathcal{M} = 26, \mathcal{F} = 1024, \mathcal{S} = 100$	Default Resnet50 (Table 1 in [73])	96.74%	99.50%	98.13%	0.9957	0.001
	D17	$\mathcal{M} = 39, \mathcal{F} = 1024, \mathcal{S} = 70$	"	96.55%	96.95%	96.75%	0.9912	0.001
	D18	$\mathcal{M} = 39, \mathcal{F} = 512, \mathcal{S} = 100$	"	94.70%	96.46%	95.58%	0.9891	0.002

frames along with 96.71% accuracy, 94.09% specificity and 99.33% sensitivity (system C34).

Deep architectures have produced a higher AUCs and lower σ_{AUC} than the shallow classifiers on accelerometer-based classification task. Figure 8 shows the mean ROC curves for the optimal LR, SVM, MLP, CNN, LSTM and Resnet50, whose configurations are shown in Table 3 and the mean AUCs were calculated over the 14 cross-validation folds. The Resnet50 classifier is superior to all other classifiers over a wide range of operating points (Fig. 8).

6.2 Audio-Based Cough Detection

To place the performance of the accelerometer-based cough detection presented in the previous section into perspective, we have performed a matching set of experiments, this time using the audio signals to perform audio-based cough detection. These experiments are based on precisely the same events as the acceleration experiments, since our corpus contains both audio and acceleration signals for each.

Table 6 shows the best-three configurations for each of the six classifier architectures in systems D1 to D18. Again,

the results indicate that shallow classifiers (LR, SVM and MLP) achieve good classification scores.

LR achieved the highest AUC of 0.9129 with σ_{AUC} of 0.003 when using 26 MFCCs, 512 sample long frames and extracting 100 frames (system D1). The system has also generated the specificity of 87.52%, sensitivity of 87.71% and an accuracy of 87.61%. The SVM achieved an AUC of 0.9066 with σ_{AUC} of 0.003 for 26 MFCCs, 1024 sample long frames and extracting 120 frames (system D4). The system has also generated the specificity of 86.75%, sensitivity of 86.91% and an accuracy of 86.83%. An MLP has produced the highest AUC of 0.9254 with σ_{AUC} of 0.002 for 39 MFCCs, 2048 sample long frames and extracting 120 frames (system D7). The system has also generated the specificity of 89.47%, sensitivity of 90.10% and an accuracy of 89.78%. This is the best performance achieved by the shallow classifiers.

Again, the DNN classifiers have outperformed the shallow classifiers by a large margin. The best LSTM classifier has produced the highest AUC of 0.9932 with σ_{AUC} of 0.002 for 26 MFCCs, 1024 sample long frames and extracting 70 frames (system D10). The system has also generated the specificity of 94.57%, sensitivity of 96.59% and an accuracy of 95.58%. The best CNN classifier has produced the highest AUC of

