

HeAR - Health Acoustic Representations

Sebastien Baur^{1*}, Zaid Nabulsi^{1*}, Wei-Hung Weng¹, Jake Garrison¹, Louis Blankemeier¹, Sam Fishman¹, Christina Chen¹, Sujay Kakarmath¹, Minyoi Maimbolwa², Nsala Sanjase², Brian Shuma², Yossi Matias¹, Greg S. Corrado¹, Shwetak Patel¹, Shravya Shetty¹, Shruthi Prabhakara¹, Monde Muyoyeta², Diego Ardila¹

¹Google Research, USA

²TB department, Center of Infectious Disease Research in Zambia, Zambia

HEALTH_ACOUSTIC_REPRESENTATIONS@GOOGLE.COM

Abstract

Health acoustic sounds such as coughs and breaths are known to contain useful health signals with significant potential for monitoring health and disease, yet are underexplored in the medical machine learning community. The existing deep learning systems for health acoustics are often narrowly trained and evaluated on a single task, which is limited by data and may hinder generalization to other tasks. To mitigate these gaps, we develop HeAR, a scalable self-supervised learning-based deep learning system using masked autoencoders trained on a large dataset of 313 million two-second long audio clips. Through linear probes, we establish HeAR as a state-of-the-art health audio embedding model on a benchmark of 33 health acoustic tasks across 6 datasets. By introducing this work, we hope to enable and accelerate further health acoustics research.

Keywords: health acoustics, respiratory sounds, cough, acoustic vitals, audio sensing

1. Introduction

Acoustic non-semantic attributes of speech can enable machine learning models to perform paralinguistic tasks, including emotion recognition, speaker identification, and dementia detection (Shor et al., 2022). Cerebrovascular and neurodegenerative diseases like stroke, Parkinson’s, Alzheimer’s, cerebral palsy and amyotrophic lateral sclerosis (ALS) may also be detected and monitored using non-semantic patterns of speech, such as articulation, resonance, and phonation (Boschi et al., 2017). Non-semantic acoustic signals related to health are not confined solely to conversational speech data. Health-related acoustic cues, originating from the respiratory system’s airflow, including sounds like coughs and breathing patterns can be harnessed for health monitoring purposes. For example, clinicians use sounds such as “whoop”-like coughing to diagnose pertussis (Pramono et al., 2016), and agonal breathing for detecting acute cardiovascular events. Such health sounds can also be collected via ambient sensing technologies on ubiquitous devices such as mobile phones (Zimmer et al., 2022), which may augment healthcare workers in low-medium income countries (LMICs) with improved screening capabilities.

With advancements in deep learning, neural networks are now able to learn high-quality general representations directly from raw speech data (Zhang et al., 2022), and use them for various semantic and non-semantic speech-related tasks (Peplinski et al., 2020; Shor et al., 2022; Shor and Venugopalan, 2022). Health acoustics, specifically non-semantic respiratory sounds, also have potential as biomarkers to detect various respiratory diseases (Alqudaihi et al., 2021). However, current machine learning (ML) systems for health acoustics are task-specific and may not generalize well to out-of-distribution (OOD) settings (D’Amour et al., 2022) and are often limited by data quantity.

Recently, self-supervised learning (SSL) has demonstrated potential for building robust and capable systems by learning general representations from large, unlabeled sources (Balestriero et al., 2023; Chen et al., 2020; He et al., 2022). There is extensive progress on learning general, universal representations in vision (Dosovitskiy et al., 2020), language (Chowdhery et al., 2022) and speech (Zhang et al., 2023). Such

* These authors contributed equally

approaches are also used for learning representations of biomedical language, medical images and even physiological waveforms (Belyaeva et al., 2023; Singhal et al., 2023; Xu et al., 2023). However, these approaches are still underexplored in the field of health acoustics.

To demonstrate the potential of SSL on the underexplored health acoustic modality, we introduce **HeAR: Health Acoustic Representations**, a self-supervised generative learning-based system trained on a large dataset of two-second long audio clips for learning low-dimensional representations that can transfer well across health acoustic tasks and generalize to OOD data. We benchmark HeAR on a diverse set of health acoustic tasks spanning 13 health acoustic event detection tasks, 14 cough inference tasks, and 6 spirometry inference tasks, across 6 datasets and demonstrate that simple linear classifiers trained on top of our representations outperform the state-of-the-art on many tasks.

2. Related Works

SSL has emerged as a critical ML paradigm to learn general representations from large unannotated datasets. SSL can have numerous training objectives including contrastive, such as SimCLR (Chen et al., 2020), BYOL (Grill et al., 2020), and generative, like masked autoencoder (MAE) (He et al., 2022). In recent years, data-driven audio SSL has made great progress, specifically for semantic speech. From CPC (Oord et al., 2018), Wav2vec 2.0 (Baevski et al., 2020), BigSSL (Zhang et al., 2022), AudioMAE (Huang et al., 2022), to BEST-RQ (Chiu et al., 2022) and Universal Speech Model (USM) (Zhang et al., 2023), researchers utilize massive unlabeled data from the Internet to train SSL-based audio encoders to learn better speech representations. There are also studies focusing on non-semantic speech. For example, TRILL (Shor et al., 2020) uses a triplet loss as the training objective, TRILLsson (Shor and Venugopalan, 2022) and FRILL (Peplinski et al., 2020) further distill the TRILL encoder to make it smaller and faster. Researchers also adopted different neural network architectures, such as Conformer (Shor et al., 2022; Srivastava et al., 2022), and Slowfast NFNet (Kazakos et al., 2021; Wang et al., 2022) to develop performant audio encoders. In our work, we adopt a generative SSL framework (MAE) but focus on non-semantic health acoustics, which is relatively underexplored yet useful for healthcare applications.

ML for non-semantic health acoustics is also an emerging research area for health. There are a growing number of studies using respiratory sounds for health monitoring and disease detection. For example, cough sound patterns can be used as a biomarker to identify coughers (Whitehill et al., 2020), detect various respiratory diseases, such as COVID-19 (Coppock et al., 2021; Laguarta et al., 2020; Schuller et al., 2020), bronchitis, bronchiolitis, pertussis (Bales et al., 2020), obstructive versus restrictive lung diseases (Rudraraju et al., 2020), and tuberculosis (TB) (Tracey et al., 2011; Larson et al., 2012; Botha et al., 2018; Rudraraju et al., 2020; Pahar et al., 2021; Zimmer et al., 2022; Sharma et al., 2024). However, these works focus on developing a single task-specific system trained in a supervised learning framework that may not generalize as well to new settings. Our work instead introduces a system trained without supervision on a large and diverse unlabelled corpus, which may generalize better to unseen distributions and new tasks (Radford et al., 2021).

3. Methods

HeAR comprises three main components: a data curation step (including a health acoustic event detector), a general purpose training step to develop an audio encoder (embedding model), and a task-specific evaluation step that adopts the trained embedding model for various downstream tasks. The system is designed to encode two-second long audio clips and generate audio embeddings for use in downstream tasks. Figure 1 illustrates these high-level components of the system.

The health acoustic event detector is a multilabel classification convolutional neural network (CNN) that identifies the presence of any of six types of non-speech health acoustic events in two-second audio clips : coughing, baby coughing, breathing, throat clearing, laughing, and speaking. The detector is described in Appendix A.

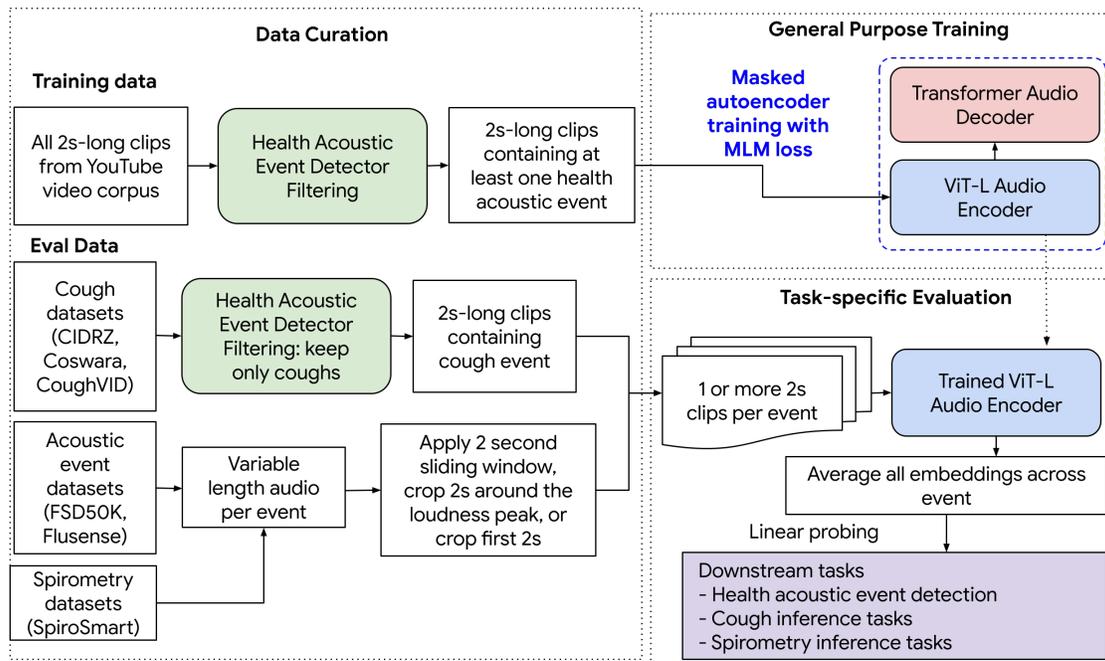


Figure 1: HeAR system overview.

