

SEMI-SUPERVISED PULMONARY  
AUSCULTATION ANALYSIS WITH CROSS  
PSEUDO SUPERVISION

**Jieruei Chang**

Princeton High School  
Princeton, NJ, United States  
March 2024

## ABSTRACT

Pulmonary auscultation analysis is a crucial part of the diagnosis of respiratory illnesses, yet it remains largely a manual task dependent on the skills of individual physicians. Efforts have been made to use machine learning to automate detection of abnormal lung sounds in auscultation analysis, but typical labeled datasets require experts to engage in time-consuming annotation. In contrast, unlabeled data is relatively easy to obtain. In this study, I use cross pseudo supervision to leverage a small amount of labeled audio data and a larger amount of unlabeled data to perform automated auscultation analysis. I show that this method significantly outperforms fully supervised models trained only on the labeled data.

# 1 INTRODUCTION AND PROBLEM DEFINITION

“Prevention is better than cure.” That is what Erasmus declared five hundred years ago, but those five words still ring true today. To prevent the spread of a pandemic, safe yet effective methods of diagnosis must be developed. In particular this project focuses on the diagnosis of respiratory diseases, for the last three years have shown the lightning pace at which they can proliferate. Pulmonary auscultation, or the act of listening to lung sounds with a stethoscope, provides vital information to aid in respiratory disease diagnosis by detecting the inhalation-exhalation cycle as well as abnormalities such as crackles, rhonchi, stridor, or wheezes [1]. Auscultation is non-invasive, making it a very safe procedure; but despite the simplicity of the method, it is effective because respiratory diseases directly affect how air passes through the lungs (for example by blocking or contracting airways), so the sounds produced vary as a direct byproduct of the underlying condition. It has been shown that pneumonia resulting from COVID-19 show distinctive auscultation characteristics [2]; in addition, a review of 28 studies involving a total of 2,032 patients showed that these respiratory sounds are useful indicators of respiratory illnesses such as asthma, cystic fibrosis, and bronchiolitis [3].

Although there have been great advances in healthcare throughout the last century, pulmonary auscultation largely still requires the expertise of medical professionals, due to the high degree of variability present in lung recordings and the difficulty in distinguishing normal from abnormal sounds. In recent decades, electronic stethoscopes such as the Littmann 3200 have been developed for long-duration auscultation to aid in patient monitoring, but they still require human expertise for analysis; additionally, as physicians may not have the time or auditory acuity to properly analyze the recordings produced, manual analysis can lead to incorrect diagnoses of respiratory diseases and significantly impact patient outcomes [4].

This motivates the development of machine learning models that can automate the process of detecting adventitious sounds in electronic stethoscope recordings as shown in Fig. 1. The figure shows that this is not a traditional single-label classification problem. It is not enough to simply say that “the patient is experiencing crackles,” because labels can overlap and the temporal position of these sounds is important. This is because different diseases cause adventitious sounds to appear in differing positions in the inhalation-exhalation cycle. For example, crackle sounds occurring in the middle of the inhalation phase are a symptom of interstitial lung fibrosis and pneumonia, while crackle sounds appearing at the beginning of inhalation and during exhalation are a symptom of chronic bronchitis. Together, this makes my problem a multi-label temporal segmentation problem.

In addition, traditional machine learning models require large datasets of labeled data. While lung sound recordings are relatively easy to obtain in a hospital setting, labeling the recordings requires the expertise of

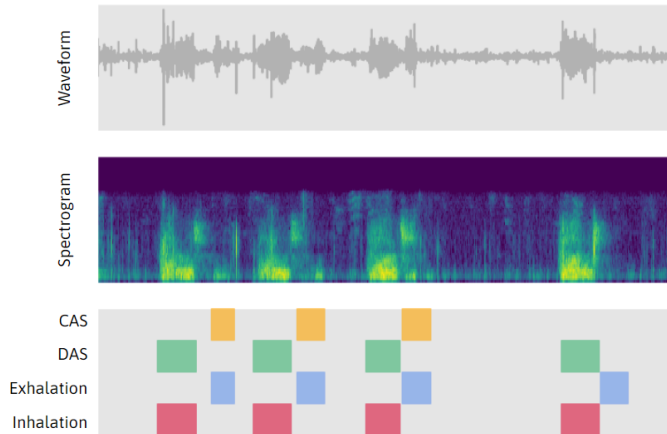


Figure 1: Illustration of lung sound recording and associated labels. From top to bottom: audio waveform, spectrogram, and annotated labels (CAS: continuous adventitious sound, DAS: discontinuous adventitious sound).

medical experts and large amounts of annotation time. Therefore, I apply semi-supervised techniques to this problem by training models on a small set of labeled auscultation data and a larger set of unlabeled data, to show the viability of such an approach in lung auscultation analysis.

Specifically, I propose a semi-supervised pulmonary auscultation analysis method based on cross pseudo supervision, using two differently initialized networks that supervise each other on unlabeled data. Furthermore, since a model that is trained from scratch will first have to learn to understand audio data before it can learn to predict respiratory sounds, I leverage model backbones pretrained with existing large-scale audio datasets.

Therefore, my contributions are as follows:

- I apply semi-supervised learning techniques using two networks to the task of pulmonary auscultation analysis. To the best of my knowledge this is the first time such techniques have been used for this problem. Particularly, I show that semi-supervised learning with cross pseudo supervision outperforms purely supervised baselines. This is also the first time that cross pseudo supervision has been employed in the audio domain.
- I show that pretraining models on large-scale audio datasets can further improve performance in pulmonary auscultation analysis.

## 2 RELATED WORK

### 2.1 SEMI-SUPERVISED LEARNING

Semi-supervised techniques allow models to leverage a large unlabeled dataset along with a smaller number of labeled examples, which is useful in scenarios such as my problem setting, where obtaining labeled data is expensive but unlabeled data is relatively abundant. One such technique is self-training, where a model predicts pseudo labels for unlabeled data, which is then used to retrain the model. Another technique is label propagation, a graph-based method that assigns labels to unlabeled data on the assumption that similar data has similar labels.