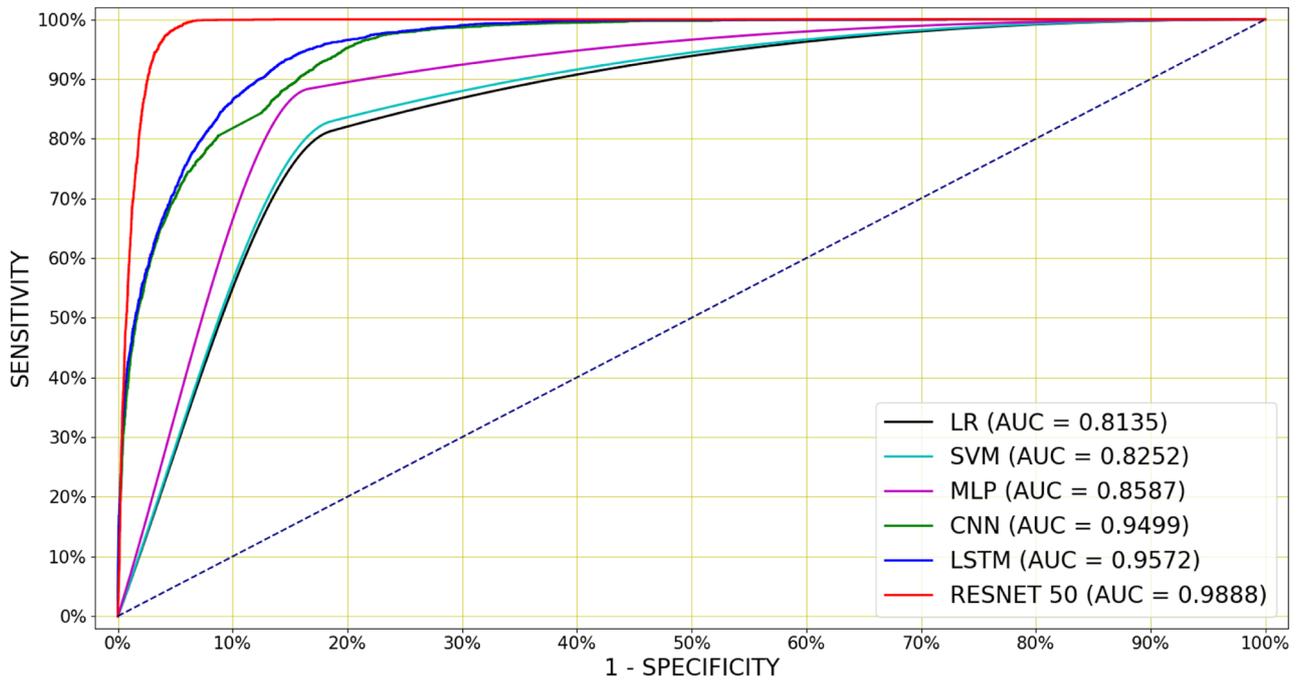


Figure 8 Mean ROC curves for accelerometer-based cough detection, for the best performing classifiers whose hyperparameters are mentioned in Table 3. A Resnet50 has performed the best outper-

forming all the other classifiers over a wide range of operating points by achieving the AUC of 0.9888 and the accuracy of 96.71%.