Similarly to AudioMAE (Huang et al., 2022), we used a MAE (He et al., 2022) to learn audio representations by training an autoencoder to reconstruct masked 16x16 spectrogram patches (units: time and frequency). 75% of input patches are masked out and encoded by ViT-L (Dosovitskiy et al., 2020). Learnable mask tokens are added to the sequence of encoded tokens, and an 8-layer transformer decoder is tasked with reconstructing the missing patches, by minimizing the L_2 distance between normalized masked patches and its predictions. Note that we did not experiment with self-attention as done by Huang et al. (2022). The model is trained using the AdamW optimizer for 950k steps (~ 4 epochs) with global batch size 4096, and hyperparameters from Huang et al. (2022). Learning rate follows a cosine decay schedule, starting at $4.8e-4$, following the commonly used linear batch scaling rule (Goyal et al., 2017).

Datasets For training, we curate a dataset, YT-NS (YouTube Non-Semantic), consisting of two-second long audio clips extracted from three billions public non-copyrighted YouTube videos using the health acoustic event detector (described in Appendix A), totalling 313.3 million two-second clips or roughly 174k hours of audio. We chose a two-second window since most events we cared about were shorter than that. The HeAR audio encoder is trained solely on this dataset.

We benchmark HeAR both on general health acoustic event classification and on cough inference tasks. For general health acoustic event classification, we use FSD50K (Fonseca et al., 2021) and FluSense (Al Hossain et al., 2020). FluSense contains human-annotated timestamped labels, which we use to extract short labeled audio clips (50% are shorter than 2 seconds and 90% shorter than 7 seconds). FSD50K does not have such timestamps, and some of its clips can be far longer than the duration of clips seen during training. For this reason, we crop two-second audio clips around the loudness peak, computed as the average sound amplitude in dB. We found that this resulted in higher performance for all models on all detection tasks.

For cough inference tasks, we use CoughVID (Orlandic et al., 2021), Coswara (Bhattacharya et al., 2023), and a prospective tuberculosis-specific dataset from the Centre for Infectious Disease Research in Zambia (CIDRZ) where cough audio recordings and chest X-rays were obtained from a cohort of symptomatic patients. Considering that the recording quality may be affected by various environmental factors, we obtain CIDRZ audio from microphones of varying quality under a study data collection protocol that controls the

environmental factors. Three devices were used to collect the cough audio, representing three mobile phone “tiers” (for the purposes of this study) with different costs and potentially varying recording quality: Pixel3a, GalaxyA12, and GalaxyA22. We refer to the datasets as CIDRZ low-tier, CIDRZ mid-tier, and CIDRZ high-tier respectively. We use recordings from all devices to train the linear probes since they provide more data, and focus evaluation results in the main text on one device (low-tier) for brevity. Results for evaluation on mid-tier and high-tier datasets are available in Appendix C. We also include evaluation on SpiroSmart, one pulmonary testing dataset described in Garrison (2018), which includes spirometry efforts paired with audio recordings of forced expiratory efforts from patients of chronic obstructive pulmonary disease (COPD) clinics around the world. We describe the evaluation datasets, including train/validation/test split sizes, in Appendix Table B1, and summary statistics of the unpublished CIDRZ dataset are listed in Appendix B and Appendix Table B2.

Baseline Models We consider several state-of-the-art audio encoder baselines for comparison: (1) TRILL (Shor et al., 2020), a publicly available ResNet50-based encoder trained on an AudioSet subset with speech labels by optimizing triplet loss, (2) FRILL (Peplinski et al., 2020), a publicly available MobileNet-based distilled version of TRILL for mobile devices, (3) BigSSL-CAP12 (Shor et al., 2022), a Conformer-based encoder trained on YouTube 900k-hour speech and LibriLight with a wav2vec 2.0 objective, and (4) CLAP, a CNN-based audio encoder trained using multimodal contrastive learning on a mixture of datasets that include FSD50K (Elizalde et al., 2023). We also investigated different Mel spectrogram-based models using spectrograms or Mel-frequency cepstral coefficients (MFCCs) as features, as well as a randomly initialized MAE. We didn’t include them in the final evaluation since their performance was very low and they did not contribute meaningfully to our benchmark.

Evaluation on Downstream Tasks To evaluate the quality of representations learned by HeAR and compare it to other encoders’, we train linear probes. More specifically, we encode all recordings from all datasets using TRILL, FRILL, BigSSL-12, HeAR, and CLAP, and we train separate linear or logistic regressions (with a cross-validated ridge penalty) to predict available labels on these datasets (Köhn, 2015). When this information is available, the cross-validation procedure groups recordings of the same individual, the same audio clip, and the same site within the same fold, to stick as close as possible to the out-of-sample evaluation scenario. Once the regularization coefficient is chosen, we evaluate the performance of that linear model on a held-out validation dataset for which all audio recordings come from sources not seen for training HeAR or cross-validating the linear probes. The validation datasets were used to experiment with various data preprocessing schemes, especially for datasets that had clips longer than YT-NS (up to 30x longer). Finally, once the best preprocessing scheme has been identified on the validation datasets, we compute the performance of the best linear models on our held-out test dataset, which is also disjoint from the train and validation datasets. We set up 13 health acoustic event detection tasks from two datasets, 14 cough inference tasks from three datasets, and 6 spirometry tasks from one dataset. For FluSense and FSD50K, we train a separate linear probe for each task, predicting whether the specific audio event occurs in the clip. For the cough inference tasks, we train a separate linear probe for each task and for each dataset, predicting the specific label from a two-second audio recording of a cough. Cough inference tasks include identifying three types of chest X-ray (CXR) findings (unspecified abnormalities, presence of focal or multifocal lung opacities, and pleural effusion), two diagnostic tasks (COVID on two datasets and tuberculosis on another one), and identifying demographics and lifestyle factors (smoking status, sex, age, BMI). Spirometry tasks include estimation of forced expiratory volume (FEV1), forced vital capacity (FVC), the FEV1/FVC ratio, peak flow (PEF), total exhale duration (FET), and sex classification. We report either the area under the receiver operating characteristic curve (AUROC) or average precision (AP) for one-versus-rest classification tasks. We used the DeLong method to compute the 95% confidence intervals (CIs) of AUROC (DeLong et al., 1988). For regression tasks, we report mean absolute error with bootstrapped 95% confidence intervals.

4. Results

Across a range of 33 tasks on 6 datasets, HeAR achieved the highest performance among all models, as measured by mean reciprocal rank (0.708, see Table 1), reaching the highest rank on 17 tasks (3 out of 13 health acoustic event detection tasks in Table 2 and Figure 2, 10 out of 14 cough inference tasks in Table 3 and Figure 3, and 5 out of 6 spirometry tasks in Table 4 and Figure 4).

Task group	TRILL	FRILL	BigSSL-CAP12	HeAR	CLAP (48k)
All	0.322	0.273	0.419	0.708	0.555
Health acoustic detection	0.225	0.235	0.438	0.538	0.846
Cough	0.423	0.305	0.373	0.812	0.370
Spirometry	0.298	0.281	0.486	0.833	0.356

Table 1: Mean reciprocal ranks on groups of tasks.

For the health acoustic detection tasks, CLAP performs best overall (mAP=0.691, MRR=0.846), which might be expected since FSD50K was used in its training procedure. HeAR has the second highest performance (mAP=0.658, MRR=0.538) and highest among models that haven’t used FSD50K for training. We observed on FSD50K that the performance of HeAR degraded significantly with sequence length, which we hypothesize to be due to the use of fixed sinusoidal positional encodings, which are known (Kazemnejad et al., 2024) to generalize poorly to unseen longer sequence lengths. Cropping the loudest two-second clip improved the performance of all models, and especially HeAR.

Dataset	Binary Classification Task	Summary statistics for test split (label)	Metric	TRILL	FRILL	BigSSL-CAP12	HeAR	CLAP (48k)
FSD50K + FluSense	All	N/A	mAP	0.494	0.516	0.613	0.658	0.691
FSD50K	Breathing	227 / 10004 (Y/N)	AP	0.301 [0.242, 0.365]	0.336 [0.276, 0.399]	0.365 [0.294, 0.434]	0.434 [0.365, 0.496]	0.467 [0.394, 0.538]
	Cough	106 / 10125 (Y/N)	AP	0.450 [0.356, 0.547]	0.452 [0.359, 0.543]	0.658 [0.568, 0.742]	0.621 [0.513, 0.719]	0.751 [0.673, 0.821]
	Laughter	253 / 9978 (Y/N)	AP	0.438 [0.379, 0.495]	0.425 [0.365, 0.483]	0.673 [0.622, 0.726]	0.680 [0.624, 0.732]	0.715 [0.664, 0.762]
	Respiratory sounds	380 / 9851 (Y/N)	AP	0.539 [0.489, 0.587]	0.535 [0.489, 0.580]	0.629 [0.583, 0.675]	0.670 [0.624, 0.716]	0.702 [0.652, 0.749]
	Sneeze	61 / 10170 (Y/N)	AP	0.361 [0.260, 0.471]	0.448 [0.340, 0.559]	0.570 [0.445, 0.685]	0.650 [0.537, 0.746]	0.912 [0.843, 0.964]
	Speech	785 / 9446 (Y/N)	AP	0.430 [0.397, 0.466]	0.418 [0.384, 0.452]	0.567 [0.533, 0.603]	0.534 [0.498, 0.572]	0.599 [0.568, 0.629]
FluSense	Breathing	35 / 1624 (Y/N)	AP	0.147 [0.088, 0.246]	0.233 [0.135, 0.359]	0.357 [0.238, 0.506]	0.336 [0.236, 0.464]	0.371 [0.238, 0.539]
	Cough	430 / 721 (Y/N)	AP	0.903 [0.881, 0.922]	0.892 [0.870, 0.912]	0.954 [0.941, 0.965]	0.974 [0.966, 0.982]	0.963 [0.949, 0.974]
	Gasp	106 / 1780 (Y/N)	AP	0.466 [0.384, 0.570]	0.587 [0.499, 0.694]	0.653 [0.568, 0.734]	0.608 [0.518, 0.695]	0.701 [0.606, 0.789]
	Sneeze	160 / 1367 (Y/N)	AP	0.648 [0.579, 0.714]	0.661 [0.589, 0.727]	0.810 [0.753, 0.860]	0.788 [0.720, 0.848]	0.825 [0.770, 0.877]
	Sniffle	158 / 1413 (Y/N)	AP	0.718 [0.654, 0.778]	0.667 [0.599, 0.736]	0.720 [0.648, 0.792]	0.852 [0.799, 0.893]	0.841 [0.783, 0.889]
	Speech	483 / 593 (Y/N)	AP	0.949 [0.937, 0.961]	0.949 [0.937, 0.960]	0.983 [0.976, 0.989]	0.972 [0.962, 0.981]	0.973 [0.962, 0.983]
	Throat-Clearing	18 / 2057 (Y/N)	AP	0.070 [0.022, 0.171]	0.099 [0.029, 0.255]	0.035 [0.026, 0.046]	0.436 [0.243, 0.644]	0.169 [0.090, 0.314]

