
Necessity of Uncertainty Quantification for Audio-driven Healthcare Diagnosis

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep learning excels in analyzing multi-modal signals for healthcare diagnostics
2 but lacks the ability to quantify confidence in the predictions, which can lead to
3 overconfident, erroneous diagnoses. In this work, we propose to predict model
4 output independently and estimate the corresponding uncertainty. We present a
5 unified audio-driven disease detection framework incorporating uncertainty quan-
6 tification (UQ). This is achieved using a Dirichlet density approximation for model
7 prediction and independent kernel distance learning in feature latent space for
8 UQ. This approach requires minimum modifications to existing audio encoder
9 architectures and is extremely parameter efficient compared to k-ensemble mod-
10 els. The uncertainty-aware model improves prediction reliability by producing
11 confidence scores that closely match the accuracy values. Evaluations using the
12 largest publicly available respiratory disease datasets demonstrate the advantage of
13 the proposed framework in accuracy, training and inference time over ensemble
14 and dropout methods. The proposed model improves speech and audio analysis
15 for medical diagnosis by identifying and calibrating uncertainties, enabling better
16 decision-making and risk assessment. This is shown by high uncertainty scores at
17 low model accuracy.

18 1 Introduction

19 The increase in general awareness and interest in speech technologies for disease diagnosis has
20 generated significant growth in recorded public health datasets Song et al. (2023); Novikova and
21 Balagopalan ([n. d.]) across different modalities such as audio, imaging and time series (EEG). As
22 the healthcare industry increasingly embraces data-driven approaches, the accurate interpretation of
23 these subtle and complex multi-modal signals has become paramount for informed decision-making
24 and improved patient outcomes. However, for these models to be useful in practical implementation,
25 the outputs of such models must be explainable for medical decision making Miller (2019). Multi-
26 modal medical datasets have been extensively researched for the task of disease diagnosis, symptom
27 identification and monitoring Kulkarni et al. (2023); Wang and Wang (2022); Bae et al. (2023).
28 Popularly, large-scale convolutional neural network (CNN) architectures Demir et al. (2020) such as
29 ResNet Gairola et al. (2021); Bengs et al. ([n. d.]) trained on spectrogram images of audio inputs are
30 used for this task. Recently, direct waveform speech encoders (Wav2Vec Baevski et al. (2020), and
31 PASE Ravanelli et al. (2020)) have shown improved speech feature representations for respiratory
32 monitoring Kulkarni et al. (2023). After featurisation, a classification layer followed by softmax is
33 used to produce output scores. However, fixed softmax scores may result in fundamentally incorrect
34 outputs without indicating that the estimate is uncertain. Thus, achieving a statistically nuanced
35 understanding of model outputs via uncertainty quantification (UQ) is crucial in safety-critical
36 applications such as disease detection.

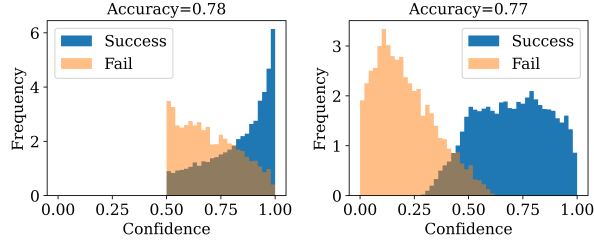


Figure 1: Calibration histograms for speech-driven COVID classifier (left) and uncertainty aware model (right) coloured according to prediction accuracy

37 This can be illustrated with a simple example of speech-driven disease detection. Input audio can
 38 be either "healthy" or COVID-"positive." A softmax-based classifier gives scores that express the
 39 likelihood of two different classes. Figure 1 (left) shows a histogram of the softmax scores coloured
 40 according to the correctness of the predicted output. The plot shows that irrespective of correctness of
 41 the prediction, the output confidence is always greater than 50%. The confidence score for two inputs
 42 (one predicted correctly and another incorrectly) lying on a vertical will be exactly same (healthy =
 43 0.89, positive = 0.11) and (healthy = 0.11, positive = 0.89). Without UQ model,
 44 there is no way to decide the reliability of either prediction based on just softmax probabilities. An
 45 independent UQ estimate can quantify high uncertainty for false predictions, as shown in Figure 1
 46 (right). An uncertainty-aware audio classification model enables 1) prediction of confidence scores
 47 independent of model outputs and 2) calibration of model such that estimated uncertainty closely
 48 follows model accuracy.