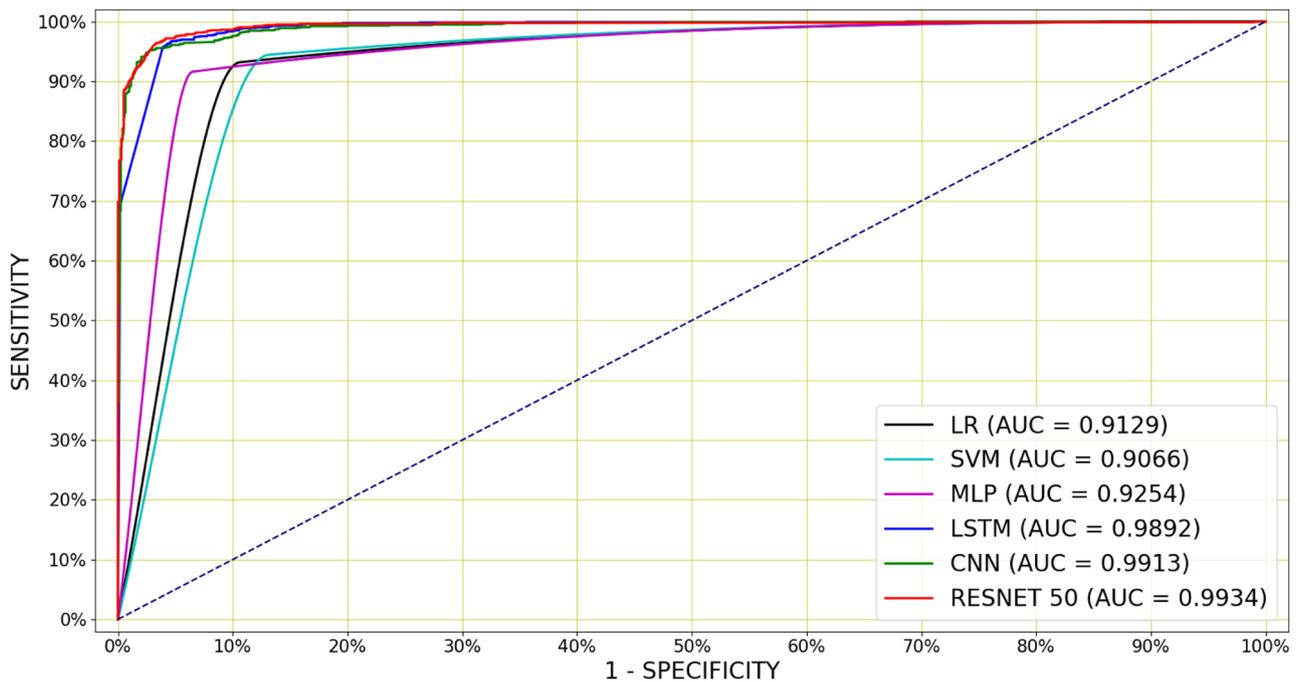


Figure 9 Mean ROC curves for audio-based cough detection, for the best performing classifiers whose hyperparameters are mentioned in Table 3. The best performance has been achieved from a

Resnet50 and similar performances have been achieved from a CNN and LSTM. The best Resnet50 produces the AUC of 0.9957 and the accuracy of 98.13%.

0.9944 with σ_{AUC} of 0.002 while features were extracted for 26 MFCCs, 1024 sample long frames and extracting 100 frames (system D13). The system has also generated the specificity of 93.24%, sensitivity of 97.88% and an accuracy of 95.56%. However, the highest AUC of 0.9957 has been achieved again from a Resnet50 classifier with a σ_{AUC} of 0.001 for 26 MFCCs, 1024 sample (i.e. 46.44 msec) long frames and extracting 100 frames from the entire event (system D16). This system has also achieved a specificity of 96.74% and a sensitivity of 99.5% along with the accuracy of 98.13%.

Again, deep architectures have produced higher AUCs and lower σ_{AUC} than the shallow classifiers on audio-based classification task. These best results for audio-based classification are shown in Fig. 9. Table 6 also indicate that the number of MFCCs has been varied between 13 and 65, although the best performance was achieved by using 26 and 39 MFCCs. Using the frame length of 1024 and extracting 100 frames from the events has provided the best performance for most of the classifiers.

7 Discussion

The results shown in Tables 4, 5 and 6 indicate that audio-based cough detection is consistently more accurate than accelerometer-based classification. However, it is interesting to note that the performances offered by the two alternatives are fairly close. In fact, the deep architectures like Resnet50 offer almost equal performance for audio-based and accelerometer-based cough detection. It also seems that the CNN and LSTM find it easier to classify cough events based on audio rather than accelerometer signal. We postulate that this is due to the limited range of time and frequency information contained in accelerometer data, which is in turn due to the lower accelerometer sampling rate.

For acceleration signals, the extraction of 10 frames each with a length of 640 ms produced the best result. For the audio, the extraction of 100 frames each with a length of 46.44 ms provided the optimal performance. These audio frame lengths are close to those traditionally used for feature extraction in automatic speech recognition. We also note that the performances of the deep classifiers are consistently better than those offered by the baseline shallow classifiers for both types of signals. Although the datasets differ, our system also appears to improve on recent work using the accelerometer integrated into a smartwatch [75] and distinguishing cough from other audio events such as sneeze, speech and noise [7].

8 Conclusion and Future Work

We have demonstrated that an automatic non-invasive machine learning based cough detector is able to accurately discriminate between the accelerometer and audio signals due to coughing

and due to other movements as captured by a consumer smart-phone attached to a patient's bed.