Table 2: Performance comparison on health acoustic event detection tasks on FSD50K and FluSense datasets. Due to high class imbalance, and following conventional reporting on those datasets, we report average precision (AP), together with bootstrapped 95% confidence intervals. We also include mean average precision on the top row. Note that “Respiratory sounds” in FSD50K is the respiratory sound class excluding the other five FSD50K health acoustic event classes in the table.

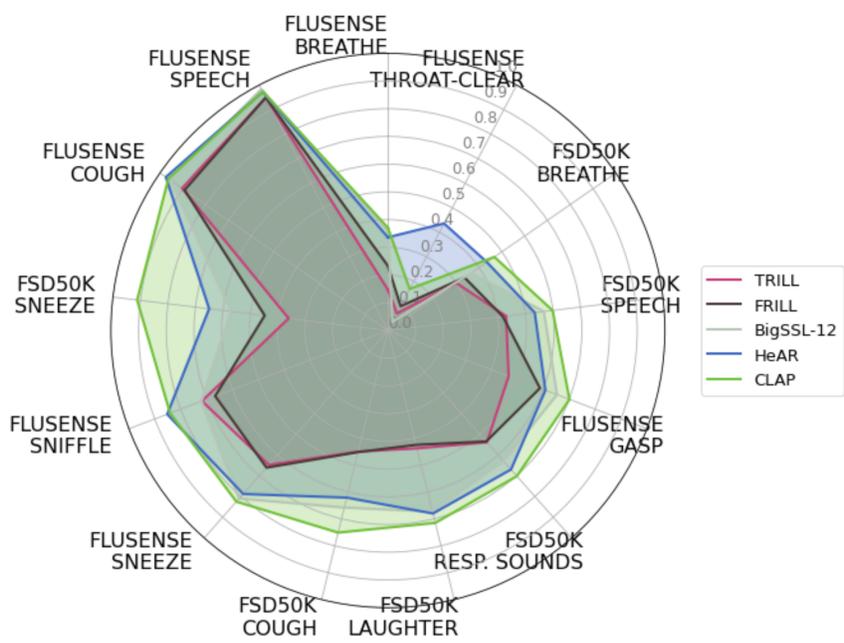


Figure 2: Radar plot of the performance comparison on health acoustic event detection tasks on FSD50K and FluSense datasets.

When evaluated on cough inference tasks, HeAR performed better than the baselines across 10/14 tasks, including demographics, lifestyle, and COVID tasks. On TB and CXR tasks, its performance is comparable to the best performing model. A summary of results is listed in Table 3. Appendix Table C1 compares the performance of all models across recording devices on CIDRZ. In particular, the performance of HeAR on the CXR most balanced tasks (lung opacities and unspecified abnormalities) exhibits the lowest variation (at most 1% AUROC difference between the best and worst recording devices) and it reaches the highest performance for mid-tier and high-tier datasets, while TRILL and FRILL, while scoring highest AUROC on the low-tier dataset (Table 3), also exhibit the highest variation, with up to 12% drop between best and worst recording devices (Appendix Table C1).

Dataset	Task	Summary statistics for test split (label)	Metric	TRILL	FRILL	BigSSL-CAP12	HeAR	CLAP
CIDRZ (Pixel3a)	Focal / multi focal lung opacities	61 / 204 (Y/N)	AUROC [DeLong 95% CI]	0.809 [0.747, 0.870]	0.800 [0.740, 0.860]	0.747 [0.672, 0.821]	0.794 [0.728, 0.861]	0.760 [0.690, 0.830]
CIDRZ (Pixel3a)	Abnormal CXR	64 / 201 (Y/N)	AUROC [DeLong 95% CI]	0.815 [0.757, 0.874]	0.778 [0.712, 0.844]	0.739 [0.664, 0.814]	0.763 [0.695, 0.830]	0.734 [0.658, 0.810]
CIDRZ (Pixel3a)	Pleural effusion	20 / 244 (Y/N)	AUROC [DeLong 95% CI]	0.683 [0.553, 0.812]	0.688 [0.562, 0.813]	0.684 [0.548, 0.819]	0.610 [0.465, 0.755]	0.748 [0.629, 0.866]
CIDRZ (Pixel3a)	Tuberculosis	24 / 240 (Y/N)	AUROC [DeLong 95% CI]	0.652 [0.520, 0.784]	0.648 [0.523, 0.772]	0.659 [0.533, 0.786]	0.739 [0.636, 0.841]	0.740 [0.627, 0.853]
CIDRZ (Pixel3a)	Sex	151 / 114 (F/M)	AUROC [DeLong 95% CI]	0.933 [0.901, 0.965]	0.928 [0.894, 0.961]	0.936 [0.909, 0.964]	0.974 [0.958, 0.990]	0.907 [0.872, 0.942]
CoughVID	Sex	1031 / 1924 (F/M)	AUROC [DeLong 95% CI]	0.850 [0.835, 0.866]	0.848 [0.832, 0.863]	0.872 [0.858, 0.887]	0.897 [0.884, 0.910]	0.821 [0.805, 0.838]
Coswara	Sex	174 / 478 (F/M)	AUROC [DeLong 95% CI]	0.920 [0.894, 0.947]	0.917 [0.891, 0.942]	0.937 [0.917, 0.956]	0.979 [0.965, 0.993]	0.892 [0.862, 0.923]
CIDRZ (Pixel3a)	Smoking status	65 / 198 (ever / never)	AUROC [DeLong 95% CI]	0.822 [0.762, 0.883]	0.811 [0.747, 0.874]	0.840 [0.786, 0.895]	0.877 [0.833, 0.921]	0.808 [0.750, 0.866]
Coswara	Smoking status	43 / 58 (current / never)	AUROC [DeLong 95% CI]	0.627 [0.518, 0.735]	0.619 [0.509, 0.729]	0.587 [0.476, 0.698]	0.631 [0.523, 0.739]	0.619 [0.509, 0.729]
CoughVID	COVID status	172 / 2237 (Y/N)	AUROC [DeLong 95% CI]	0.636 [0.565, 0.708]	0.634 [0.562, 0.706]	0.663 [0.596, 0.730]	0.710 [0.647, 0.774]	0.618 [0.540, 0.696]
Coswara	COVID status	63 / 470 (Y/N)	AUROC [DeLong 95% CI]	0.622 [0.581, 0.663]	0.615 [0.574, 0.657]	0.611 [0.570, 0.652]	0.645 [0.603, 0.687]	0.624 [0.584, 0.665]
CIDRZ (Pixel3a)	BMI	$\mu=22.9 \text{ kg/m}^2$ $\sigma=6.01 \text{ kg/m}^2$	Mean absolute error [95% bootstrapped CI]	3.861 [3.362, 4.458]	3.860 [3.359, 4.465]	3.875 [3.366, 4.452]	3.818 [3.328, 4.397]	3.836 [3.337, 4.433]
CIDRZ (Pixel3a)	Age	$\mu=35.6 \text{ yr}$ $\sigma=13.1 \text{ yr}$	Mean absolute error [95% bootstrapped CI]	10.590 [9.733, 11.509]	10.959 [10.125, 11.855]	10.009 [9.157, 10.944]	9.316 [8.550, 10.123]	10.775 [9.934, 11.644]
Coswara	Age	$\mu=32.7 \text{ yr}$ $\sigma=12.3 \text{ yr}$	Mean absolute error [95% bootstrapped CI]	9.994 [9.479, 10.570]	10.013 [9.484, 10.590]	9.665 [9.156, 10.247]	8.742 [8.269, 9.277]	10.133 [9.592, 10.743]

Table 3: Performance comparison on cough inference tasks. We report AUROC (with DeLong 95% confidence intervals) for binary classification tasks and mean absolute error (with bootstrapped confidence intervals) for regression tasks (age and BMI).

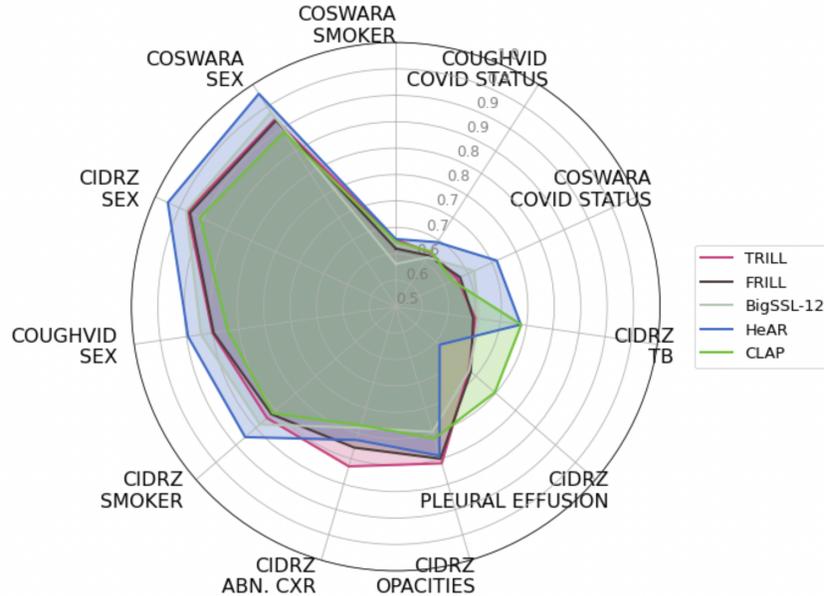


Figure 3: Radar plot of the performance comparison on cough inference tasks.

