

©20XX IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

# WAKE-COUGH: COUGH SPOTTING AND COUGHER IDENTIFICATION FOR PERSONALISED LONG-TERM COUGH MONITORING

*M Pahar*<sup>1</sup>, *M Klopper*<sup>2</sup>, *B Reeve*<sup>2</sup>, *R Warren*<sup>2</sup>, *G Theron*<sup>2</sup>, *A Diacon*<sup>3</sup>, *T R Niesler*<sup>1</sup>

<sup>1</sup>Department of Electrical and Electronic Engineering, Stellenbosch University, South Africa

<sup>2</sup>SAMRC Centre for Tuberculosis Research, Faculty of Medicine and Health Sciences, Stellenbosch University, South Africa

<sup>3</sup>TASK Applied Science, Cape Town, South Africa

{mpahar, marisat, byronreeve, rw1, gtheron, ahd, trn}@sun.ac.za

## ABSTRACT

We present ‘wake-cough’, an application of wake-word spotting to coughs using Resnet50 and identifying coughers using i-vectors, for the purpose of a long-term, personalised cough monitoring system. Coughs, recorded in a quiet ( $73\pm 5$  dB) and noisy ( $34\pm 17$  dB) environment, were used to extract i-vectors, x-vectors and d-vectors, used as features to the classifiers. The system achieves 90.02% accuracy from an MLP to discriminate 51 coughers using 2-sec long cough segments in the noisy environment. When discriminating between 5 and 14 coughers using longer (100 sec) segments in the quiet environment, this accuracy rises to 99.78% and 98.39% respectively. Unlike speech, i-vectors outperform x-vectors and d-vectors in identifying coughers. These coughs were added as an extra class in the Google Speech Commands dataset and features were extracted by preserving the end-to-end time-domain information in an event. The highest accuracy of 88.58% is achieved in spotting coughs among 35 other trigger phrases using a Resnet50. Wake-cough represents a personalised, non-intrusive, cough monitoring system, which is power-efficient as using wake-word detection method can keep a smartphone-based monitoring device mostly dormant. This makes wake-cough extremely attractive in multi-bed ward environments to monitor patient’s long-term recovery from lung ailments such as tuberculosis and COVID-19.

**Index Terms**— x-vector, i-vector, d-vector, cougher identification, keyword spotting

## 1. INTRODUCTION

Wake-words (WW) are used as trigger phrases which enable keyword spotting (KWS) systems to initiate certain tasks such as speech recognition, providing a bridge between the end-user and either the device or the cloud [1]. For example, some widely-used trigger phrases for voice assistants on smart devices are: Google’s ‘OK Google’, Apple’s ‘Hey Siri’, Amazon’s ‘Alexa’, Microsoft’s ‘Hey Cortana’ [2] and they are highly sensitive in both quiet and noisy environment [3], making them extremely useful in hands-free situations like driving [4]. Coughing is the forceful expulsion of air to clear up the

airway and a common symptom of respiratory diseases, such as tuberculosis (TB) [5], asthma [6], pertussis [7], COVID-19 [8], which can be identified using machine learning classifiers. To successfully implement cough as a WW in commercial smartphones, it is necessary to accurately identify the cougher [9] in a noisy and quiet environment and the cough among various other commonly used trigger phrases [10].

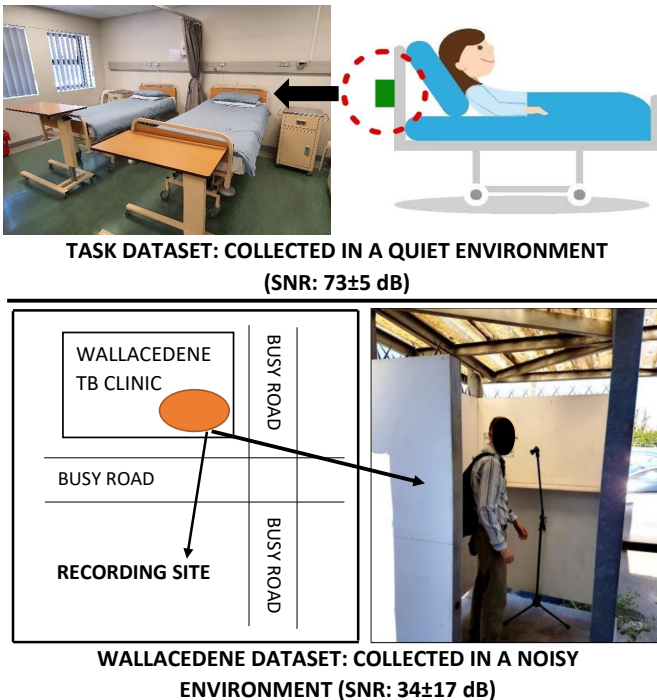
Vocal audio such as speech can be identified using i-vectors, which present a low-dimensional speaker and channel dependant space using factor analysis [11]. The performance can be improved by using x-vectors [12] and d-vectors [13], which use the data augmentation and DNN based embeddings to map speaker embeddings.