In the last few years, deep learning based approaches have emerged that utilize two networks. Mean Teacher [5] is an approach involving a student network and a teacher network, where the teacher network is the exponential moving average of past student networks. During training, the student network first undergoes supervised training with labeled data, and then is trained with unlabeled data by using pseudo labels generated by the teacher network. Guided Collaborative Training (GCT) [6] is a complex approach that essentially uses two differently-initialized networks and enforces a consistency loss between their predictions, encouraging them to produce similar outputs.

More recently, Cross Pseudo Supervision (CPS) [7] also uses two differently-initialized networks, but instead of directly using network predictions to calculate the consistency loss, it generates pseudo labels from those predictions and calculates consistency loss between each network's prediction and the other network's pseudo labels. Despite its simplicity, the CPS approach has been shown to outperform other semi-supervised methods, including both Mean Teacher and GCT, in image-based semantic segmentation tasks.

### 2.2 AUTOMATIC PULMONARY AUSCULTATION ANALYSIS AND SOUND EVENT DETECTION

Many machine learning methods have been applied to the task of pulmonary auscultation analysis. Early machine-learning based methods used multilayer perceptrons [8] and convolutional neural networks [9]. Fernando et al. [10] used temporal convolutional networks (TCNs), creating a lightweight yet performant classifier. Hsu et al. [11] used CNNs coupled with LSTMs and GRUs due to the nature of the audio domain, since recurrent networks are effective on sequential data where there are strong relationships between previous and future events.

There exists prior work that uses a similar problem setting to this report. Chamberlain et al. [12] used a two-stage approach to use labeled data in conjunction with unlabeled data, by building a feature extractor

trained on unlabeled data that generates embeddings for a classifier trained on labeled data. Lang et al. [13] use a graph-based method similar to label propagation to leverage unlabeled data. My approach, using two networks that supervise each other with unlabeled data, enables semi-supervised learning while using labeled and unlabeled data in tandem, rather than in two distinct components.

The Detection and Classification of Audio Scenes and Events (DCASE) Challenge [14] includes a similar task, semi-supervised sound event detection. This is a similar problem, which involves determining the sounds present in a recording and their respective start and end times. One approach uses sound separation techniques that aim to separate individual sound sources from an audio mixture; by isolating specific sound events from the mixture, it becomes easier to detect and classify them [15]. Additionally, some researchers have modified and applied the Detection Transformer (DETR) approach to audio event detection [16]. These are interesting approaches that are outside the scope of this study, but may be relevant in future research. However, my problem is more difficult because all sounds recorded are human-produced, so the class separation is smaller (i.e., there is a larger disparity between human speech and the hum of a vacuum cleaner than between the sounds of inhalation and exhalation). Therefore, distinguishing respiratory sounds requires a model capable of detecting much more nuanced differences in audio data.

### 2.3 NETWORK PRETRAINING

In recent years, pretraining has proven to be a powerful technique for enhancing the performance of deep neural networks in various domains. The theory is that a network can be first trained on a larger dataset to learn to create useful representations from data, before being fine-tuned for specific tasks. Notably, pretraining on large-scale image datasets such as ImageNet [17] has shown success in improving the performance of image-based machine learning models, as well as accelerating model convergence during training. Models such as ResNet50 [18] pretrained on ImageNet are commonly used as backbones in computer vision models.

Similarly, AudioSet [19] is a large-scale dataset for audio-based applications, enabling the creation of pretrained audio neural networks (PANNs) [20] as feature extractors for the audio domain.

## 3 APPROACH

My model uses a twin network architecture, with each network consisting of a feature extractor generating embeddings from audio spectrogram data, and a multi-label classifier that detects respiratory sound events using those embeddings. The two networks supervise each other on both labeled and unlabeled data; the proposed overall network structure is shown in Fig. 2.

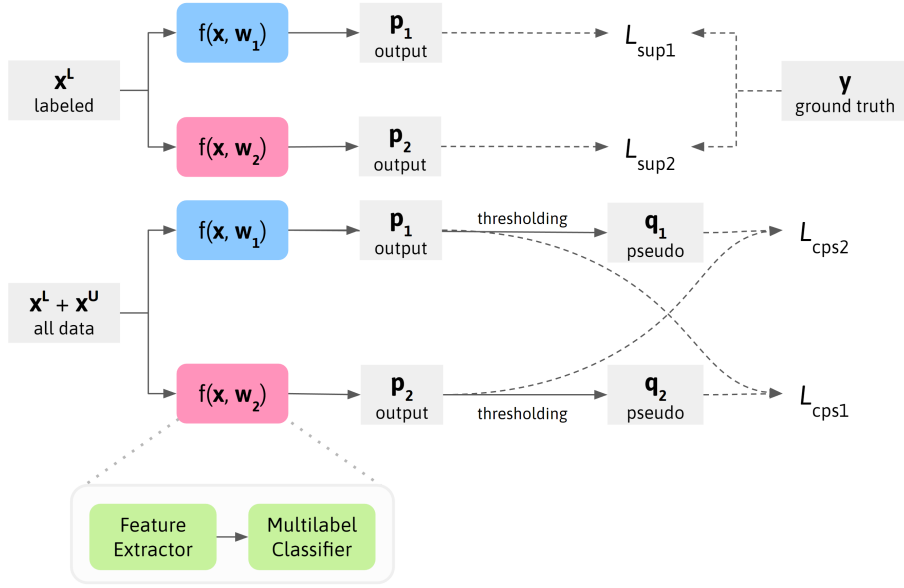


Figure 2: Overall Model Architecture.