On SpiroSmart, HeAR performed better than our baselines on 4/5 lung function tasks and on sex classification. A summary of results is listed in Table 4.

Task	Summary statistics for test split	Metric	TRILL	FRILL	BigSSL-CAP12	HeAR	CLAP (48k)
FEV1	108 patients $\mu=1.71L$ $\sigma=0.914L$	Mean absolute error [95% bootstrapped CI]	0.488 [0.408, 0.568]	0.481 [0.404, 0.565]	0.479 [0.404, 0.558]	0.418 [0.351, 0.491]	0.518 [0.442, 0.604]
FVC	108 patients $\mu=2.29L$ $\sigma=0.915L$	Mean absolute error [95% bootstrapped CI]	0.559 [0.482, 0.641]	0.548 [0.471, 0.631]	0.536 [0.457, 0.614]	0.476 [0.409, 0.547]	0.561 [0.483, 0.641]
FEV1/FVC	108 patients $\mu=0.717$ $\sigma=0.167$	Mean absolute error [95% bootstrapped CI]	0.087 [0.073, 0.102]	0.090 [0.075, 0.105]	0.083 [0.069, 0.098]	0.083 [0.070, 0.097]	0.086 [0.073, 0.103]
PEF	108 patients $\mu=4.78 L/s$ $\sigma=2.32 L/s$	Mean absolute error [95% bootstrapped CI]	1.199 [1.022, 1.375]	1.318 [1.117, 1.498]	1.319 [1.130, 1.506]	1.147 [0.956, 1.343]	1.388 [1.192, 1.581]
FET	108 patients $\mu=6.68s$ $\sigma=2.47s$	Mean absolute error [95% bootstrapped CI]	1.540 [1.254, 1.937]	1.541 [1.251, 1.932]	1.452 [1.171, 1.819]	1.508 [1.212, 1.878]	1.371 [1.101, 1.703]
Sex (male/female)	108 patients 50 / 58 (M/F)	AUROC [DeLong 95% CI]	0.914 [0.858, 0.970]	0.914 [0.858, 0.970]	0.931 [0.882, 0.980]	0.934 [0.891, 0.978]	0.878 [0.809, 0.947]

Table 4: Performance comparison on spirometry tasks. We report mean absolute error for regression tasks (with bootstrapped 95% confidence intervals), and AUROC (with DeLong 95% confidence intervals) for binary classification tasks (sex).

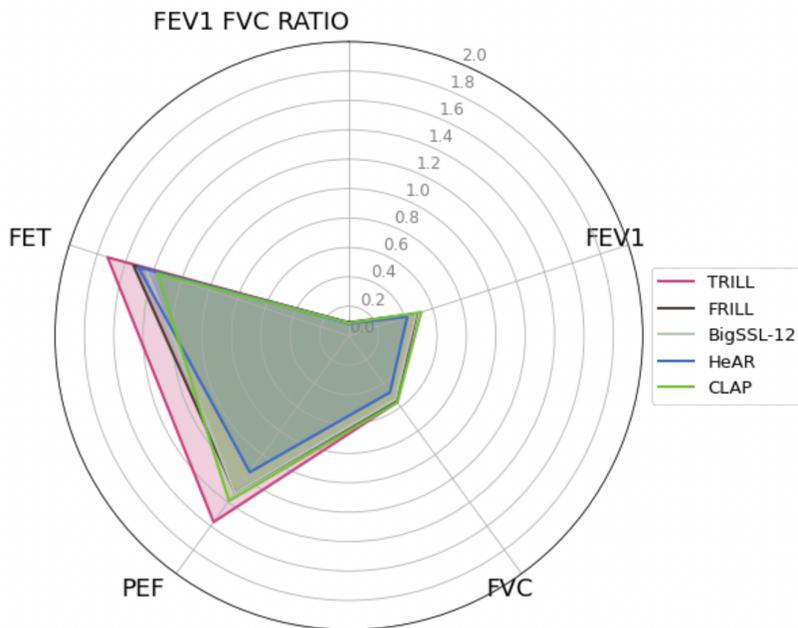


Figure 4: Radar plot of the performance comparison on spirometry tasks.

5. Discussion

In this work, we develop and evaluate the HeAR system that integrates a health acoustic event detector and a generative learning-based audio encoder (MAE) to learn health acoustic representations. The audio encoder is trained on YT-NS, without the requirement of human or expert data curation. We then demonstrate the quality of acoustic representations learned from that system via health acoustic event detection tasks, cough inference tasks, and spirometry inference tasks. In particular, advances in classification of tuberculosis from cough sounds could help risk-stratify patients needing X-ray screening or further testing in environments where chest X-rays are scarce or unavailable. The potential to monitor lung function from smartphone audio recordings could help develop easy-to-use and ubiquitous COPD screening tools and doctors monitor more closely their patients’ lung function evolution.

Self-supervised learning has achieved significant success in various domains, leveraging vast unlabeled data to train robust and generalizable encoders (Radford et al., 2021; Zhang et al., 2022, 2023; Devlin et al., 2018; Raffel et al., 2020; Yu et al., 2022). These encoders learn representations that have demonstrated increased robustness to distribution shift, improved transferability, and greater data efficiency (Radford et al., 2021). To the best of our knowledge, our work presents the first application of such large-scale self-supervised learning to health acoustic tasks. Our experiments reveal that increased pretraining data enhances downstream performance across diverse tasks (see Appendix D), consistent with prior findings (Radford et al., 2021). Notably, HeAR consistently achieves superior performance on diverse health-relevant tasks (inference of medical conditions and medically-relevant quantities from recordings of coughs or exhalations), as shown in Tables 1, 3, and 4. Since it is likely that users would use microphones not seen during training, robustness to recording devices is valuable. We show in Appendix C that HeAR’s performance on CIDRZ tasks remains stable across recording devices. In addition, when training linear probes on two devices and evaluating the probes on the third device (i.e., held-out device), HeAR performs better than other models (MRR=0.745 vs. second-best being CLAP with MRR=0.497), indicating potential for real-world applications. We hypothesize that the scale and diversity of recording devices in YT-NS contribute to making HeAR embeddings more device-agnostic. In addition, we find that HeAR is more data efficient than the baselines, sometimes reaching the same level of performance when trained on as little as 6.25% of the amount of training data (see Appendix F). This is particularly relevant to instances where labeled training data is scarce, which is unfortunately commonplace in health research. Publicly available datasets are a highly valuable resource but remain scarce and typically have few participants, making it difficult to leverage modern deep learning techniques. Our approach addresses data scarcity by enabling models to achieve adequate performance with fewer training examples than traditional methods require.

On performance of these models, we should note that the reported performance for these tasks leverage linear probes and frozen embeddings, rather than fine-tuning the whole neural network (Köhn, 2015). While this is common practice to evaluate models trained with SSL objectives and provides a fair comparison of different audio encoders, it may not yield the best performance for any given task. Performance could potentially be improved by including patient metadata as additional features of the linear probes or by fine tuning the full model. For almost all tasks, confidence intervals are very wide due to our small datasets’ sizes and high class imbalance. This makes it hard to draw statistically significant conclusions, and further validation is required for such models to become part of clinically useful tools. Some datasets, like CIDRZ and SpiroSmart, are also specific samples of the population with high disease prevalence, so performance of the models trained on those datasets may not generalize to a healthier population, and further validation is required to estimate clinical usefulness of such tools for a general population. Importantly, though not unexpected since fundamental frequencies differ on average between sexes, the fact that the representations contain information about sex should be accounted for in future model development based on these representations (Weng et al., 2024). For example, model performance should be examined stratified by sex, and any biases corrected as appropriate. It is also important to note that the performance reported on all tasks is not directly comparable to the literature because (1) prior dataset splits may not be described (Bhattacharya et al., 2023), and (2) we use linear probing rather than full fine-tuning (Köhn, 2015).

Other factors similarly affect generalization of these insights from these benchmark datasets. Evaluation on FSD50K and FluSense may not be representative of actual acoustic health events detection performance. FluSense audio clips are samples of Youtube videos, and we found that HeAR’s pretraining dataset YT-NS includes 172 of those videos (114 / 30 / 28 in train / validation / test), representing potentially 1394 clips (987 / 318 / 89 in train / validation / test). This could artificially inflate HeAR’s performance on FluSense. Similarly, BigSSL-CAP12 was trained on 900k hours of speech data from Youtube, which could also include samples from FluSense and inflate its measured performance on that dataset. On FSD50K, the timestamps of the acoustic events of interest are not available. Our coarse approximation of this information (the two-second around the loudest part of the clip) is likely to affect the labels and this is not accounted for in our training and evaluation procedures. Besides, this dataset was also used for training CLAP. These two reasons make it an imperfect benchmark.