We have trained and evaluated six classifiers including three shallow classifiers: logistic regression (LR), support vector machine (SVM) and multilayer perceptron (MLP) and three deep neural network (DNN) classifiers: convolutional neural networks (CNN), long-short-term-memory (LSTM) networks, and a 50-layer residual-based neural network architecture (Resnet50). A specially-compiled corpus of manually-annotated acceleration and audio events, including approximately 6000 cough and 68000 non-cough events such as sneezing, throat-clearing and getting in and out of the bed, gathered from 14 adult male patients in a small TB clinic was used to train and evaluate these classifiers by using a leave-one-out cross-validation scheme. For accelerometer-based classification, the best system uses a Resnet50 architecture and produces an AUC of 0.9888 as well as a 96.71% accuracy, 94.09% specificity and 99.33% sensitivity while features were extracted from ten 32 sample (320 msec) long frames. This demonstrates that it is possible to discriminate between cough events and other non-cough events by using very deep architectures such as a Resnet50; based on signals gathered from an accelerometer that is not attached to the patient's body, but rather to the headboard of the patient's bed.

We have also compared this accelerometer-based cough detection with audio-based cough detection for the same cough and non-cough events. For audio-based classification, the best result has also been achieved from a Resnet50 with the highest AUC of 0.9957. This shows that the accelerometer-based cough detection is almost equally accurate as audio-based classification while using very deep architectures such as a Resnet50. Shallow classifiers and DNN such as CNN and LSTM however perform better in classifying cough events on audio signals rather than accelerometer signals, as audio signal carries more dynamic and diverse frequency content.

Accelerometer-based detection of cough events has successfully been considered before due to its lower sampling rates and lesser demand of high processing power, however only by using sensors worn by the subjects, which is intrusive and can be inconvenient in some respects. This study shows that excellent discrimination is also possible when the sensor is attached to the patient's bed, thus providing a less intrusive and more convenient solution. Furthermore, since the use of acceleration signal avoids the need to gather audio, privacy is inherently protected. Therefore, the use of a bed-mounted accelerometer inside an inexpensive consumer smartphone may represent a more convenient, cost-effective and readily accepted method of long-term patient cough monitoring.

In the future, we will be attempting to optimise some of the Resnet50 metaparameters and fuse both audio and accelerometer signal to achieve higher specificity and accuracy in cough detection. We are also in the process of applying

the proposed system in an automatic non-invasive cough monitoring system. We also note, the manually annotated cough events sometimes contains multiple bursts of cough onsets and we are currently investigating automatic methods that allow such bursts within a cough event to be identified.

Acknowledgements We would like to thank the South African Centre for High Performance Computing (CHPC) for providing computational resources on their Lengau cluster and Lungenliga St. Gallen-Appenzell, Switzerland for providing funding for this research. We thank M J Mlungisi for his invaluable support in data collection, and gratefully acknowledge the support of Telkom South Africa.