49 In this work, we present a novel framework for uncertainty-aware disease detection using speech
 50 and non-speech inputs through quantification and disentanglement of sample uncertainty and model
 51 calibration. The framework comprises of a probabilistic classification head on top of a self-supervised
 52 audio encoder and model uncertainties are quantified using a feature distance-based metric. A
 53 training scheme is proposed to optimize uncertainty estimation independent of model prediction or
 54 classification training. A novel formulation of learnable transformation matrix in latent space is used
 55 to maximise feature space diversity for distance calculation. Evaluations show that the uncertainty-
 56 aware model produces low confidence scores at low accuracy values, thus improving output reliability.
 57 Experiments on the largest public respiratory disease datasets show that the proposed UQ model
 58 is generalizable, computationally efficient at training and enables fast evaluation during inference
 59 without sacrificing classification performance. Specifically, our contributions are as follows -

- 60 • Advocate the use of a probabilistic classifier in place of softmax scores to quantify irreducible
 61 uncertainties inherent in learning problem for audio-driven disease diagnosis and medical
 62 decision making
- 63 • Emphasize the necessity of model calibration for reducible uncertainties in audio-driven
 64 disease diagnosis. we show that combining probabilistic classifier simple k-ensembles (even
 65 with small k=5) significantly improves model calibration score
- 66 • Propose a novel single inference method of uncertainty quantification with minimal changes
 67 to large encoder models for high-fidelity datasets such as audio and speech. The proposed
 68 model performs as well as k-ensembles at a fraction of compute and memory costs

69 To best of our knowledge, this is the first systematic study of uncertainties quantification and model
 70 calibration associated with audio driven disease diagnosis.

71 2 Model

72 Lets denote $a(t) \in \mathcal{A}$ as an input audio waveform and $(y_j = y_j + \epsilon_j)$ is its corresponding noisy label
 73 which takes a value from label space $j \in \{1, \dots, \mathbf{J}\}$ and ϵ_j is the label noise due to data gathering
 74 process or the noise inherent to the mapping problem $G : \mathcal{A} \rightarrow \mathbf{J}$. We decompose above function
 75 mapping as $G = h \circ f$, where, $f : \mathcal{A} \rightarrow \mathbf{R}^n$ indicates a deep audio feature encoder. The feature
 76 encoder gives embedding vectors $X_w(a) \in \mathcal{R}^d$. The uncertainty aware classification head $h : X \rightarrow y$
 77 gives a prediction over class labels $P[y|x] = h(X)$.

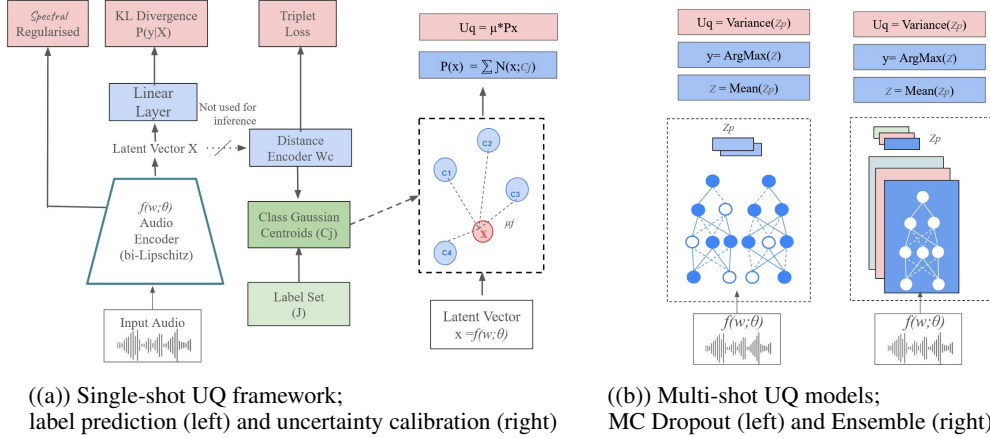


Figure 2: Proposed framework for uncertainty quantification (UQ) of audio driven disease detection

78 The proposed uncertainty quantification (UQ) framework, illustrated in Figure 2, consists of two
 79 parts:

- 80 1. A probabilistic classifier h trained to output concentration parameters of Dirichlet distribution
 81 over the softmax layers. This classifier head is used on top of a regularised deep audio
 82 feature extractor (f), which produces latent embedding X .
- 83 2. An uncertainty aware calibration training to estimate UQ as a function of feature space
 84 density. We use a novel learnable Mahalanobis distance-based metric, which ensures the
 85 latent space is bi-Lipschitz continuous and captures a measure of data distribution.