The ad hoc preprocessing of FSD50K clips stems from the inability of most models to successfully generalize to inputs larger than what they have seen at training time. This is particularly acute in the case of HeAR, since its audio encoder is a transformer using 2-dimension fixed positional encoding trained on short two-second clips. That duration was chosen since it is typically longer than many of these health acoustic events. For tasks where the objective is to infer information about the participants (including all tasks from CIDRZ, Coswara, CoughVID, and SpiroSmart), only being able to process such short clips is sufficient. However, for detection tasks in longer audio clips like FSD50K’s and FluSense, this may not be sufficient (needle-in-a-haystack problem). The purpose of including those datasets in the evaluation procedure was to show that the same type of approach was useful for other types of sounds beyond cough alone, but specific adjustments would likely be necessary for reaching a more satisfying performance. Examples include modeling approaches not using transformers, for example with CNN and contrastive learning objectives, or using more appropriate positional encoding schemes (Kazemnejad et al., 2024).

Most models that we have experimented with are quite large and could not reasonably run on a smartphone. For reducing battery use, improving latency, preserving privacy, and making sure that this kind of technology can be used in the field with limited internet access, further research is needed to make sure that those models can run locally on-device, using techniques like distillation (Hinton et al., 2015) or quantization (Gholami et al., 2022).

We hope that our research can spur further research in the field of ML for health acoustic research. Individuals interested in getting access to use HeAR or the CIDRZ cough dataset can email health_acoustic_representations@google.com to be notified when available.

Acknowledgments

We thank Yun Liu, Rory Pilgrim, Timo Kohlberger, Eduardo Fonseca, Aren Jansen, Dan Ellis, Ryan Ehrlich, Marc Wilson from Google Research, and Luyu Wang, Lucas Smaira, Eric Lau, Chung-Cheng Chiu, Basil Mustafa from Google DeepMind for their guidance, technical support and critical feedback. We also thank Solomon Chifwamba, Pauline Musumali, Kachimba Shamaoma, Seke Muzazu, Francesca Silwamba from the Centre for Infectious Disease Research in Zambia. We also appreciate the CoughVID and Project Coswara teams for making their respective datasets publicly available, and the Google Research team for software and hardware infrastructure support. CoughVID and Coswara are licensed under a [Creative Commons Attribution 4.0 International \(CC BY 4.0\) License](https://creativecommons.org/licenses/by/4.0/) and follow the Disclaimer of Warranties and Limitation of Liability in the license.

References

Forsad Al Hossain, Andrew A Lover, George A Corey, Nicholas G Reich, and Tauhidur Rahman. Flusense: a contactless syndromic surveillance platform for influenza-like illness in hospital waiting areas. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 4(1):1–28, 2020.

- Kawther S Alqudaihi, Nida Aslam, Irfan Ullah Khan, Abdullah M Almuhaideb, Shikah J Alsunaidi, Nehad M Abdel Rahman Ibrahim, Fahd A Alhaidari, Fatema S Shaikh, Yasmine M Alsenbel, Dima M Alalharith, et al. Cough sound detection and diagnosis using artificial intelligence techniques: challenges and opportunities. *Ieee Access*, 9:102327–102344, 2021.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- Charles Bales, Muhammad Nabeel, Charles N John, Usama Masood, Haneya N Qureshi, Hasan Farooq, Iryna Posokhova, and Ali Imran. Can machine learning be used to recognize and diagnose coughs? In *2020 International Conference on e-Health and Bioengineering (EHB)*, pages 1–4. IEEE, 2020.
- Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, et al. A cookbook of self-supervised learning. *arXiv preprint arXiv:2304.12210*, 2023.
- Anastasiya Belyaeva, Justin Cosentino, Farhad Hormozdiari, Cory Y McLean, and Nicholas A Furlotte. Multimodal llms for health grounded in individual-specific data. *arXiv preprint arXiv:2307.09018*, 2023.
- Debarpan Bhattacharya, Neeraj Kumar Sharma, Debottam Dutta, Srikanth Raj Chetupalli, Pravin Mote, Sriram Ganapathy, C Chandrakiran, Sahiti Nori, KK Suhail, Sadhana Gonuguntla, et al. Coswara: A respiratory sounds and symptoms dataset for remote screening of sars-cov-2 infection. *Scientific Data*, 10(1):397, 2023.
- Veronica Boschi, Eleonora Catricala, Monica Consonni, Cristiano Chesi, Andrea Moro, and Stefano F Cappa. Connected speech in neurodegenerative language disorders: a review. *Frontiers in psychology*, 8:269, 2017.
- GHR Botha, Grant Theron, RM Warren, Marisa Klopper, Keertan Dheda, PD Van Helden, and TR Niesler. Detection of tuberculosis by automatic cough sound analysis. *Physiological measurement*, 39(4):045005, 2018.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- Chung-Cheng Chiu, James Qin, Yu Zhang, Jiahui Yu, and Yonghui Wu. Self-supervised learning with random-projection quantizer for speech recognition. In *International Conference on Machine Learning*, pages 3915–3924. PMLR, 2022.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Harry Coppock, Alex Gaskell, Panagiotis Tzirakis, Alice Baird, Lyn Jones, and Björn Schuller. End-to-end convolutional neural network enables covid-19 detection from breath and cough audio: a pilot study. *BMJ innovations*, 7(2), 2021.
- Alexander D’Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D Hoffman, et al. Underspecification presents challenges for credibility in modern machine learning. *The Journal of Machine Learning Research*, 23(1):10237–10297, 2022.
- Elizabeth R DeLong, David M DeLong, and Daniel L Clarke-Pearson. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*, pages 837–845, 1988.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic Font, and Xavier Serra. Fsd50k: an open dataset of human-labeled sound events. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30: 829–852, 2021.
- Jake Garrison. *Spiro AI: Smartphone Based Pulmonary Function Testing*. PhD thesis, 2018.
- Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pages 291–326. Chapman and Hall/CRC, 2022.
- Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *arXiv preprint arXiv:2207.06405*, 2022.
- Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. Slow-fast auditory streams for audio recognition. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 855–859. IEEE, 2021.
- Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy. The impact of positional encoding on length generalization in transformers. *Advances in Neural Information Processing Systems*, 36, 2024.
- Arne Köhn. What’s in an embedding? analyzing word embeddings through multilingual evaluation. *EMNLP*, 2015.
- Jordi Laguarda, Ferran Hueto, and Brian Subirana. Covid-19 artificial intelligence diagnosis using only cough recordings. *IEEE Open Journal of Engineering in Medicine and Biology*, 1:275–281, 2020.
- Sandra Larson, Germán Comina, Robert H Gilman, Brian H Tracey, Marjory Bravard, and José W López. Validation of an automated cough detection algorithm for tracking recovery of pulmonary tuberculosis patients. 2012.

- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Lara Orlandic, Tomas Teijeiro, and David Atienza. The coughvid crowdsourcing dataset, a corpus for the study of large-scale cough analysis algorithms. *Scientific Data*, 8(1):156, 2021.
- Madhurananda Pahar, Marisa Klopper, Byron Reeve, Rob Warren, Grant Theron, and Thomas Niesler. Automatic cough classification for tuberculosis screening in a real-world environment. *Physiological Measurement*, 42(10):105014, 2021.
- Jacob Peplinski, Joel Shor, Sachin Joglekar, Jake Garrison, and Shwetak Patel. Frill: A non-semantic speech embedding for mobile devices. *arXiv preprint arXiv:2011.04609*, 2020.
- Renard Xaviero Adhi Pramono, Syed Anas Imtiaz, and Esther Rodriguez-Villegas. A cough-based algorithm for automatic diagnosis of pertussis. *PLoS one*, 11(9):e0162128, 2016.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- Gowrisree Rudraraju, ShubhaDeepti Palreddy, Baswaraj Mamidgi, Narayana Rao Sripada, Y Padma Sai, Naveen Kumar Vodnala, and Sai Praveen Haranath. Cough sound analysis and objective correlation with spirometry and clinical diagnosis. *Informatics in Medicine Unlocked*, 19:100319, 2020.
- Björn W Schuller, Harry Coppock, and Alexander Gaskell. Detecting covid-19 from breathing and coughing sounds using deep neural networks. *arXiv preprint arXiv:2012.14553*, 2020.
- Manuja Sharma, Videlis Nduba, Lilian N Njagi, Wilfred Murithi, Zipporah Mwongera, Thomas R Hawn, Shwetak N Patel, and David J Horne. Tbscreen: A passive cough classifier for tuberculosis screening with a controlled dataset. *Science Advances*, 10(1):eadi0282, 2024.
- Joel Shor and Subhashini Venugopalan. Trillsson: Distilled universal paralinguistic speech representations. *arXiv preprint arXiv:2203.00236*, 2022.
- Joel Shor, Aren Jansen, Ronnie Maor, Oran Lang, Omry Tuval, Felix de Chaumont Quitry, Marco Tagliasacchi, Ira Shavitt, Dotan Emanuel, and Yinnon Haviv. Towards learning a universal non-semantic representation of speech. *arXiv preprint arXiv:2002.12764*, 2020.
- Joel Shor, Aren Jansen, Wei Han, Daniel Park, and Yu Zhang. Universal paralinguistic speech representations using self-supervised conformers. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3169–3173. IEEE, 2022.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, pages 1–9, 2023.
- Sangeeta Srivastava, Yun Wang, Andros Tjandra, Anurag Kumar, Chunxi Liu, Kritika Singh, and Yatharth Saraf. Conformer-based self-supervised learning for non-speech audio tasks. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8862–8866. IEEE, 2022.