### 3.1 CROSS PSEUDO SUPERVISION

As shown in Fig. 2, Cross Pseudo Supervision [7] enables semi-supervised learning using two networks  $f$  that are structurally identical but differently initialized with weights  $w_1$  and  $w_2$ . It enforces consistency regularization by having the networks each generate pseudo labels from unlabeled data, which are then used to supervise the other network. Whereas [7] uses one-hot pseudo labels for semantic segmentation, this is infeasible in my problem setting because multiple classes of sounds may be present at the same time. Instead, I generate pseudo labels via thresholding. The CPS loss is then calculated as follows:

$$\begin{aligned} \mathcal{L}_{cps} = \frac{1}{N} \sum_{n=1}^N & \ell_{bce}(f(\mathbf{x}_n, w_1), T_{\tau}(f(\mathbf{x}_n, w_2))) \\ & + \ell_{bce}(f(\mathbf{x}_n, w_2), T_{\tau}(f(\mathbf{x}_n, w_1))), \end{aligned} \quad (1)$$

where  $f(\mathbf{x}_n, w_k)$  is the output of a network initialized with weights  $w_k$  and given input  $\mathbf{x}_n$ ,  $N$  is the number of samples (both labeled and unlabeled),  $T_{\tau}(\theta)$  is a thresholding function that returns 1 for  $\theta > \tau$  and 0 otherwise, and  $\ell_{bce}$  is the weighted binary cross-entropy loss function, defined as

$$\begin{aligned} \ell_{bce}(\mathbf{p}, \mathbf{q}) = & -\frac{1}{CD} \sum_{c=1}^C W_c \sum_{d=1}^D \mathbf{q}_{c,d} \cdot \log(\mathbf{p}_{c,d}) \\ & +(1 - \mathbf{q}_{c,d}) \cdot \log(1 - \mathbf{p}_{c,d}), \end{aligned} \quad (2)$$

where  $C$  is the number of output classes,  $D$  is the time dimension,  $W_c$  is the weight of class  $c$  in the loss calculation,  $\mathbf{p}$  is the output probabilities, and  $\mathbf{q}$  is the labels (which are pseudo labels in the case of CPS loss). I use binary cross-entropy instead of the cross-entropy loss used in [7] due to the multi-label problem setting. To deal with large class imbalances, I weight the binary cross entropy losses for each class using  $W_c = F_{max}/F_c$ , where  $F_{max}$  is the number of time frames in the labeled dataset containing the most common class and  $F_c$  is the number of time frames containing class  $c$ .

The supervised loss is defined as

$$\begin{aligned} \mathcal{L}_{sup} = & \frac{1}{N_{lab}} \sum_{n=1}^{N_{lab}} [\ell_{bce}(f(\mathbf{x}_n, w_1), \mathbf{y}_n) \\ & + \ell_{bce}(f(\mathbf{x}_n, w_2), \mathbf{y}_n)], \end{aligned} \quad (3)$$

where  $N_{lab}$  is the number of labeled samples in the dataset and  $\mathbf{y}$  is the ground truth labels. The final loss is the weighted sum of supervised loss and CPS loss:

$$\mathcal{L} = \mathcal{L}_{sup} + \lambda \mathcal{L}_{cps}, \quad (4)$$

where  $\lambda$  is a weighting parameter to balance these two losses. The cross supervision using pseudo labels drives the networks towards the same extreme when their outputs are similar, improving the clarity of the decision boundary. If the network outputs are dissimilar, this effect is negated, so unlabeled samples with inconsistent pseudo labels are de-emphasized. As the training process continues and the pseudo labels improve, the unlabeled data effectively helps to expand the training dataset.

### 3.2 FEATURE EXTRACTOR AND MULTI-LABEL CLASSIFIER

Log-mel spectrograms generated from audio data are fed into a feature extractor module, which produces an intermediate representation that is passed to a multi-label classifier module as shown in Fig. 3. I use the CNN architecture from PANN [20] with or without AudioSet pretraining as my feature extractor. Since PANN is designed primarily for clip-level audio source classification and not framewise sound event detection, I



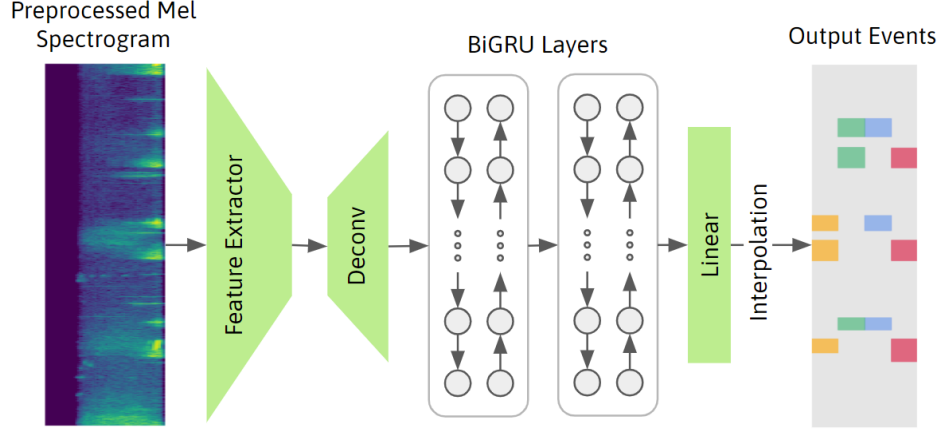


Figure 3: Single Network Architecture.

modify the pretrained model with deconvolutional layers to increase the model’s time resolution. The multi-label classifier module consists of two bidirectional gated recurrent unit (BiGRU) layers coupled to a linear output layer, followed by an interpolator to match the dimensions of the outputs and labels. GRUs are well-suited for time-domain tasks due to their ability to capture long-range temporal dependencies and patterns in sequential data [11]. GRUs are described with the formulas

$$r_t = \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{t-1} + b_{hr}) \quad (5)$$

$$z_t = \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{t-1} + b_{hz}) \quad (6)$$

$$k_t = \tanh(W_{ik}x_t + b_{in} + r_t \odot (W_{hk}h_{t-1} + b_{hk})) \quad (7)$$

$$h_t = (1 - z_t) \odot k_t + z_t \odot h_{t-1}, \quad (8)$$

where  $x_t$  is the input at time  $t$ ,  $h_t$  is the hidden state/output at time  $t$ , and  $\odot$  is the Hadamard product.  $r_t$  is a “reset gate” that controls how much of the hidden state to discard when creating the new candidate hidden state  $k_t$ .  $z_t$  is an “update gate” that controls the weighting of the linear interpolation between the previous hidden state  $h_{t-1}$  and the new candidate hidden state  $k_t$ .  $W$  and  $b$  are learnable weight and bias terms respectively. This structure is also shown graphically in Fig. 4.