Coughers have been identified using x-vectors on natural coughs in an open world environment for 8 male and 8 female subjects after implementing data augmentation to address the issue of background noise [14] and using d-vectors on forced coughs [15]. Here, we identify both natural and forced coughs among other trigger phrases in the Google Speech commands dataset [16] while also identifying the coughers in noisy and quiet environment using i-vectors, x-vectors and d-vectors.

## 2. DATASET PREPARATION

For the cougher identification task, two datasets which will be referred to as TASK and Wallacedene (Table 1), were both manually annotated using ELAN [17]. The TASK dataset was collected inside a multi-bed ward at a 24h tuberculosis (TB) clinic near Cape Town, South Africa and contains natural coughs [18]. A plastic enclosure, attached to the bed-frames, holds a Samsung Galaxy J4 smartphone connected to a BOYA BY-MM1 cardioid microphone (Figure 1) and the distance from the cougher and the microphone was between 30 and 150 cm. The dataset includes 6000 cough events, sampled at 22.05 kHz and collected from 14 adult male patients over a 6 month period, totalling 3.16 hours of cough audio with an average SNR of  $73 \pm 5$  dB. No other information of the patients was collected due to ethical constraints. Wallacedene dataset was collected inside an outdoor booth next to a busy primary health clinic in Wallacedene, near Cape Town,

South Africa representing a real-world environment where a TB test would likely to be deployed [19] (Figure 1). Patients were asked to count from 1 to 10, then cough, take a few deep breaths, and cough again, thus producing forced coughs. These counts were used as speech to provide a baseline to compare the performance of cougher identification. The audio, sampled at 44.1 kHz, was recorded using a RØDE M3 condenser microphone from 38 males and 13 females, keeping a 10 to 15 cm gap between the microphone and the patients. The environmental noise was present in both cough and speech, having the average SNR of 34 dB and 33 dB respectively with a standard deviation of 17 dB (Table 1).



**Fig. 1. Data collection process for cougher identification:** TASK, containing only coughs, was collected in a quiet environment. Wallacecene, containing both cough and speech (counting from 1 to 10), was collected in a noisy environment.

Table 1 shows that the TASK dataset is less-noisy, contains much longer cough audio for each subject, whereas Wallacecene dataset is noisier but has cough and speech audio from a larger number of subjects. All audio recordings were downsampled to 16 kHz, required for kaldi ASR system [20].

For cough spotting, we randomly selected 3795 coughs from the TASK and Wallacecene datasets. Each cough was normalised to a 1-sec duration by either trimming or padding with silence and downsampled to 16 kHz. These ‘cough’ events were added as an extra class to the 2nd version of Google Speech Commands dataset containing 1-sec long 109,624 events, sampled at 16 kHz, belonging to 35 classes [16]. These events were mixed with the background noises (Section 5.8 of [16]) with a randomly selected SNR between 73 and 34 dB (Table 1). A subset of this dataset was also cre-