86 In the subsequent sections, we describe these two component of the proposed UQ framework

87 2.1 Probabilistic Classifier

88 A deterministic softmax classifier only outputs a single scaled vector $s(x)$ corresponding a input x
 89 such that $\sum_j s(x) = 1$. In contrast, the probabilistic classifier head is trained to predict a vector of
 90 concentration parameters $\alpha = (\alpha_1 \dots \alpha_J)$ one for each class label $j \in J$, and a strength parameter
 91 $\alpha_0 := \sum_j (\alpha_j)$. This set of concentration parameters define a Dirichlet distribution $\text{Dir}(\alpha)$ with
 92 probability density given by equation 1, where $\Gamma(\cdot)$ denotes *Gamma* function.

93 This is used to sample a class probability vector \mathbf{p} as a random vector $\mathbf{p} \sim \text{Dir}(\alpha)$, At the inference
 94 time, a sample from Dirichlet distribution gives indicative probability p_j of input x belonging to class
 95 j . The expected probability (mean) and the variance for a single input x is given by

$$\text{Dir}(\mathbf{p}|\alpha) = \frac{\Gamma(\alpha_0)}{\prod_{c=1}^C \Gamma(\alpha_c)} \prod_{c=1}^C p_c^{\alpha_c - 1} \leftrightarrow \begin{cases} \mu(x) := \mathbf{E}[p_j] = \frac{\alpha_j}{\alpha_0} \\ \sigma^2(x) := \frac{\alpha_j(\alpha_0 - \alpha_j)}{\alpha_0(\alpha_0 + 1)} \end{cases} \quad (1)$$

96 Thus, the classifier head is a model with uncertainty that outputs two quantities corresponding to label
 97 distribution, the mean $\mu(x)$ and the variance $\sigma(x)$. The sampling based output stems from key insight
 98 that softmax based classifier cannot capture output categorical probability but a distribution over
 99 categorical softmax (i.e. Dirichlet) can be used to formulate deep learning as evidence acquisition
 100 problem Sensoy et al. (2018); DeVries and Taylor (2018).

101 The classification head is trained using unweighted combination of negative log likelihood term \mathcal{L}^{NLL}
 102 and a KL-divergence term, following the Sensoy et al. (2018); Bachstein et al. (2019). Appendix
 103 covers Loss function derivations and final expressions. Upon training the classifier model using
 104 above loss, we obtain predictive distribution parameters - mean $\mu(x)$ and variance $\sigma(x)$. However
 105 this quantity only gives the output label probability of a given input for a fixed model. Considering the
 106 original function mapping problem $G : \mathcal{A} \rightarrow \mathbf{J}$ and the decomposition $G = h \circ f$, the probabilistic

107 classifier h is a single sample from a possibly large intractable hypothesis space \mathcal{H} . Further, the
 108 audio encoder f is parametrised by a set weights \mathbf{W} . In the supervised setting a point estimate
 109 of vector W is obtained by empirical risk maximisation of an objective function. In Bayesian
 110 modelling Lakshminarayanan et al. (2016); Gal and Uk (2016), uncertainties in this point estimate,
 111 are computed by assuming that the weights w follow a prior distribution $Pr(w)$. Subsequently, the
 112 model training process leads to posterior distribution $P(w|D)$. The trained model $f_w(x)$ uses this
 113 posterior distribution to calculate the estimated output y . The measure of uncertainty, UQ, is given by
 114 the expected value and variance of the prediction $f_w(x)$ over the posterior density distribution of w .
 115 However, for high-dimensional datasets such as audio and speech, accurate modelling of the density
 116 function $P(w|D)$ is impossible, given the complexity and non-linear nature of weights w of audio
 117 classification models Hernández-Lobato and Adams (2015).

118 Figure 2(b) shows two approaches for approximating the intractable posterior density $P(w|D)$ by
 119 introducing diversity into model evaluations. The feature encoder generates a fixed and deterministic
 120 encoding vector X for the input audio signal. Model uncertainty is quantified by analysing the
 121 variance of the outputs obtained through multiple forward passes of diverse models. In **Monte Carlo**
 122 **(MC) dropout** Gal and Ghahramani (2016); Xiao and Wang (2019), probabilistic ($p = 0.1$) dropout
 123 layers between non-linear layers of the network are activated during inference resulting in variable
 124 outputs. Whereas, in **Deep ensemble** Lakshminarayanan et al. (2016), k -different models ($k = 5, 15$)
 125 are trained using different subsets of the dataset. The ensemble prediction is the average soft-max
 126 outputs from the individual models.