The bidirectional aspect of a BiGRU allows it to consider past and future information simultaneously, which further improves its ability to find relationships between frames in the time domain; this gives the model more contextual information when trying to differentiate between difficult classes. A BiGRU essentially contains two GRU cells, one of which runs forward through the time domain and one of which runs backward through time.

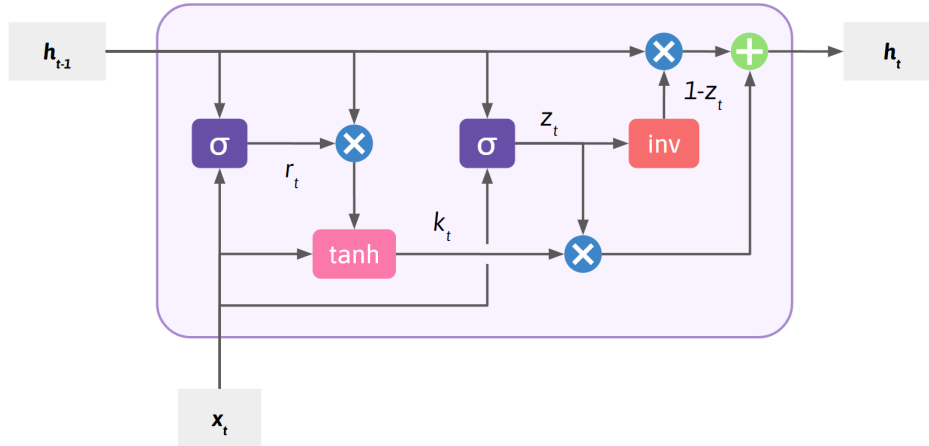


Figure 4: Gated Recurrent Unit Structure.

The structure of a two-layer BiGRU is shown in Fig. 5. In each layer, one cell is recurrent in the forward direction, and one cell is recurrent in the backward direction. The hidden states from each cell in the first layer is passed to the corresponding cell in the second layer as shown in the figure. The outputs from both the forward and backward GRUs are concatenated to produce the final output.

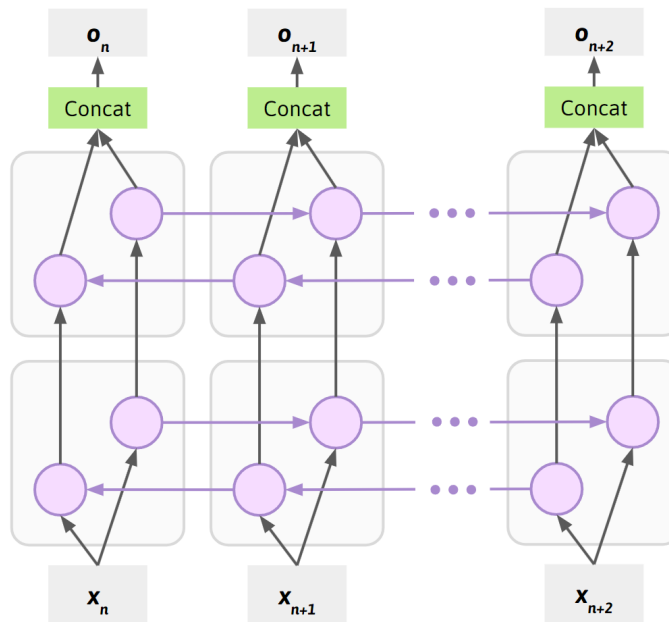


Figure 5: Bidirectional GRU Structure.

## 4 EXPERIMENTAL RESULTS

### 4.1 SETUP

#### 4.1.1 DATASETS

I train and evaluate my models on *HF\_Lung\_VI* [11], the largest publicly available dataset of lung sound recordings, as opposed to smaller datasets such as *ICBHI* [21]. It comprises 9,765 strongly labeled 15-second sound recordings from 279 patients. The annotations contain 34,095 inhalation labels, 18,349 exhalation labels, 13,883 continuous adventitious sound labels (comprising 4,740 rhonchi labels, 8,457 wheeze labels, and 686 stridor labels), and 15,606 discontinuous adventitious (crackling) sound labels. The dataset contains recordings from two different machines: a Littmann 3200 electronic stethoscope and an HF-Type-1, a recording device custom-built by the dataset authors. I only use recordings from the Littmann as only 18 out of 279 patients were recorded using the HF-Type-1 device, so the samples generated are not as representative of the wider population.

#### 4.1.2 EVALUATION

I evaluate performance using the area under the receiver operating characteristic curve (AUC) and the F1-score. The receiver operating characteristic (ROC) curve plots true positive rate with respect to the false positive rate at different binary classification thresholds; the AUC is the integral of the ROC curve. The F1-score is defined as

$$F_1 = \frac{TP}{TP + \frac{1}{2}(FP + FN)}, \quad (9)$$

where  $TP$  is the number of true positives,  $FP$  is the number of false positive, and  $FN$  is the number of false negative. For each model, I calculate AUC and F1 for each of the four classes: inhalation, exhalation, continuous adventitious sound (CAS), and discontinuous adventitious sound (DAS). Following the dataset evaluation setup [11], I group rhonchi, wheeze, and stridor into a single class. I use the train/test split provided by the dataset authors, so there is no cross-contamination between training and testing data (i.e., each patient’s data is present only in the training or in the testing set). I use a validation split comprising 10% of the total training data, also taking care to avoid cross-contamination. In total, 2,929 recordings (labeled and unlabeled) are used for training, 325 are used for validation, and 1,250 are used for evaluation.

### 4.1.3 DATA PREPROCESSING

To reduce noise in the audio data, I first run raw audio waveforms through a tenth-order Butterworth high pass filter to remove frequency data below 100Hz, as that is the lower frequency bound of lung sounds.