**Table 1. Data used in cougher & speaker identification:** TASK and Wallacecene datasets contain different noise level.

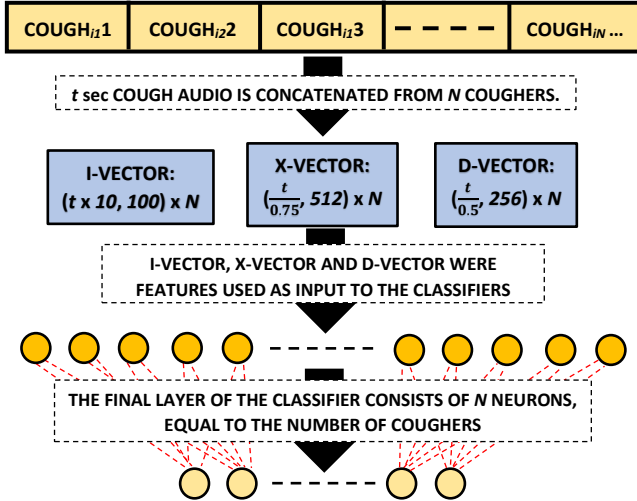
Dataset	Subjects	Events	Avg SNR	Avg Length
<i>Cougher identification</i>				
TASK	14	6000	73±5 dB	1.87±0.2 sec
Wallacedene	51	1358	34±17 dB	0.77±0.1 sec
<i>Speaker identification</i>				
Wallacedene	51	510	33±17 dB	0.99±0.2 sec

ated with only 42,341 events belonging to 10 classes, which can be used as commands in IoT or robotics [16]. For spotting cough as a trigger phrase, we note these two datasets as SC-36 and SC-11, containing 36 and 11 classes respectively.

### 3. FEATURE EXTRACTION

For cougher identification, We have extracted x-vectors and i-vectors using extractors pre-trained on the under-resourced languages [21], which are spoken by the subjects in the TASK and Wallacecene datasets (Figure 2).  $t$ -sec long audio from each of  $N$  coughers are concatenated by following the data preparation requirements of Kaldi ASR toolkit [20]. i-vectors are generated for each non-overlapping 0.1 sec audio from each utterance ID, with a dimension of  $(t \times 10, 100)$  for each cougher [11]. Unique x-vectors are generated for each 1.5 sec of utterance with 0.75 sec overlap, having a dimension of  $(1, 512)$  [12]. Thus for each  $t$  sec long audio from each cougher, there are x-vectors of dimension  $(\frac{t}{0.75}, 512)$ . We have also extracted d-vectors using extractor pre-trained on VCC 2018, VCTK, Librispeech, and CommonVoice English datasets and generalized using end-to-end loss function [13]. Every  $t$  sec cough is split into non-overlapping 0.5 sec audio clips, thus producing d-vectors of dimension  $(\frac{t}{0.5}, 256)$  for every cougher, suggesting that the i-vectors have a higher dimensionality than x-vectors and d-vectors. The number of subjects ( $N$ ) and the cough-time ( $t$ ) were the hyperparameters in cougher identification task (Table 3). For speakers, we used all counts, having only  $N$  as the hyperparameter. For TASK and Wallacecene datasets,  $N$  has been varied between 5 & 14 and 5 & 51 respectively with a step size of 5.

For spotting cough as a trigger phrase, we have extracted STFT, ZCR and kurtosis from overlapping frames ( $\mathcal{F}$ ) of the audio, where the frame overlap is computed to ensure that the audio signal is always divided into exactly  $\mathcal{S}$  frames, so that the entire audio event is always captured within a fixed number of frames, allowing a fixed input dimension to be maintained while preserving the general overall temporal structure of the event. Such fixed two-dimensional features are particularly useful for the training of DNN classifiers [8]. Table 2 shows that in our experiments each audio signal is divided into between 70 and 150 frames, each between 512 and 4096 samples i.e. 32 msec and 256 msec long, varying the spectral



**Fig. 2. Feature extraction for cougher identification:**  $t$  sec long coughs ( $\text{COUGH}_{i1}, \text{COUGH}_{i2}, \text{COUGH}_{i3}, \dots$ , where,  $1 \leq i \leq N$ ) from each cougher are concatenated as they appear in the audio recording for  $N$  coughers. i-vector, x-vector and d-vector are extracted from this  $t \times N$  sec long audio and they are fed to the classifiers with  $N$  neurons at the final layer to distinguish the cougher using cross-validation.

information of each event in SC-11 and SC-36 datasets.