References

- Korpáš, J., Sadloňová, J., & Vrabec, M. (1996). Analysis of the cough sound: an overview. *Pulmonary Pharmacology*, 9, 261–268.
- Knocikova, J., Korpas, J., Vrabec, M., & Javorka, M. (2008). Wavelet analysis of voluntary cough sound in patients with respiratory diseases. *Journal of Physiology and Pharmacology*, 59, 331–40.
- Botha, G., Theron, G., Warren, R., Klopper, M., Dheda, K., Van Helden, P., & Niesler, T. (2018). Detection of tuberculosis by automatic cough sound analysis. *Physiological Measurement*, 39, 045005.
- Al-khassaweneh, M., & Bani Abdelrahman, R. (2013). A signal processing approach for the diagnosis of asthma from cough sounds. *Journal of Medical Engineering & Technology*, 37, 165–171.
- Pramono, R. X. A., Imtiaz, S. A., & Rodriguez-Villegas, E. (2016). A cough-based algorithm for automatic diagnosis of pertussis. *PLoS One*, 11, e0162128.
- Carfì, A., Bernabei, R., Landi, F., et al. (2020). Persistent symptoms in patients after acute COVID-19. *JAMA*, 324, 603–605.
- Miranda, I. D., Diacon, A. H., & Niesler, T. R. (2019). A comparative study of features for acoustic cough detection using deep architectures. In *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)* (pp. 2601–2605).
- Laguarta, J., Hueto, F., & Subirana, B. (2020). COVID-19 Artificial Intelligence Diagnosis using only Cough Recordings. *IEEE Open Journal of Engineering in Medicine and Biology*, 1, 275–281.
- Pahar, M., Klopper, M., Warren, R., & Niesler, T. (2021). COVID-19 cough classification using machine learning and global smartphone recordings. *Computers in Biology and Medicine*, 135, 104572.
- Chen, F., Adcock, J., & Krishnagiri, S. (2008). Audio privacy: reducing speech intelligibility while preserving environmental sounds. In *Proceedings of the 16th ACM International Conference on Multimedia* (pp. 733–736).
- Tung, Y.-C., & Shin, K. G. (2019). Exploiting Sound Masking for Audio Privacy in Smartphones. In *Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security* (pp. 257–268).
- Xia, S., & Jiang, X. (2020). PAMS: Improving Privacy in Audio-Based Mobile Systems. In *Proceedings of the 2nd International Workshop on Challenges in Artificial Intelligence and Machine Learning for Internet of Things* (pp. 41–47).
- Liaqat, D., Nemati, E., Rahman, M., & Kuang, J. (2017). A method for preserving privacy during audio recordings by filtering speech. In *2017 IEEE Life Sciences Conference (LSC)* (pp. 79–82). IEEE.
- Mehigan, T. J. (2009). Harnessing accelerometer technology for inclusive mobile learning. In *Proceedings of the 11th International Conference on Human-Computer Interaction with Mobile Devices and Services* (pp. 1–2).
- Mohammadi, H., Samadani, A.-A., Steele, C., & Chau, T. (2019). Automatic discrimination between cough and non-cough accelerometry signal artefacts. *Biomedical Signal Processing and Control*, 52, 394–402.
- Munyard, P., Busst, C., Logan-Sinclair, R., & Bush, A. (1994). A new device for ambulatory cough recording. *Pediatric Pulmonology*, 18, 178–186.
- Coyle, M., Derchak, P. A., & Myers, L. J. (2010). Systems and methods for monitoring cough. US Patent 7,727,161.
- Fan, J., Comina, G., Gilman, R., Lopez, J., & Tracey, B. H. (2014). Cough monitoring for pulmonary tuberculosis using combined microphone/accelerometer measurements. *The Journal of the Acoustical Society of America*, 135, 2268–2268.
- Pavesi, L., Subburaj, S., & Porter-Shaw, K. (2001). Application and validation of a computerized cough acquisition system for objective monitoring of acute cough: a meta-analysis. *Chest*, 120, 1121–1128.
- Paul, I. M., Wai, K., Jewell, S. J., Shaffer, M. L., & Varadan, V. V. (2006). Evaluation of a new self-contained, ambulatory, objective cough monitor. *Cough*, 2, 7.
- Chan, J. Y. M., Tunnell, S. A., & Jacobs, J. A. L. (2014). Systems, methods and kits for measuring cough and respiratory rate using an accelerometer. US Patent App. 13/783,257.
- Hirai, K., Tabata, H., Hirayama, M., Kobayashi, T., Oh, Y., & Mochizuki, H. (2015). A new method for objectively evaluating childhood nocturnal cough. *Pediatric Pulmonology*, 50, 460–468.
- Drugman, T., Urbain, J., Bauwens, N., Chessini, R., Valderrama, C., Lebecque, P., & Dutoit, T. (2013). Objective Study of Sensor Relevance for Automatic Cough Detection. *IEEE Journal of Biomedical and Health Informatics*, 17, 699–707.
- Pahar, M., Miranda, I., Diacon, A., & Niesler, T. (2021). Deep Neural Network based Cough Detection using Bed-mounted Accelerometer Measurements. In *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 8002–8006).
- Amoh, J., & Odame, K. (2015). DeepCough: A deep convolutional neural network in a wearable cough detection system. In *2015 IEEE Biomedical Circuits and Systems Conference (BioCAS)* (pp. 1–4). IEEE.
- Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., & Sloetjes, H. (2006). ELAN: a professional framework for multimodality research. In *5th International Conference on Language Resources and Evaluation (LREC 2006)*.
- Van Hulse, J., Khoshgoftaar, T. M., & Napolitano, A. (2007). Experimental perspectives on learning from imbalanced data. In *Proceedings of the 24th International Conference on Machine Learning* (pp. 935–942).
- Krawczyk, B. (2016). Learning from imbalanced data: open challenges and future directions. *Progress in Artificial Intelligence*, 5, 221–232.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16, 321–357.
- Lemaître, G., Nogueira, F., & Aridas, C. K. (2017). Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *The Journal of Machine Learning Research*, 18, 559–563.
- Windmon, A., Minakshi, M., Bharti, P., Chellappan, S., Johansson, M., Jenkins, B. A., & Athilingam, P. R. (2018). Tussiswatch: A smart-phone system to identify cough episodes as early symptoms of chronic obstructive pulmonary disease and congestive heart failure. *IEEE Journal of Biomedical and Health Informatics*, 23, 1566–1573.
- Bingham, C., Godfrey, M., & Tukey, J. (1967). Modern techniques of power spectrum estimation. *IEEE Transactions on Audio and Electroacoustics*, 15, 56–66.