I further conduct noise reduction using nonstationary spectral gating [22]. This calculates a short-time Fourier transform (STFT) of the noisy signal to create a spectrogram representation. The spectrogram is smoothed across the time domain to create a moving average of each frequency band, which is used to create a noise mask that marks areas of the spectrogram that have sufficiently high amplitude. The noise mask is smoothed with a 64ms filter and applied to the original spectrogram with

$$S_{\text{denoised}} = p(M \odot S_{\text{raw}}) + (1 - p)S_{\text{raw}}, \quad (10)$$

where  $M$  is the noise mask,  $S_{\text{raw}}$  is the original spectrogram,  $S_{\text{denoised}}$  is the final output,  $\odot$  denotes the Hadamard product, and  $p$  controls the amount of noise suppression. I set  $p$  to 0.8 rather than 1.0 to avoid removing any signal that appears to be noise but is in fact a respiratory sound. An inverse short-time Fourier transform is then applied to  $S_{\text{denoised}}$  to obtain a noise-reduced audio waveform.

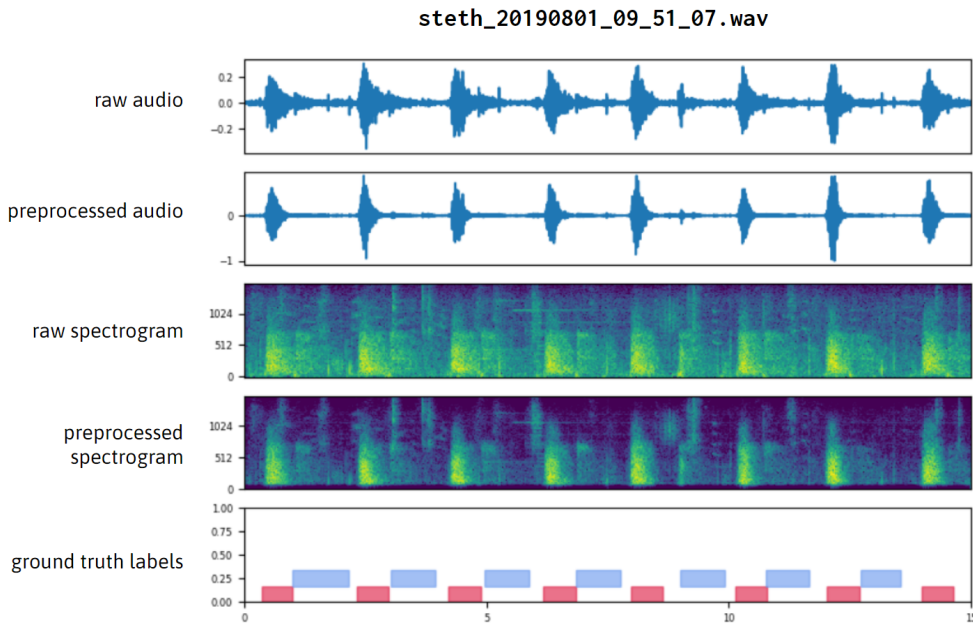


Figure 6: Effect of Audio Preprocessing.

I then generate log-mel spectrograms from these audio waveforms using 64 Mel bins, an STFT window of 256 samples and a hop size of 80 samples, then apply min-max normalization. In general, this preprocessing is effective in reducing the hospital white noise that is present in most of the recordings; this is reflected in the spectrogram outputs, as the less salient parts of the spectrogram are dimmed, better emphasizing the ground truth labels. Fig. 6 shows that the inhalation (red) and exhalation (blue) sounds are still clearly visible in the preprocessed spectrogram, whereas the background noise is largely removed.

#### 4.1.4 DATA AUGMENTATION

Data augmentation is a technique that can be used to artificially increase the size of a limited training dataset and reduce model overfitting. Whereas image data is amenable to augmentations such as linear translation and rotation, this generally is not possible with audio data as shifting the data upwards or downwards would create changes in the frequency domain, possibly changing the type of respiratory sound. Instead, the data is augmented by shifting it forwards and backwards in the time domain, rolling it such that data at the end of a segment is moved to the front and vice versa. This is allowed because audio data is time-invariant; a crackle after four seconds of recording should sound the same as a crackle after ten seconds.

#### 4.1.5 IMPLEMENTATION

Models are implemented using the PyTorch framework. I initialize the weights of convolutional and linear layers with Kaiming initialization and initialize BiGRUs with orthogonal initialization. For pretrained models, the feature extractor is initialized with the weights of a PANN trained on 8kHz data from AudioSet. I train using the AdamW optimizer with a weight decay of 0.001, a polynomial learning rate (initialized at 0.0001) with power 0.9, and a batch size of 16. I set  $\lambda$  to 1.0 in all experiments. Because Cross Pseudo Supervision generally makes the model training less stable, for CPS models I first start with 50 epochs of non-CPS training, then take the model with the best sum of F1 scores on the validation set to continue with 50 epochs of CPS training using both labeled and unlabeled training data. All models tend to converge by the 30<sup>th</sup> epoch. Detailed model specifications are listed in Table 1.

The first main section, the *Feature Extractor*, is an off-the-shelf AudioSet-pretrained CNN model. Since this architecture significantly reduces the time dimension to just 46 (from the original 1501), I designed a series of deconvolutional layers (in reality a fractionally-strided convolution) to somewhat compensate for the reduced time dimension. The deconvolutional layers intentionally do not fully restore the time dimension, as such fine-grained time resolution is not necessary and doing so would substantially increase the training/execution time cost due to the recurrent nature of the GRU.

The GRU itself is two-layered, as small-scale experiments showed that one GRU layer was insufficient to properly model the data distribution and three GRU layers were prone to overfitting. The final linear layer generates the multi-label temporal segmentation, using a sigmoidal function for output instead of softmax normalization. The interpolator uses a simple nearest-neighbor algorithm to match the shape of the model output and the ground truth so the loss can be calculated.