**Table 2. Feature extraction hyperparameters.** Table 4 and 5 show classification results for these hyperparameters.

Hyperparameter	Description	Range
<i>Cougher identification</i>		
Subject ( $N$ )	no. of coughers or speakers	5 to 51 with step of 5
Cough-time ( $t$ )	cough from each subject	2, 5 to 100 with step of 5
<i>Cough spotting</i>		
Frame length ( $\mathcal{F}$ )	used to extract features	$2^k, k = 9, \dots, 12$
No. of frames ( $\mathcal{S}$ )	extracted from audio	$10 \times k, k = 7, 10, 12, 15$

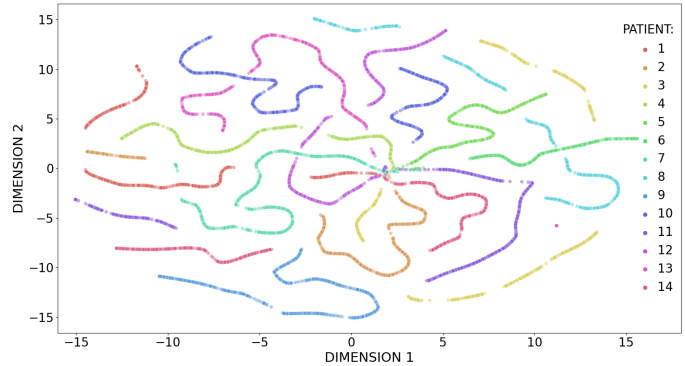
LR, LDA, SVM and MLP were used in identifying coughers and CNN, LSTM and Resnet50 were used in spotting coughs as a trigger phrase. Table 3 lists the hyperparameters considered for these classifiers and they were optimised using 5-fold cross-validation and the standard deviation among the outer folds is noted as  $\sigma_{ACC}$  in Table 4. For Resnet50, the 50-layer architecture in [22] has been used.

#### 4. RESULTS AND DISCUSSION

Table 4 shows the results using the best two features for both TASK (less-noisy) and Wallacedene (noisier) datasets. The highest accuracy (99.78%) has been achieved by an MLP when using i-vectors to identify coughers from 100 sec ( $t = 100$ ) long cough collected from each of 5 coughers. By increasing the number of coughers to 10 and 14, the performance of the MLP classifier decreased to 98.87% and 98.39% respectively for i-vectors (Table 4 and Figure 4).

**Table 3. Classifier hyperparameters** used in both identifying ‘coughers’ and ‘cough’ as a trigger phrase for KWS.

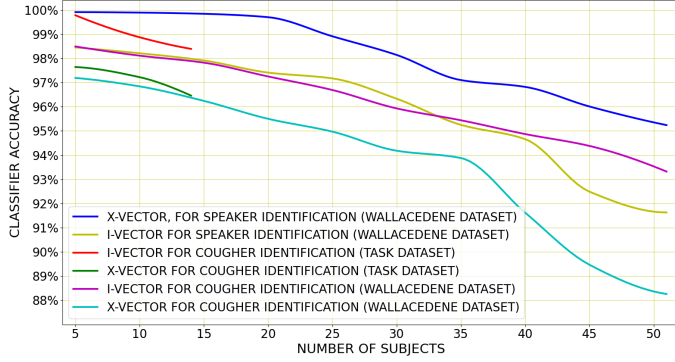
	Hyperparameters	Classifier	Range
coughers	Regularisation	LR & SVM	$10^i$ where $i = -7, \dots, 7$
	$l1$ penalty	LR	0 to 1 in steps of 0.05
	$l2$ penalty	LR, MLP	0 to 1 in steps of 0.05
	Kernel coeff.	SVM	$10^i$ where $i = -7, \dots, 7$
	No. of neurons	MLP	70 to 150 in steps of 20
cough	Batch size	CNN & LSTM	$2^k$ where $k = 6, 7, 8$
	No. of epochs	CNN & LSTM	10 to 200 in steps of 20
	No. of conv filters	CNN	$3 \times 2^k$ where $k = 3, 4, 5$
	kernel size	CNN	2 and 3
	Dropout rate	CNN & LSTM	0.1 to 0.5 in steps of 0.2
	Dense layer size	CNN & LSTM	$2^k$ where $k = 4, 5$
	LSTM units	LSTM	$2^k$ where $k = 6, 7, 8$
	Learning rate	LSTM	$10^k$ where $k = -2, -3, -4$



**Fig. 3. The t-SNE cluster of i-vectors extracted from 2-sec long cough audio from 14 coughers in TASK dataset.** The MLP produces 95.11% accuracy using these i-vectors in discriminating 14 coughers (Table 4).