33. Jung, T.-P., Makeig, S., Stensmo, M., & Sejnowski, T. J. (1997). Estimating alertness from the EEG power spectrum. *IEEE Transactions on Biomedical Engineering*, *44*, 60–69.
34. Durak, L., & Arikan, O. (2003). Short-time fourier transform: two fundamental properties and an optimal implementation. *IEEE Transactions on Signal Processing*, *51*, 1231–1242.
35. Sinha, R. (2003). Artificial neural network detects changes in electro-encephalogram power spectrum of different sleep-wake states in an animal model of heat stress. *Medical and Biological Engineering and Computing*, *41*, 595–600.
36. Liang, B., Iwnicki, S., & Zhao, Y. (2013). Application of power spectrum, cepstrum, higher order spectrum and neural network analyses for induction motor fault diagnosis. *Mechanical Systems and Signal Processing*, *39*, 342–360.
37. Levinson, N. (1947). The wiener rms (root mean square) error criterion in filter design and prediction. *Journal of Mathematical Physics*, *25*, 261–278.
38. Lux, R. L., Sower, C. T., Allen, N., Etheridge, S. P., Tristani-Firouzi, M., & Saarel, E. V. (2014). The application of root mean square electrocardiography (RMS ECG) for the detection of acquired and congenital long qt syndrome. *PLoS One*, *9*, e85689.
39. Gilmore, J., Islam, M., Duncan, J., Natu, R., & Martinez-Duarte, R. (2017). Assessing the importance of the root mean square (RMS) value of different waveforms to determine the strength of a dielectrophoresis trapping force. *Electrophoresis*, *38*.
40. Zhang, X., Barkhaus, P. E., Rymer, W. Z., & Zhou, P. (2013). Machine learning for supporting diagnosis of amyotrophic lateral sclerosis using surface electromyogram. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, *22*, 96–103.
41. Yuan, X., Tan, Q., Lei, X., Yuan, Y., & Wu, X. (2017). Wind power prediction using hybrid autoregressive fractionally integrated moving average and least square support vector machine. *Energy*, *129*, 122–137.
42. Lepine, J., Rouillard, V., & Sek, M. (2017). On the use of machine learning to detect shocks in road vehicle vibration signals. *Packaging Technology and Science*, *30*, 387–398.
43. Ren, L., Cui, J., Sun, Y., & Cheng, X. (2017). Multi-bearing remaining useful life collaborative prediction: A deep learning approach. *Journal of Manufacturing Systems*, *43*, 248–256.
44. Wei Han, Cheong-Fat Chan, Chiu-Sing Choy, & Kong-Pang Pun (2006). An efficient MFCC extraction method in speech recognition. In *IEEE International Symposium on Circuits and Systems*.
45. Pahar, M., & Smith, L. S. (2020). Coding and Decoding Speech using a Biologically Inspired Coding System. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)* (pp. 3025–3032). IEEE.
46. Chatzarrain, H., Arcelus, A., Goubran, R., & Knoefel, F. (2011). Feature extraction for the differentiation of dry and wet cough sounds. In *2011 IEEE International Symposium on Medical Measurements and Applications* (pp. 162–166). IEEE.
47. Pahar, M., Klopper, M., Reeve, B., Warren, R., Theron, G., & Niesler, T. (2021). Automatic cough classification for tuberculosis screening in a real-world environment. *Physiological Measurement*, *42*, 105014.
48. Pahar, M., Klopper, M., Warren, R., & Niesler, T. (2022). COVID-19 detection in cough, breath and speech using deep transfer learning and bottleneck features. *Computers in Biology and Medicine*, *141*, 105153.
49. Azmy, M. M. (2017). Feature extraction of heart sounds using velocity and acceleration of MFCCs based on support vector machines. In *2017 IEEE Jordan Conference on Applied Electrical Engineering and Computing Technologies (AEECT)* (pp. 1–4).