- Brian H Tracey, Germán Comina, Sandra Larson, Marjory Bravard, José W López, and Robert H Gilman. Cough detection algorithm for monitoring patient recovery from pulmonary tuberculosis. In *2011 Annual international conference of the IEEE engineering in medicine and biology society*, pages 6017–6020. IEEE, 2011.
- Luyu Wang, Pauline Luc, Yan Wu, Adria Recasens, Lucas Smaira, Andrew Brock, Andrew Jaegle, Jean-Baptiste Alayrac, Sander Dieleman, Joao Carreira, et al. Towards learning universal audio representations. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4593–4597. IEEE, 2022.
- Yuxuan Wang, Pascal Getreuer, Thad Hughes, Richard F Lyon, and Rif A Saurous. Trainable frontend for robust and far-field keyword spotting. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5670–5674. IEEE, 2017.
- Wei-Hung Weng, Andrew Sellergen, Atilla P Kiraly, Alexander D’Amour, Jungyeon Park, Rory Pilgrim, Stephen Pfohl, Charles Lau, Vivek Natarajan, Shekoofeh Azizi, et al. An intentional approach to managing bias in general purpose embedding models. *The Lancet Digital Health*, 6(2):e126–e130, 2024.
- 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)*, pages 896–900. IEEE, 2020.
- Shawn Xu, Lin Yang, Christopher Kelly, Marcin Sieniek, Timo Kohlberger, Martin Ma, Wei-Hung Weng, Attila Kiraly, Sahar Kazemzadeh, Zakkai Melamed, et al. Elixr: Towards a general purpose x-ray artificial intelligence system through alignment of large language models and radiology vision encoders. *arXiv preprint arXiv:2308.01317*, 2023.
- Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- Yu Zhang, Daniel S Park, Wei Han, James Qin, Anmol Gulati, Joel Shor, Aren Jansen, Yuanzhong Xu, Yanping Huang, Shibo Wang, et al. Bigssl: Exploring the frontier of large-scale semi-supervised learning for automatic speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1519–1532, 2022.
- Yu Zhang, Wei Han, James Qin, Yongqiang Wang, Ankur Bapna, Zhehuai Chen, Nanxin Chen, Bo Li, Vera Axelrod, Gary Wang, et al. Google USM: Scaling automatic speech recognition beyond 100 languages. *arXiv preprint arXiv:2303.01037*, 2023.
- Alexandra J Zimmer, César Ugarte-Gil, Rahul Pathri, Puneet Dewan, Devan Jaganath, Adithya Cattamanchi, Madhukar Pai, and Simon Grandjean Lapierre. Making cough count in tuberculosis care. *Communications medicine*, 2(1):83, 2022.

Appendix A. Health Acoustic Event Detector

The health acoustic event detector is trained on two publicly available datasets (FSD50K and Flusense) and a proprietary health acoustic dataset. FSD50K contains over 50K audio clips (over 100 hours) annotated using AudioSet ontology (Fonseca et al., 2021), and FluSense is the subset of AudioSet dataset including sounds related to flu illnesses, which has seven labels with enough samples for running our cross-validation procedure: breathe, cough, gasp, sneeze, snuffle, speech, throat-clear, and “etc” (everything else) (Al Hossain et al., 2020). The private dataset is collected from a variety of sources. The detector uses the audio clips with labels such as “coughing”, “sneezing”, and “breathing” for training.

The detector first converts and resamples the audio to mono channel 16 kHz sampling rate, then crops the audio into two-second log-mel spectrogram features with 48 frequency bins ranging from 125Hz to 7.5kHz with per-channel energy normalization (PCEN) (Wang et al., 2017). These features are passed into a small convolutional neural network (CNN). The loss is balanced binary cross entropy, and the output of CNN is the logits for each prediction class. Detection yield for each event class in YouTube is listed in Table A1. Two classes identified by the detector (sneezing and snoring) were not used for filtering YouTube samples because they were not deemed reliable enough.

Sound Type	Yield (number of two-second audio clips)
Coughing	50,414,000
Baby coughing	1,411,000
Breathing	31,534,560
Throat clearing	4,095,000
Laughing	102,826,000
Speaking	123,024,000

Table A1: Detection yield for each health acoustic event from three billions YouTube clips.

Appendix B. Evaluation Datasets

The details of five evaluation datasets are listed in Table B1. Note that some datasets include more than one recording per participant. When this happens, predictions obtained on recordings from a given individual are averaged. Metrics in all tables (in particular, Tables 1,2,3,C1) are computed on a per-participant level (and not per-recording).

For FluSense, there may be several acoustic events occurring in each clip. When that happens, there will be several labeled crops extracted from a given clip. Predictions in that case are not aggregated at the clip level, hence the identical counts in both columns for that dataset.

Dataset	Tasks	Number of recordings used for training linear probes	Number of recordings (participants) used for validation	Number of recordings (participants) for final evaluation	Reference
FSD50K	Health acoustic event (6 tasks)	32652	8313 (8313)	10231 (10231)	Fonseca et al. (2021)
Flusense	Health acoustic event (7 tasks)	7535	1779 (1779)	2360 (2360)	Al Hossain et al. (2020)
CoughVID	COVID, sex (2 tasks)	44249	15083 (4095)	15123 (2955)	Orlandic et al. (2021)
Coswara	COVID, sex, smoking status, age (4 tasks)	10230	4285 (531)	5846 (652)	Bhattacharya et al. (2023)
CIDRZ	TB, sex, smoking status, age, BMI, 3 CXR findings (8 tasks, 3 different devices)	8210	663 (86) (low-tier) 854 (84) (mid-tier) 857 (86) (high-tier)	2107 (265) (low-tier) 2460 (265) (mid-tier) 2546 (272) (high-tier)	Dataset collected for this study. Details in the CIDRZ section below.
SpiroSmart	FEV1, FVC, FEV1/FVC ratio, PEF, FET, sex (6 tasks)	13239	696 (544)	772 (108)	Garrison (2018)

Table B1: Evaluation datasets statistics.

CIDRZ TB Dataset Description The CIDRZ TB dataset is collected by the Centre for Infectious Disease Research in Zambia. The study was approved by the University of Zambia Biomedical Ethics Committee, and all participants provided written informed consent prior to enrollment in the study. Adults who had symptoms suggestive of TB, were identified as close contacts of TB patients, or were newly diagnosed with HIV were recruited at three clinical sites (Chawama, Chainda-South, and Kanyama) in Zambia (trial NCT05139940). Audio recordings of cough sounds were obtained from 599 consented patients. To ensure robustness across different microphones, the sounds were recorded by four devices: Zoom H2N microphone (high quality audio recorder), Samsung GalaxyA22 (high-tier phone), Samsung GalaxyA12 (mid-tier phone), and Pixel3a (low-tier phone). In this work, we focused on recordings from the three phone microphones. The audio clips are recorded by the [Android application Audio Recorder](#) and encoded in the wav file format with 24-bit PCM, sampling rate of 192 kHz, and in stereo channels under a quiet environment. Cohort details are listed in Table [B2](#).

To collect cough sounds, the participant was asked to remove his/her mask and generate four cough events (three single coughs and one sequence of multiple coughs). There is a 10-15 seconds gap between cough events to enable a return to “baseline” before the next cough.

The CXR and the corresponding CXR finding annotations of the CIDRZ TB Dataset have been collected along with the cough sound collection.

HEAR - HEALTH ACOUSTIC REPRESENTATIONS

		All	Train	Tune	Test
Sample size (%)		599 (100.0)	229 (38.2)	89 (14.9)	281 (46.9)
Site (%)	Chawama	281 (46.9)	0	0	281 (100.0)
	Chainda-South	34 (5.7)	26 (11.4)	8 (9.0)	0
	Kanyama	284 (47.4)	203 (88.6)	81 (91.0)	0
Female (%)		297 (49.6)	99 (43.2)	41 (46.1)	157 (55.9)
Age [IQR]		35.0 [27.0,45.0]	36.0 [28.0,46.0]	37.0 [29.0,44.0]	33.0 [25.0,45.0]
BMI [IQR]		21.0 [19.0,24.0]	20.0 [18.0,24.0]	20.0 [18.0,23.0]	21.0 [19.0,25.0]
Positive TB (%)		92 (15.4)	46 (20.1)	18 (20.2)	28 (10.0)
Positive HIV (%)		217 (36.2)	85 (37.1)	36 (40.4)	96 (34.2)
Cough duration (%)	1 - 2 weeks	227 (37.9)	78 (34.1)	28 (31.5)	121 (43.1)
	3 - 4 weeks	24 (4.0)	5 (2.2)	3 (3.4)	16 (5.7)
	<1 week	121 (20.2)	52 (22.7)	17 (19.1)	52 (18.5)
	>4 weeks	155 (25.9)	68 (29.7)	28 (31.5)	59 (21.0)
Productive cough (%)		455 (76.0)	181 (79.0)	68 (76.4)	206 (73.3)
Hemoptysis (%)		65 (10.9)	25 (10.9)	12 (13.5)	28 (10.0)
Chest pain (%)		375 (62.6)	149 (65.1)	63 (70.8)	163 (58.0)
Short of breath (%)		204 (34.1)	84 (36.7)	35 (39.3)	85 (30.2)
Fever (%)		217 (36.2)	86 (37.6)	38 (42.7)	93 (33.1)
Night sweat (%)		245 (40.9)	96 (41.9)	36 (40.4)	113 (40.2)
Weight loss (%)		374 (62.4)	157 (68.6)	49 (55.1)	168 (59.8)
Previous TB (%)	0	497 (83.0)	187 (81.7)	70 (78.7)	240 (85.4)
	1	92 (15.4)	38 (16.6)	17 (19.1)	37 (13.2)
	2	8 (1.3)	3 (1.3)	2 (2.2)	3 (1.1)
	3	2 (0.3)	1 (0.4)	0	1 (0.4)
Tobacco use (%)	Current	127 (21.2)	55 (24.0)	22 (24.7)	50 (17.8)
	Stopped	48 (8.0)	16 (7.0)	11 (12.4)	21 (7.5)
	Never	422 (70.5)	158 (69.0)	56 (62.9)	208 (74.0)
Cigarettes per day (%)	No	472 (78.8)	174 (76.0)	67 (75.3)	231 (82.2)
	1 - 10	93 (15.5)	37 (16.2)	16 (18.0)	40 (14.2)
	11 - 20	20 (3.3)	12 (5.2)	3 (3.4)	5 (1.8)
	>20	14 (2.3)	6 (2.6)	3 (3.4)	5 (1.8)

Table B2: CIDRZ cohort descriptive statistics per split. Metadata field varies; the following table reports data where available.