All classifiers performed well in identifying both coughers and speakers on the noisier Wallacedene dataset. The speaker identification is used as the baseline and Table 4 shows that using x-vectors produced better classification scores than using i-vectors for speaker identification, also found by others [12]. The highest accuracy (99.91%) has been achieved from the MLP while discriminating only 5 speakers using x-vectors. This accuracy drops to 98.14% using MLP while differentiating 30 speakers and it drops further to 95.24% while discriminating all 51 speakers in Wallacedene dataset. For a lower number of coughers such as 5, MLP classifier has achieved the highest accuracy of 98.49% using i-vectors. The accuracy of MLP has dropped to 97.82%, 96.69%, 94.87% and 93.32% and the  $\sigma_{ACC}$  has increased sharply while the number of coughers is increased to 15, 25, 40 and 51 respectively. These scores show that although cougher identification is not as accurate as speaker identification, the performance is close, especially for the smaller number of subjects.

The results also show that cougher identification on less-noisy TASK dataset is more accurate than noisier Wallace-



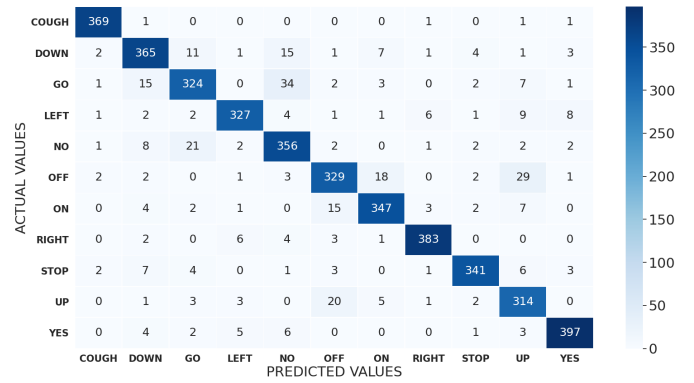
**Fig. 4. Classifier performance.** The accuracies from the MLP classifier decrease while discriminating more subjects.

dene dataset. Although longer coughs from each subject improve the classifier accuracy in general, similar performance is achieved (accuracy of 95.11% & 90.02% on less-noisy & noisy data) for coughs as short as only 2 sec (Figure 3). Although the performance is close, i-vectors performed better than x-vectors in cougher identification. MLP is the classifier of choice and it shows lower  $\sigma_{ACC}$  across the cross-validation folds for the less-noisy data than noisier data. d-vectors are outperformed by i-vectors and x-vectors for both the speech and cough, as also found by [23], thus excluded from Table 4.

**Table 4. Classifier accuracies in identifying coughers for both TASK and Wallacedene (WD) datasets.**

Dataset	$N$	$t$	Feature	LR	LDA	SVM	MLP	$\sigma_{ACC}$
TASK	5	100	i-vector	98.91%	98.87%	99.44%	99.78%	0.0007
		100	x-vector	96.71%	96.73%	97.54%	97.64%	0.0009
		80	i-vector	97.54%	97.88%	98.19%	98.87%	0.0006
		80	x-vector	96.31%	96.24%	96.55%	97.22%	0.0005
	10	2	i-vector	94.41%	94.51%	94.55%	<b>95.11%</b>	0.0005
		100	i-vector	96.46%	96.71%	97.48%	98.39%	0.0006
		100	x-vector	97.26%	97.54%	98.78%	96.46%	0.0008
		WD (Cougher)	5	20	i-vector	97.23%	97.19%	97.77%
20	x-vector			95.54%	95.97%	96.72%	97.19%	0.0078
15	20		i-vector	97.16%	97.14%	97.31%	97.82%	0.0061
	20		x-vector	95.41%	95.30%	95.72%	96.24%	0.0068
25	20		i-vector	95.04%	95.18%	95.94%	96.69%	0.0072
	20		x-vector	93.31%	93.55%	94.07%	94.97%	0.0082
40	20		i-vector	93.38%	93.62%	94.09%	94.87%	0.0091
	20		x-vector	90.23%	90.07%	90.97%	91.62%	0.0102
WD (Speaker)	5	2	i-vector	89.26%	89.38%	89.22%	<b>90.02%</b>	0.0178
		20	i-vector	90.27%	90.49%	91.89%	93.32%	0.0301
	20	x-vector	84.61%	84.74%	85.83%	88.26%	0.0247	
	30	—	x-vector	98.57%	98.64%	99.48%	99.91%	0.0018
		—	i-vector	97.21%	97.17%	97.70%	98.45%	0.0027
	51	—	x-vector	96.81%	96.85%	97.42%	98.14%	0.0081
—		i-vector	94.81%	94.87%	95.18%	96.33%	0.0078	
51	—	x-vector	99.44%	99.44%	99.44%	95.24%	0.0229	
	—	i-vector	90.01%	90.05%	90.34%	91.63%	0.0274	