Table 1: Single Network Parameter Specifications

<b>Component</b>	<b>Output Shape</b>
<i><b>Feature Extractor</b></i>	
Log-mel Spectrogram (64 mel bins, 1501 frames)	[16, 64, 1501]
Convolution Block 1: ( $3 \times 3$ , BN, ReLU) $\times 2$	[16, 64, 1501, 64]
$2 \times 2$ Average Pooling	[16, 64, 750, 32]
Convolution Block 2: ( $3 \times 3$ , BN, ReLU) $\times 2$	[16, 128, 750, 32]
$2 \times 2$ Average Pooling	[16, 128, 375, 16]
Convolution Block 3: ( $3 \times 3$ , BN, ReLU) $\times 2$	[16, 256, 375, 16]
$2 \times 2$ Average Pooling	[16, 256, 187, 8]
Convolution Block 4: ( $3 \times 3$ , BN, ReLU) $\times 2$	[16, 512, 187, 8]
$2 \times 2$ Average Pooling	[16, 512, 93, 4]
Convolution Block 5: ( $3 \times 3$ , BN, ReLU) $\times 2$	[16, 1024, 93, 4]
$2 \times 2$ Average Pooling	[16, 1024, 46, 2]
Convolution Block 6: ( $3 \times 3$ , BN, ReLU) $\times 2$	[16, 2048, 46, 2]
<i><b>Deconvolution</b></i>	
$2 \times 2$ Deconvolution	[16, 1024, 93, 5]
$2 \times 2$ Deconvolution	[16, 512, 187, 11]
$2 \times 2$ Deconvolution	[16, 256, 375, 23]
Linear	[16, 375, 256]
<i><b>Classifier</b></i>	
Bidirectional GRU $\times 2$	[16, 375, 2048]
Linear	[16, 375, 4]
Interpolator	[16, 1500, 4]
<b>Total Parameters</b>	<b>113,536,275</b>

Table 2: Ablation study (models trained on 1/8 of labeled data; higher numbers are better)

Component				AUC (%)				F1 (%)			
CNN	BiGRU	Pretrain	CPS	Inhale	Exhale	DAS	CAS	Inhale	Exhale	DAS	CAS
✓				89.48	74.22	79.30	84.36	70.23	32.61	35.30	32.17
✓	✓			92.64	79.10	78.37	86.49	73.62	38.57	34.98	34.08
✓	✓	✓		92.44	80.77	83.79	92.65	73.43	40.52	44.19	49.69
✓	✓	✓	✓	<b>93.59</b>	<b>83.99</b>	<b>87.64</b>	<b>93.88</b>	<b>75.80</b>	<b>45.76</b>	<b>49.16</b>	<b>53.57</b>

## 4.2 RESULTS

### 4.2.1 ABLATION STUDY

The results of adding each individual component of my auscultation analysis model are shown in Table 2; overall, they each generally improve the model performance. The best-performing model is the one with all three additions to the most basic CNN architecture: BiGRU recurrent layers, AudioSet pretraining, and semi-supervised learning with Cross Pseudo Supervision.

The addition of BiGRU layers produces significant improvement on inhalation and exhalation. Improvement is particularly evident on the exhalation class, with an 4.88% (percentage points) increase in AUC and an 5.96% increase in F1. This demonstrates that the BiGRU allows the model to exploit time-wise relationships due the strong causal correlation between inhalation and exhalation (i.e., exhalation events always follow inhalation events). Pretraining on AudioSet results in a large improvement in both AUC and F1 on the DAS and CAS classes, showing that this pretraining improves the model generalizability in low-data situations, as these are the two classes with the fewest number of training examples. The DAS performance improves by 9.21% on the F1-score and the CAS improves by 15.61% on F1-score. Cross Pseudo Supervision leads to substantial increases across all classes, showing its ability to leverage the unlabeled data to expand the training dataset via pseudo-labeling. On the AUC, the inhale, exhale, DAS, and CAS classes increase in performance by 1.15%, 3.22%, 3.85%, and 1.23% respectively; on the F1-score, they increase by 2.37%, 5.24%, 4.97%, and 3.88% respectively.

Table 3: Comparison of F1 Scores across various partitions of labeled data for baseline and proposed models.

Partition	Model	Inhalation	Exhalation	DAS	CAS
1/2	Baseline	75.87%	46.59%	46.24%	42.06%
	Proposed	<b>76.65%</b>	<b>49.90%</b>	<b>50.46%</b>	<b>55.63%</b>
1/4	Baseline	74.52%	44.19%	41.13%	40.80%
	Proposed	<b>75.50%</b>	<b>47.09%</b>	<b>50.55%</b>	<b>53.95%</b>
1/8	Baseline	73.62%	38.57%	34.98%	34.08%
	Proposed	<b>75.80%</b>	<b>45.76%</b>	<b>49.16%</b>	<b>53.57%</b>
1/16	Baseline	71.62%	37.80%	30.24%	25.69%
	Proposed	<b>74.31%</b>	<b>43.19%</b>	<b>44.67%</b>	<b>53.10%</b>

#### 4.2.2 COMPARISON WITH BASELINES

I compare my semi-supervised proposed model to the supervised baseline model (similar in architecture to *CNN+BiGRU*, the best model evaluated in [11]) on 1/2, 1/4, 1/8, and 1/16 labeled data partitions using the AUC in Fig. 7 and using the F1-score in Table 3. The proposed model outperforms the baseline in all partitions. Observe that supervised baseline performance drops significantly as the amount of labeled data decreases, while the performance of the proposed model does not degrade as dramatically, which shows that the semi-supervised approach becomes more effective at smaller labeled data partitions. This is particularly evident on the CAS/DAS classes; the baseline performance for DAS drops by 23 percentage points on the AUC whereas the performance of the proposed model drops only by 3 percentage points, and the baseline performance for CAS drops by 16 percentage points on the F1 whereas the proposed approach only drops by 3 percentage points. On a few occasions, a proposed model trained on lesser amounts of labeled data slightly outperforms a proposed model trained on more labeled data; this is most likely an artifact of randomness during the training procedure. However, the size of the total training dataset (both labeled and unlabeled) remains constant when using semi-supervised learning, so the amount of unlabeled data decreases as the amount of labeled data increases. This then underscores how effectively the proposed method capitalizes on the unlabeled data as it almost matches the contribution of the labeled data in enhancing model performance.