Appendix C. CIDRZ TB Dataset Performance Per Recording Device Type

The performance of BigSSL-CAP12, HeAR, and CLAP is stable across recording devices on most tasks, while TRILL’s and FRILL’s vary significantly between low-tier and high-tier (Table C1). MRR for those tasks are 0.381, 0.274, 0.386, 0.786, and 0.456 for TRILL, FRILL, BigSSL-CAP12, HeAR, and CLAP, respectively.

Evaluation Recording Device	Task	TRILL	FRILL	BigSSL-12	HeAR	CLAP
Pixel3a	Focal / multi focal lung opacities	0.809 [0.747, 0.870]	0.800 [0.740, 0.860]	0.747 [0.672, 0.821]	0.794 [0.728, 0.861]	0.760 [0.690, 0.830]
GalaxyA12	Focal / multi focal lung opacities	0.743 [0.671, 0.814]	0.728 [0.655, 0.802]	0.760 [0.688, 0.831]	0.802 [0.734, 0.870]	0.746 [0.673, 0.820]
GalaxyA22	Focal / multi focal lung opacities	0.719 [0.646, 0.792]	0.686 [0.611, 0.762]	0.746 [0.672, 0.820]	0.806 [0.744, 0.867]	0.758 [0.686, 0.830]
Pixel3a	Abnormal CXR	0.815 [0.757, 0.874]	0.778 [0.712, 0.844]	0.739 [0.664, 0.814]	0.763 [0.695, 0.830]	0.734 [0.658, 0.810]
GalaxyA12	Abnormal CXR	0.730 [0.655, 0.806]	0.734 [0.662, 0.807]	0.734 [0.662, 0.807]	0.763 [0.691, 0.836]	0.743 [0.673, 0.813]
GalaxyA22	Abnormal CXR	0.729 [0.657, 0.801]	0.725 [0.653, 0.797]	0.732 [0.658, 0.806]	0.768 [0.701, 0.836]	0.746 [0.674, 0.818]
Pixel3a	Pleural effusion	0.683 [0.553, 0.812]	0.688 [0.562, 0.813]	0.684 [0.548, 0.819]	0.610 [0.465, 0.755]	0.748 [0.629, 0.866]
GalaxyA12	Pleural effusion	0.673 [0.564, 0.781]	0.632 [0.481, 0.784]	0.635 [0.499, 0.771]	0.625 [0.476, 0.774]	0.567 [0.408, 0.725]
GalaxyA22	Pleural effusion	0.637 [0.521, 0.754]	0.567 [0.433, 0.701]	0.713 [0.587, 0.838]	0.634 [0.493, 0.776]	0.730 [0.614, 0.846]
Pixel3a	Sex	0.933 [0.901, 0.965]	0.928 [0.894, 0.961]	0.936 [0.909, 0.964]	0.974 [0.958, 0.990]	0.907 [0.872, 0.942]
GalaxyA12	Sex	0.944 [0.917, 0.971]	0.941 [0.912, 0.969]	0.931 [0.899, 0.963]	0.981 [0.969, 0.993]	0.912 [0.877, 0.947]
GalaxyA22	Sex	0.951 [0.925, 0.976]	0.950 [0.924, 0.976]	0.939 [0.910, 0.967]	0.981 [0.966, 0.996]	0.900 [0.864, 0.936]
Pixel3a	Tuberculosis	0.652 [0.520, 0.784]	0.648 [0.523, 0.772]	0.659 [0.533, 0.786]	0.739 [0.636, 0.841]	0.740 [0.627, 0.853]
GalaxyA12	Tuberculosis	0.637 [0.527, 0.747]	0.631 [0.503, 0.758]	0.680 [0.571, 0.790]	0.720 [0.617, 0.824]	0.745 [0.641, 0.848]
GalaxyA22	Tuberculosis	0.675 [0.574, 0.775]	0.581 [0.457, 0.704]	0.687 [0.584, 0.791]	0.734 [0.639, 0.829]	0.794 [0.700, 0.889]
Pixel3a	Smoking status	0.822 [0.762, 0.883]	0.811 [0.747, 0.874]	0.840 [0.786, 0.895]	0.877 [0.833, 0.921]	0.808 [0.750, 0.866]
GalaxyA12	Smoking status	0.831 [0.772, 0.890]	0.830 [0.771, 0.888]	0.836 [0.781, 0.891]	0.862 [0.817, 0.907]	0.797 [0.723, 0.871]
GalaxyA22	Smoking status	0.812 [0.753, 0.871]	0.809 [0.749, 0.868]	0.835 [0.777, 0.892]	0.861 [0.811, 0.910]	0.787 [0.721, 0.853]
Pixel3a	Age	10.590 [9.733, 11.509]	10.959 [10.125, 11.855]	10.009 [9.157, 10.944]	9.316 [8.550, 10.123]	10.775 [9.934, 11.644]
GalaxyA12	Age	10.819 [9.962, 11.680]	10.898 [10.000, 11.740]	10.175 [9.366, 10.928]	9.280 [8.505, 10.049]	10.610 [9.770, 11.408]
GalaxyA22	Age	11.285 [10.429, 12.137]	11.332 [10.458, 12.216]	10.006 [9.203, 10.840]	9.554 [8.787, 10.375]	10.695 [9.881, 11.581]
Pixel3a	BMI	3.861 [3.362, 4.458]	3.860 [3.359, 4.465]	3.875 [3.366, 4.452]	3.818 [3.328, 4.397]	3.836 [3.337, 4.433]
GalaxyA12	BMI	3.853 [3.376, 4.448]	3.845 [3.358, 4.441]	3.781 [3.286, 4.369]	3.719 [3.202, 4.270]	3.786 [3.288, 4.378]
GalaxyA22	BMI	3.792 [3.324, 4.345]	3.793 [3.315, 4.342]	3.747 [3.295, 4.263]	3.695 [3.247, 4.206]	3.711 [3.250, 4.252]

Table C1: Performance comparison on cough inference tasks of CIDRZ TB Dataset per recording device type.

Appendix D. Effect of Scaling Training Data Size

We also find that scaling up the training data (YT-NS) used for training the HeAR audio encoder helps improve the linear probing performance across different downstream tasks (Figure D1).

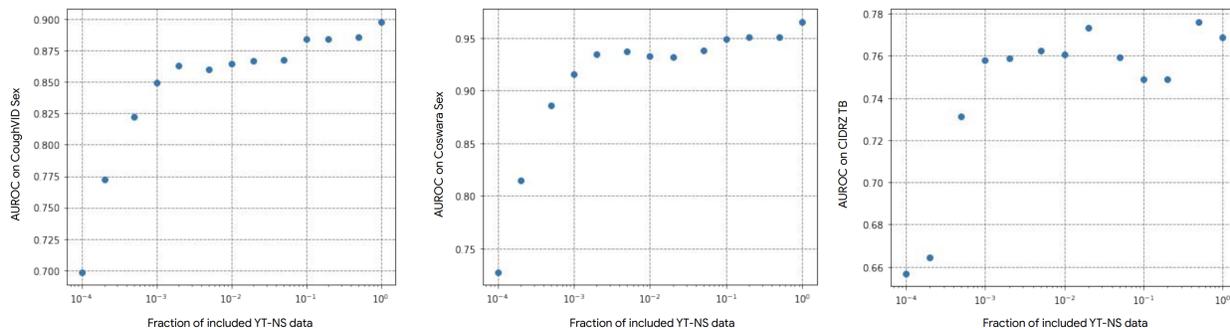


Figure D1: Scaling effect of increasing the YT-NS data size for training HeAR. We use CoughVID and Coswara sex classification, and CIDRZ tuberculosis prediction tasks as examples.

Appendix E. Generalization to Unseen Devices

Two devices are used for training the linear probes, and evaluation is done using recordings from the remaining device (i.e., out-of-distribution (OOD) device). MRR for those tasks are 0.382, 0.303, 0.358, 0.743, and 0.497 for TRILL, FRILL, BigSSL-CAP12, HeAR, and CLAP, respectively.

For all tasks, HeAR performance remains stable across all OOD devices and is consistently among the highest ranked. Other models like TRILL and FRILL have unstable performance, while CLAP and BigSSL-CAP12 are more stable but typically worse (Table E1).