Coughs were successfully spotted among other trigger phrases in both SC-11 and SC-36 datasets. Table 5 shows, although LSTM and CNN has performed well, the best performance of 92.73% accuracy ( $ACC$ ) & mean Cohen’s Kappa ( $\mathcal{K}$ ) of 0.9218 on SC-11 and 88.58% accuracy &  $\mathcal{K}$  of 0.8757 on SC-36 have been achieved from Resnet50. The confusion matrix using the best SC-11 system exhibits the high accuracies for spotting coughs in Figure 5. Table 5 also shows that the best results of CNN and Resnet50 were obtained mostly



**Fig. 5. The confusion matrix** of detecting coughs among 10 other trigger phrases in SC-11 dataset using the Resnet50.

from 1024 sample (64 msec) long frames and 100 segments.

**Table 5. Cough spotting:** The best-three results for each classifier shows Resnet50 has performed the best by achieving 92.73% & 88.58% accuracy on the SC-11 & SC-36 dataset.

Classifier	SC-11 Dataset				SC-36 Dataset			
	$\mathcal{F}$	$\mathcal{S}$	$ACC$	$\mathcal{K}$	$\mathcal{F}$	$\mathcal{S}$	$ACC$	$\mathcal{K}$
LSTM	512	150	88.09%	0.8767	512	120	80.74%	0.7937
	2048	120	87.66%	0.8614	1024	120	80.40%	0.7931
	512	70	87.09%	0.8598	512	100	80.11%	0.7902
CNN	1024	100	91.25%	0.9007	1024	120	86.74%	0.8592
	2048	100	90.72%	0.8981	1024	70	85.98%	0.8463
	1024	70	90.11%	0.8945	2048	100	85.22%	0.8411
Resnet50	1024	100	92.73%	0.9218	2048	100	88.58%	0.8777
	2048	120	92.69%	0.8733	2048	70	87.94%	0.8729
	2048	100	92.55%	0.8715	1024	120	87.68%	0.8702

## 5. CONCLUSION

We propose a system using cough as a wake-word by spotting coughs among other trigger phrases and identifying the cougher. A less-noisy and noisier dataset, containing 14 and 51 subjects respectively, was used to extract i-vector, x-vector and d-vector, to classify the cougher. The best performance was achieved from an MLP, showing coughers as many as 51 can be identified with 90.02% accuracy, using i-vectors from as short as 2 sec audio from each cougher in the noisy environment. We also found, unlike speakers, coughers were better identifiable using i-vectors. Coughs can also be spotted as wake-words using a Resnet50 on features keeping end-to-end time-domain information among 35 other keywords in Google Speech Commands dataset with 88.58% accuracy. Wake-cough represents a means of personalised, long-term cough monitoring that is non-intrusive and, due to the use of wake-word detection methods, power-efficient since a smartphone-based monitoring device can remain mostly dormant. Thus, it represents an attractive and viable means for monitoring a patient’s long-term recovery from lung ailments such as TB and COVID-19. Its ability to discriminate between coughers also makes it attractive in multi-bed ward environments in monitoring patient’s recovery process.