50. Bachu, R., Kopparthi, S., Adapa, B., & Barkana, B. D. (2010). Voiced/Unvoiced Decision for Speech Signals Based on Zero-Crossing Rate and Energy. In *Advanced Techniques in Computing Sciences and Software Engineering*, (pp. 279–282).
51. DeCarlo, L. T. (1997). On the meaning and use of kurtosis. *Psychological Methods*, *2*, 292.
52. Takahashi, G., Yamada, T., Makino, S., & Ono, N. (2016). Acoustic scene classification using deep neural network and frame-concatenated acoustic feature. *Detection and Classification of Acoustic Scenes and Events*.
53. Joder, C., Essid, S., & Richard, G. (2009). Temporal Integration for Audio Classification With Application to Musical Instrument Classification. *IEEE Transactions on Audio, Speech, and Language Processing*, *17*, 174–186.
54. Christodoulou, E., Ma, J., Collins, G. S., Steyerberg, E. W., Verbakel, J. Y., & Van Calster, B. (2019). A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *Journal of Clinical Epidemiology*, *110*, 12–22.
55. Le Cessie, S., & Van Houwelingen, J. C. (1992). Ridge estimators in logistic regression. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, *41*, 191–201.
56. Tsuruoka, Y., Tsujii, J., & Ananiadou, S. (2009). Stochastic gradient descent training for l1-regularized log-linear models with cumulative penalty. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP* (pp. 477–485).
57. Yamashita, H., & Yabe, H. (2003). An interior point method with a primal-dual quadratic barrier penalty function for nonlinear optimization. *SIAM Journal on Optimization*, *14*, 479–499.
58. Bhateja, V., Taqee, A., & Sharma, D. K. (2019). Pre-Processing and Classification of Cough Sounds in Noisy Environment using SVM. In *2019 4th International Conference on Information Systems and Computer Networks (ISCON)* (pp. 822–826). IEEE.
59. Tracey, B. H., Comina, G., Larson, S., Bravard, M., López, J. W., & Gilman, R. H. (2011). Cough detection algorithm for monitoring patient recovery from pulmonary tuberculosis. In *2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society* (pp. 6017–6020). IEEE.
60. Sharan, R. V., Abeyratne, U. R., Swarnkar, V. R., & Porter, P. (2017). Cough sound analysis for diagnosing croup in pediatric patients using biologically inspired features. In *2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)* (pp. 4578–4581). IEEE.
61. Taud, H., & Mas, J. (2018). Multilayer perceptron (MLP). *Geomatic Approaches for Modeling Land Change Scenarios*, (pp. 451–455).
62. Sarangi, L., Mohanty, M. N., & Pattanayak, S. (2016). Design of MLP Based Model for Analysis of Patient Suffering from Influenza. *Procedia Computer Science*, *92*, 396–403.
63. Liu, J.-M., You, M., Wang, Z., Li, G.-Z., Xu, X., & Qiu, Z. (2014). Cough detection using deep neural networks. In *2014 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)* (pp. 560–563). IEEE.
64. Sammut, C., & Webb, G. I. (2010). Leave-one-out cross-validation. *Encyclopedia of Machine Learning*, (pp. 600–601).
65. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2017). Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, *60*, 84–90.
66. Lawrence, S., Giles, C. L., Tsoi, A. C., & Back, A. D. (1997). Face recognition: A convolutional neural-network approach. *IEEE Transactions on Neural Networks*, *8*, 98–113.
67. Albawi, S., Mohammed, T. A., & Al-Zawi, S. (2017). Understanding of a convolutional neural network. In *2017 International Conference on Engineering and Technology (ICET)* (pp. 1–6). IEEE.
68. Qi, X., Wang, T., & Liu, J. (2017). Comparison of support vector machine and softmax classifiers in computer vision. In *2017 Second International Conference on Mechanical, Control and Computer Engineering (ICMCCE)* (pp. 151–155). IEEE.

69. Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9, 1735–1780.
70. Marchi, E., Vesperini, F., Weninger, F., Eyben, F., Squartini, S., & Schuller, B. (2015). Non-linear prediction with LSTM recurrent neural networks for acoustic novelty detection. In *2015 International Joint Conference on Neural Networks (IJCNN)* (pp. 1–7). IEEE.
71. Amoh, J., & Odame, K. (2016). Deep neural networks for identifying cough sounds. *IEEE Transactions on Biomedical Circuits and Systems*, 10, 1003–1011.
72. Sherstinsky, A. (2020). Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. *Physica D: Nonlinear Phenomena*, 404, 132306.
73. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 770–778).
74. Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., & Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In *European Conference on Computer Vision* (pp. 740–755). Springer.
75. Liaqat, D., Liaqat, S., Chen, J. L., Sedaghat, T., Gabel, M., Rudzicz, F., & de Lara, E. (2021). Coughwatch: Real-World Cough Detection using Smartwatches. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 8333–8337). IEEE.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Madhurananda Pahar received his BSc in Mathematics from University of Calcutta, India; MSc in Computing for Financial Markets & PhD in Computational Neuroscience from University of Stirling, Scotland. Currently he is working as a post-doctoral fellow in the University of Stellenbosch, South Africa. His research interests are in machine learning and signal processing for audio signals and smart sensors in bio-medicine. Currently he is involved in the

application of deep learning to the detection and classification of TB and COVID coughs in real-world environments as well as the monitoring of patient behaviour using smart sensors such as an accelerometer.



ate at the University of Stellenbosch. His research interests lie in the areas of signal processing, machine learning and VLSI design.



Andreas Diacon is a physician with specialist training in Internal Medicine and Pulmonology. He has keen interest in clinical tuberculosis research, particularly to bring new anti-tuberculosis agents to the patients that need them most quickly and safely. Originating from Switzerland he obtained his PhD from Stellenbosch University and is now leading TASK in Cape Town to advance innovative diagnostics, vaccines and treatments for tuberculosis and

other diseases common in South Africa.



Thomas Niesler obtained the B.Eng (1991) and M.Eng (1993) degrees in Electronic Engineering from the University of Stellenbosch, South Africa and a Ph.D. from the University of Cambridge, England, in 1998. He joined the Department of Engineering, University of Cambridge, as a lecturer in 1998 and subsequently the Department of Electrical and Electronic Engineering, University of Stellenbosch, in 2000, where he has been Professor since 2012. His research interests lie in the areas of signal processing,

pattern recognition and machine learning.