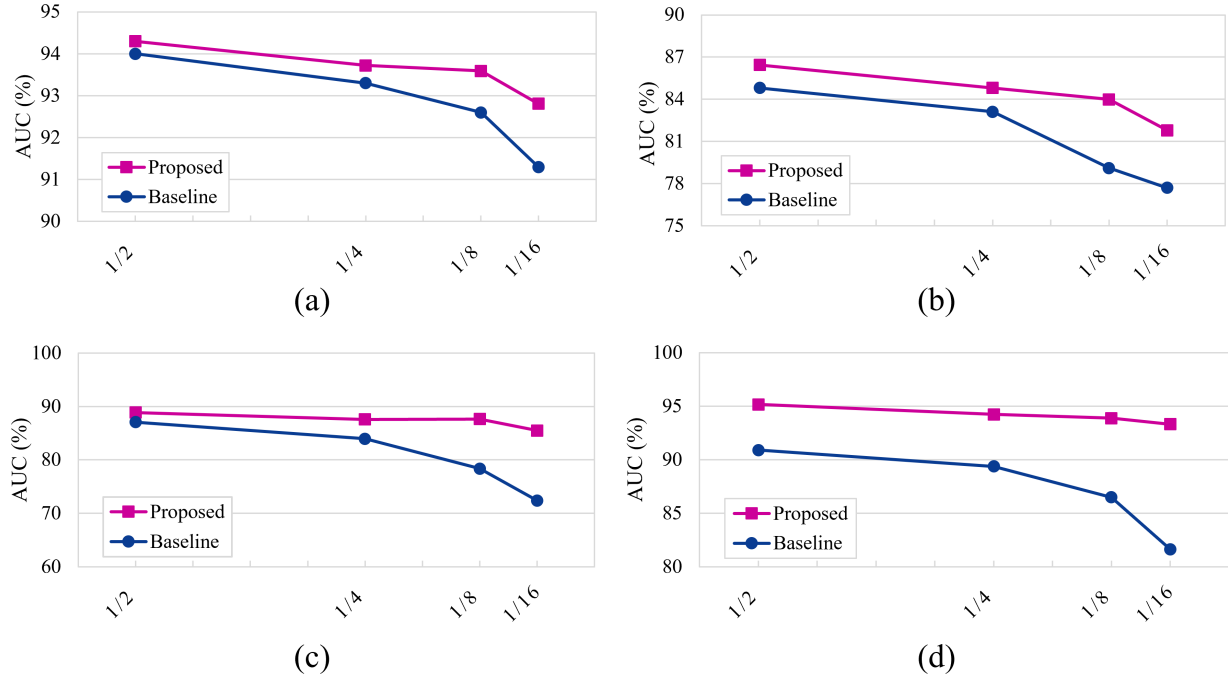


Figure 7: AUC performance of supervised baselines (blue circles) and the proposed method (pink squares) for all four classes: (a) Inhalation, (b) Exhalation, (c) DAS, (d) CAS.

#### 4.2.3 QUALITATIVE ANALYSIS

I provide some sample model outputs in Fig. 8 under two different labeled data partitions to qualitatively verify the results. In both cases, the semi-supervised model provides results much closer to the ground truth labels than the baseline. This is especially true in classes with fewer examples (namely, DAS and CAS): on the 1/2 partition, all four DAS sounds are correctly identified by the proposed model whereas the baseline model identifies two CAS sounds, and on the 1/8 partition, the proposed model is relatively accurate in labeling CAS sounds whereas the baseline model misses many (and also fails to identify two inhalation sounds). This further shows the suitability of my proposed approach in low-data problem settings.

#### 4.3 LIMITATIONS

This study has several limitations. Firstly, both labeled and unlabeled data come from the same source, so the problem setup may not fully reflect a real-world scenario. Future work can focus on using labeled and unlabeled data from different sources. Secondly, the dataset authors acknowledge that the labels are somewhat inaccurate; each sample was only annotated once, and exhalations were not labeled if the sound was not

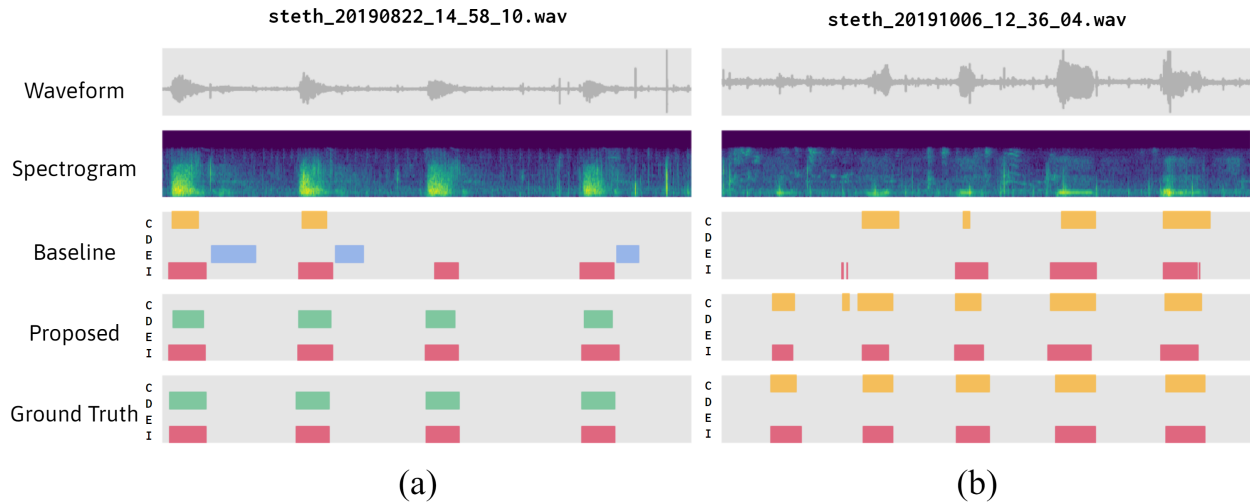


Figure 8: Qualitative results for baseline and proposed methods trained on (a) 1/2 of labeled data and (b) 1/8 of labeled data.

sufficiently clear [11]. Finally, while CPS generally leads to better performance on the evaluation metrics, its addition makes training more unstable since generated pseudo labels may be incorrect. This sometimes leads to rapid deterioration of model performance after a period of training. Although CPS intrinsically de-emphasizes inconsistent pseudo labels between the two networks, it does not safeguard against cases where both networks produce pseudo labels that are similar but incorrect. Reweighting approaches such as those proposed in [23] could enable the model to better detect and either disregard or de-emphasize poor pseudo labels.