## 6. REFERENCES

- [1] Minhua Wu, Sankaran Panchapagesan, Ming Sun, Jiacheng Gu, Ryan Thomas, Shiv Naga Prasad Vitaladevuni, Bjorn Hoffmeister, and Arindam Mandal, "Monophone-Based Background Modeling for Two-Stage On-Device Wake Word Detection," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5494–5498.
- [2] Tara Sainath and Carolina Parada, "Convolutional Neural Networks for Small-footprint Keyword Spotting," in *INTER-SPEECH*, 2015, pp. 1478–1482.
- [3] Yixin Gao, Yuriy Mishchenko, Anish Shah, Spyros Matsoukas, and Shiv Vitaladevuni, "Towards Data-Efficient Modeling for Wake Word Spotting," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7479–7483.
- [4] Guoguo Chen, Carolina Parada, and Georg Heigold, "Small-footprint keyword spotting using deep neural networks," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 4087–4091.
- [5] Renier Botha, Grant Theron, Robbin Warren, Marisa Klopper, Keertan Dheda, Paul Van Helden, and Thomas Niesler, "Detection of tuberculosis by automatic cough sound analysis," *Physiological Measurement*, vol. 39, no. 4, pp. 045005, 2018.
- [6] Mahmood Al-khassaweneh and Ra'ed Bani Abdelrahman, "A signal processing approach for the diagnosis of asthma from cough sounds," *Journal of Medical Engineering & Technology*, vol. 37, no. 3, pp. 165–171, 2013.
- [7] Renard Xaviero Adhi Pramono, Syed Anas Intiaz, and Esther Rodriguez-Villegas, "A cough-based algorithm for automatic diagnosis of pertussis," *PLOS ONE*, vol. 11, no. 9, pp. e0162128, 2016.
- [8] Madhurananda Pahar, Marisa Klopper, Robin Warren, and Thomas Niesler, "COVID-19 cough classification using machine learning and global smartphone recordings," *Computers in Biology and Medicine*, vol. 135, pp. 104572, 2021.
- [9] Fengpei Ge and Yonghong Yan, "Deep neural network based wake-up-word speech recognition with two-stage detection," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 2761–2765.
- [10] Veton Z Këpuska and TB Klein, "A novel Wake-Up-Word speech recognition system, Wake-Up-Word recognition task, technology and evaluation," *Nonlinear Analysis: Theory, Methods & Applications*, vol. 71, no. 12, pp. e2772–e2789, 2009.
- [11] Andrew Senior and Ignacio Lopez-Moreno, "Improving DNN speaker independence with i-vector inputs," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 225–229.
- [12] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur, "X-Vectors: Robust DNN Embeddings for Speaker Recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5329–5333.
- [13] Li Wan, Quan Wang, Alan Papir, and Ignacio Lopez Moreno, "Generalized End-to-End Loss for Speaker Verification," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4879–4883.
- [14] Matt Whitehill, Jake Garrison, and Shwetak Patel, "Whosecough: In-the-Wild Cougher Verification Using Multitask Learning," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 896–900.
- [15] Miao Zhang, Yixiang Chen, Lantian Li, and Dong Wang, "Speaker recognition with cough, laugh and "Wei"," in *2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2017, pp. 497–501.
- [16] Pete Warden, "Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition," *arXiv preprint arXiv:1804.03209*, 2018.
- [17] Peter Wittenburg, Hennie Brugman, Albert Russel, Alex Klassmann, and Han Sloetjes, "ELAN: a professional framework for multimodality research," in *5th International Conference on Language Resources and Evaluation (LREC 2006)*, 2006.
- [18] Madhurananda Pahar, Igor Miranda, Andreas Diacon, and Thomas Niesler, "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)*, 2021, pp. 8002–8006.
- [19] Madhurananda Pahar, Marisa Klopper, Byron Reeve, Robin Warren, Grant Theron, and Thomas Niesler, "Automatic Cough Classification for Tuberculosis Screening in a Real-World Environment," *arXiv preprint arXiv:2103.13300*, 2021.
- [20] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, "The Kaldi Speech Recognition Toolkit," in *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*. IEEE Signal Processing Society, 2011, IEEE Catalog No.: CFP11SRW-USB.
- [21] Trideba Padhi, Astik Biswas, Febe de Wet, Ewald van der Westhuizen, and Thomas Niesler, "Multilingual bottleneck features for improving ASR performance of code-switched speech in under-resourced languages," in *Proceedings of the First Workshop on Speech Technologies for Code-switching in Multilingual Communities (WSTCSMC)*, Shanghai, China, 2020.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [23] Lantian Li, Yiye Lin, Zhiyong Zhang, and Dong Wang, "Improved deep speaker feature learning for text-dependent speaker recognition," in *2015 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*. IEEE, 2015, pp. 426–429.