OOD device	Task	TRILL	FRILL	BigSSL-CAP12	HeAR	CLAP
Pixel3a	Focal / multi focal lung opacities	0.802 [0.741, 0.864]	0.759 [0.692, 0.827]	0.744 [0.669, 0.819]	0.789 [0.722, 0.856]	0.769 [0.699, 0.840]
	Abnormal CXR	0.788 [0.723, 0.852]	0.733 [0.662, 0.805]	0.734 [0.661, 0.808]	0.765 [0.699, 0.831]	0.759 [0.688, 0.830]
Pixel3a	Pleural effusion	0.672 [0.527, 0.817]	0.720 [0.599, 0.841]	0.603 [0.462, 0.744]	0.577 [0.430, 0.724]	0.657 [0.528, 0.786]
Pixel3a	Sex	0.926 [0.893, 0.958]	0.920 [0.885, 0.954]	0.937 [0.909, 0.964]	0.973 [0.956, 0.989]	0.886 [0.845, 0.926]
Pixel3a	Tuberculosis	0.669 [0.538, 0.801]	0.651 [0.522, 0.779]	0.696 [0.575, 0.816]	0.748 [0.646, 0.849]	0.659 [0.530, 0.788]
Pixel3a	Smoking status	0.816 [0.756, 0.877]	0.805 [0.742, 0.868]	0.835 [0.780, 0.889]	0.878 [0.835, 0.921]	0.758 [0.688, 0.828]
GalaxyA12	Focal / multi focal lung opacities	0.706 [0.631, 0.782]	0.730 [0.656, 0.803]	0.751 [0.678, 0.824]	0.800 [0.731, 0.869]	0.771 [0.703, 0.839]
GalaxyA12	Abnormal CXR	0.709 [0.633, 0.785]	0.679 [0.602, 0.756]	0.726 [0.653, 0.800]	0.761 [0.688, 0.835]	0.768 [0.703, 0.832]
GalaxyA12	Pleural effusion	0.591 [0.448, 0.734]	0.642 [0.533, 0.751]	0.598 [0.456, 0.741]	0.624 [0.470, 0.778]	0.674 [0.533, 0.815]
GalaxyA12	Sex	0.950 [0.924, 0.975]	0.948 [0.922, 0.974]	0.932 [0.901, 0.964]	0.981 [0.969, 0.993]	0.904 [0.866, 0.941]
GalaxyA12	Tuberculosis	0.667 [0.568, 0.766]	0.625 [0.503, 0.747]	0.699 [0.595, 0.804]	0.720 [0.621, 0.819]	0.715 [0.603, 0.826]
GalaxyA12	Smoking status	0.829 [0.771, 0.887]	0.834 [0.777, 0.891]	0.839 [0.784, 0.894]	0.871 [0.827, 0.915]	0.759 [0.682, 0.836]
GalaxyA22	Focal / multi focal lung opacities	0.719 [0.645, 0.793]	0.671 [0.593, 0.749]	0.738 [0.664, 0.812]	0.798 [0.735, 0.861]	0.771 [0.702, 0.840]
GalaxyA22	Abnormal CXR	0.711 [0.637, 0.786]	0.732 [0.661, 0.803]	0.728 [0.655, 0.802]	0.772 [0.705, 0.839]	0.774 [0.707, 0.841]
GalaxyA22	Pleural effusion	0.620 [0.491, 0.749]	0.548 [0.415, 0.682]	0.729 [0.607, 0.852]	0.679 [0.550, 0.808]	0.757 [0.641, 0.872]
GalaxyA22	Sex	0.951 [0.925, 0.977]	0.937 [0.907, 0.967]	0.944 [0.916, 0.971]	0.982 [0.967, 0.997]	0.905 [0.869, 0.940]
GalaxyA22	Tuberculosis	0.561 [0.444, 0.678]	0.592 [0.466, 0.717]	0.633 [0.525, 0.741]	0.727 [0.632, 0.822]	0.746 [0.642, 0.849]
GalaxyA22	Smoking status	0.816 [0.759, 0.873]	0.798 [0.735, 0.861]	0.833 [0.774, 0.892]	0.856 [0.806, 0.906]	0.771 [0.704, 0.838]

Table E1: Performance on cough inference tasks of CIDRZ TB Dataset with the unseen device generalization setup.

Appendix F. Data efficiency

For all cough and spirometry tasks, we train linear probes with 6.25%, 12.5%, 25%, 50%, and 100% of the data, for all models. All models use the exact same subsampled datasets for training. All probes are evaluated on 100% of the test split. For classification tasks, we make sure to subsample each label in the same way. This procedure allows us to compare how different encoders fare in different data regimes. We find that HeAR is the most data efficient model, as evidenced by its consistently higher rank across all data regimes and all tasks (Table F1, Figure F1,F2).

Training data available (%)	TRILL	FRILL	BigSSL-CAP12	HeAR	CLAP (48k)
6.25	0.409	0.392	0.358	0.789	0.331
12.5	0.436	0.310	0.390	0.840	0.304
25	0.321	0.308	0.470	0.843	0.342
50	0.392	0.358	0.384	0.790	0.358
100	0.418	0.341	0.359	0.796	0.359

Table F1: Mean reciprocal rank (MRR) across all cough (Coswara, CoughVID, CIDRZ) and spirometry (SpiroSmart) inference tasks, for varying amounts of available training data.

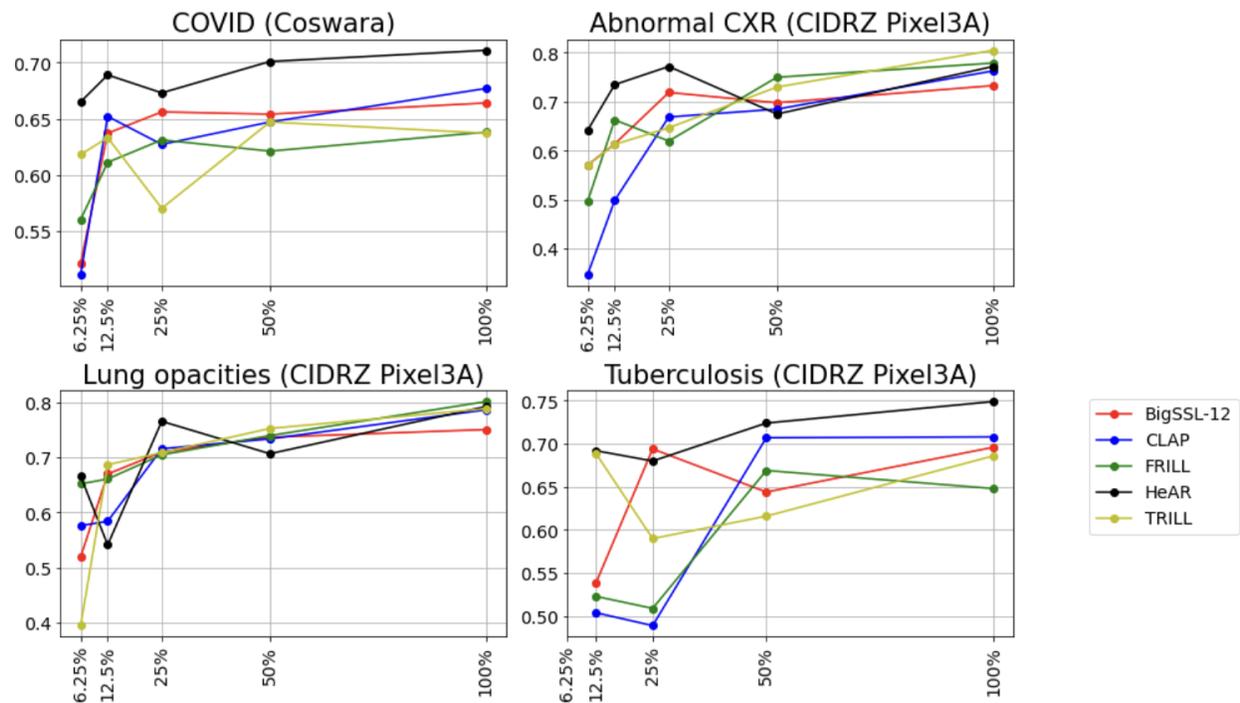


Figure F1: AUROC of all models for varying amounts of training data, for a subset of cough inference tasks.

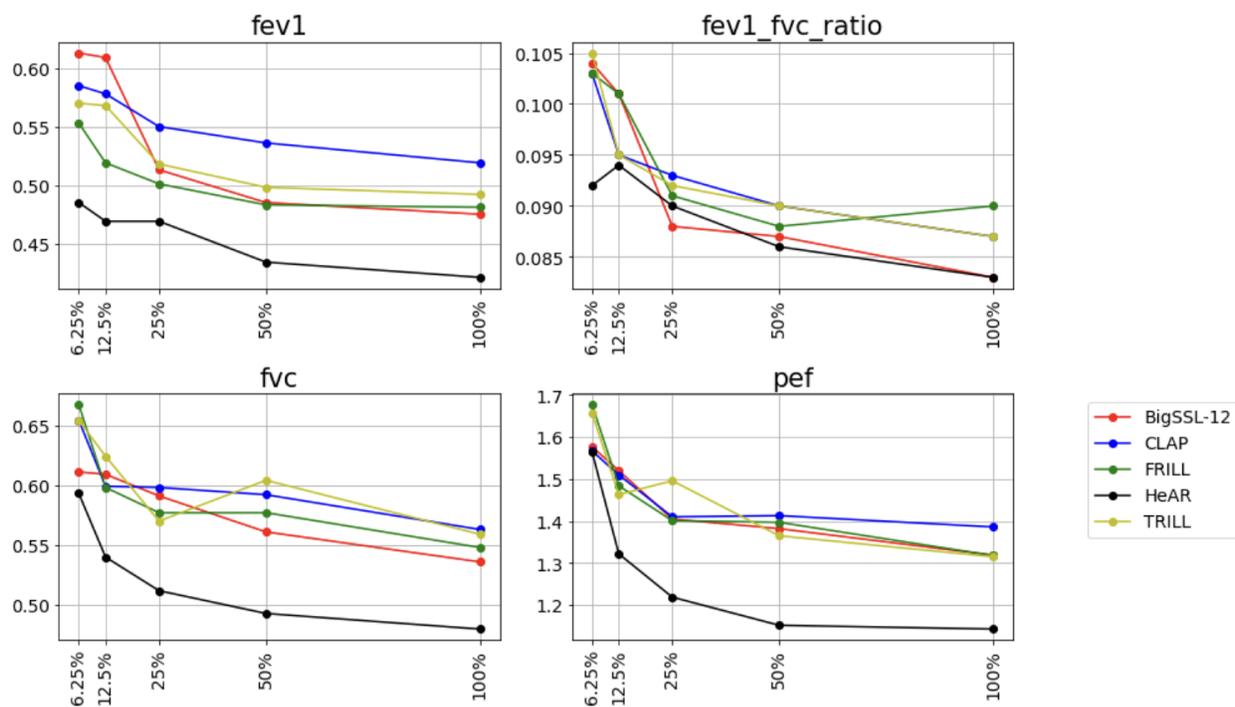


Figure F2: Mean absolute error of all models for varying amounts of training data, for a subset of spirometry inference tasks.