## 5 CONCLUSION

In this work, I have shown that Cross Pseudo Supervision, coupled with network pretraining and BiGRUs to understand temporal relationships, outperforms purely supervised baselines in semi-supervised pulmonary auscultation analysis. The improvements brought on by the usage of semi-supervised techniques show that they hold significant promise in lung auscultation analysis. The dataset used in this study was necessarily fully labeled, in order to test different partitions of labeled and unlabeled data; however, by demonstrating that unlabeled data can be effectively utilized, this study allows for the generation of much larger unlabeled datasets to further enhance the performance of pulmonary auscultation analysis systems. I am confident that continuing research into semi-supervised auscultation analysis will lead to more robust and effective diagnostic tools for respiratory conditions, benefiting both patients and healthcare providers.

## ACKNOWLEDGEMENTS

I would like to thank Prof. Ching-Chun Huang for his guidance during this project.

## REFERENCES

- [1] M. Sarkar, I. Madabhavi, N. Niranjana, and M. Dogra, "Auscultation of the respiratory system," *Annals of Thoracic Medicine*, vol. 10, no. 3, pp. 158–168, 2015.
- [2] B. Wang et al., "Characteristics of Pulmonary Auscultation in Patients with 2019 Novel Coronavirus in China," *Respiration*, vol. 99, no. 9, pp. 755–763, 2020, doi: <https://doi.org/10.1159/000509610>.
- [3] V. Abreu, A. Oliveira, J. Alberto Duarte, and A. Marques, "Computerized respiratory sounds in paediatrics: A systematic review," *Respiratory Medicine: X*, vol. 3, p. 100027, Nov. 2021.
- [4] J. Heitmann et al., "DeepBreath—automated detection of respiratory pathology from lung auscultation in 572 pediatric outpatients across 5 countries," *NPJ Digital Medicine*, vol. 6, no. 1, Jun. 2023.
- [5] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," *Neural Information Processing Systems*, 2017.
- [6] Z. Ke, D. Qiu, K. Li, Q. Yan, and R. Lau. "Guided collaborative training for pixel-wise semi-supervised learning," in *European Conference on Computer Vision (ECCV)*, 2020.
- [7] X. Chen, Y. Yuan, G. Zeng and J. Wang, "Semi-supervised semantic segmentation with cross pseudo supervision," in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Nashville, TN, USA, pp. 2613-2622, 2021.
- [8] S. Rietveld, M. Oud, and E. H. Dooijes, "Classification of asthmatic breath sounds: Preliminary results of the classifying capacity of human examiners versus artificial neural networks," *Computers and Biomedical Research*, vol. 32, no. 5, pp. 440–448, Oct. 1999.
- [9] M. Aykanat, Ö. Kılıç, B. Kurt, and S. Saryal, "Classification of lung sounds using convolutional neural networks," *EURASIP Journal on Image and Video Processing*, vol. 2017, no. 1, 2017.
- [10] T. Fernando, S. Sridharan, S. Denman, H. Ghaemmaghami and C. Fookes, "Robust and interpretable temporal convolution network for event detection in lung sound recordings," in *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 7, pp. 2898-2908, 2022.

- [11] F.-S. Hsu et al., “Benchmarking of eight recurrent neural network variants for breath phase and adventitious sound detection on a self-developed open-access lung sound database—HF\_Lung\_V1,” *PLOS ONE*, vol. 16, no. 7, p. e0254134, Jul. 2021.
- [12] D. Chamberlain, R. Kodgule, D. Ganelin, V. Miglani and R. R. Fletcher, “Application of semi-supervised deep learning to lung sound analysis,” 2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Orlando, FL, USA, 2016.
- [13] R. Lang, Y. Fan, G. Liu, and G. Liu, “Analysis of unlabeled lung sound samples using semi-supervised convolutional neural networks,” *Applied Mathematics and Computation*, vol. 411, p. 126511, 2021.
- [14] N. Turpault, R. Serizel, A. Shah, and J. Salamon, “Sound event detection in domestic environments with weakly labeled data and soundscape synthesis,” In *Workshop on Detection and Classification of Acoustic Scenes and Events*, New York City, USA, 2019.
- [15] I. Kavalerov et al., “Universal sound separation,” in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 175–179, 2019.
- [16] Z. Ye et al., “Sound Event Detection Transformer: An event-based end-to-end model for sound event detection,” *arXiv.org*, Nov. 11, 2021. <https://arxiv.org/abs/2110.02011>
- [17] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li and L. Fei-Fei, “ImageNet: A large-scale hierarchical image database,” *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [18] K. He, X. Zhang, S. Ren and J. Sun, “Deep residual learning for image recognition,” 2016 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016, pp. 770-778.
- [19] J. F. Gemmeke et al., “Audio Set: An ontology and human-labeled dataset for audio events,” 2017 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, New Orleans, LA, USA, 2017, pp. 776-780.
- [20] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, “PANNs: Large-scale pretrained audio neural networks for audio pattern recognition,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2880–2894, 2020.
- [21] B. Rocha et al., “An open access database for the evaluation of respiratory sound classification algorithms,” *Physiological Measurement*, 2019.

- [22] T. Sainburg and T. Q. Gentner, “Toward a Computational Neuroethology of Vocal Communication: From Bioacoustics to Neurophysiology, Emerging Tools and Future Directions,” *Frontiers in Behavioral Neuroscience*, vol. 15, 2021.
- [23] Z. Ren, R. Yeh, and A. Schwing, “Not all unlabeled data are equal: Learning to weight data in semi-supervised learning,” in *Neural Information Processing Systems*, 2020.