127 Combining classification head with multi-forward pass inferences in equation 1, we get a series of
 128 means $\mu_k(x)$ and variances $\sigma_k^2(x)$, where $k \in [1, K]$ are number of different ensembles or inferences
 129 of Figure 2(b). These samples are combined to form a single predictive uncertainty estimate $Var[X]$
 for input X as an empirical expectation over all inferences k . A combination of Deep Ensemble

$$\begin{aligned} \text{Var} * (\mathbf{x}) &= \frac{1}{K} \sum_k \sigma_k^2(\mathbf{x}) + \frac{1}{K} \sum_k \mu_k^2(\mathbf{x}) - \mu_*^2(\mathbf{x}) \\ &= \mathbf{E}_k[\sigma_k^2(\mathbf{x})] + \mathbf{E}_k[\mu_k^2(\mathbf{x})] - \mathbf{E}_k[\mu_k(\mathbf{x})]^2 \\ &= \underbrace{\mathbf{E}_k[\sigma_k^2(\mathbf{x})]}_{\text{Aleatoric Uncertainty}} + \underbrace{\text{Var}_k[\mu_k(\mathbf{x})]}_{\text{Epistemic Uncertainty}} \end{aligned}$$

130 and Dirichlet Probabilistic classifier gives an estimate for the Irreducible Aleatoric Uncertainty and
 131 Model Uncertainty (Epistemic). However, it is neither possible to treat each term separately nor to
 132 reduce epistemic part of uncertainty. Despite the limitations, the k -ensemble approach is shown to
 133 be the state-of-the-art for uncertainty prediction on several benchmarks Mukhoti et al. (2021). Both
 134 these methods improve performance and uncertainty estimation through model diversity but incur
 135 high computational costs during training and inference. In next section, we describe the second part
 136 of the proposed framework - an alternative to k -ensemble for quantifying approximate Epistemic
 137 uncertainty in single forward pass.
 138

139 2.2 Single Inference Uncertainty Quantification

140 In contrast to multiple feed-forward evaluation models, we propose single-shot UQ estimation using
 141 latent feature maps produced by the encoder as a representation of the class conditional distribution.
 142 A distance measure in the feature space of the model has shown to be useful for the detection of
 143 out-of-distribution examples Venkataraman et al. (2023) and uncertainty estimation Lee et al.
 144 ([n. d.]); van Amersfoort et al. (2020). However, these methods suffer from three key problems
 145 namely feature collapse van Amersfoort et al. (2020), class imbalance Venkataraman et al. (2023),
 146 smoothness and sensitivity Lee et al. ([n. d.]). We first describe the proposed single shot approach
 147 with intuitive modifications to training scheme that address the aforementioned problems.

148 The uncertainty estimation flow is shown in Figure 2(a). A centroid vector $Z \in \mathcal{R}^m$ is initialised
 149 randomly and assigned to each label class in a set of classes J . Let $X_t(i) \in \mathcal{R}^d$ be the set of audio
 150 encodings of a mini-batch during training. A distance transformation matrix $W_j(m, d)$ is initialised
 151 using a Gaussian prior per class, where d is the feature encoding dimension and $m < d$ is the
 152 size of the centroid vector. Weight matrix W_j acts as a learnable linear dimensionality reduction
 153 on feature vectors, enabling a compact representation for distance computation Ren et al. (2021);

154 Venkataramanan et al. (2023). The class-dependent nature of W_j enables class separation in latent
 155 space and is crucial for minimising the likelihood of feature collapse.

A weighted feature distance D_j between the model output and centroids is computed as:

$$D_j(X_t, Z_j) = \sqrt{\frac{\|W_j X_t - Z_j\|^2}{2m\sigma_j^2}}$$

156 where length scale σ_j is a trainable parameter and acts as class dependent normalising hyper-
 157 parameter.

158 If the matrix W is assumed to be Identity Matrix the above formulation computes Mahalanobis
 159 distance (MD) from the centroids. The learnable nature of W acts as an adaptive dimensionality
 160 reduction on the latent space X and the output WX can be expected to represent global distributions
 161 as well as class dependent local distributions.

162 During the forward pass, a class label for each sample is given by softmax of distance scores
 163 $y_i = \text{Argmin}_j Z_j(X_i)$ as the maximum correlation (minimum distance) between data point X_i and
 164 class centroids Z_j . For the UQ estimate, the set of Mahalanobis distances is normalised through the
 165 division of maximum class distance. The model uncertainty is given by mixture of the Gaussian
 166 models fitted at each class centroid $d_{\text{UQ}} = \sum_j \mathcal{N}(D_j | z_j, \sigma_j)$.

The class centroids, Z_j , are updated for every mini-batch of training using an exponential moving
 average of the feature vectors of data points corresponding to class j :

$$Z_{t+1,j} = \gamma Z_{t,j} + \frac{1}{n_j} (1 - \gamma) \sum_i (W_j X_i)$$

167 where n_j is number of samples in the j^{th} class, and γ is a hyper-parameter similar to momentum
 168 gradient descent. After each update, the class vectors are normalised such that $\|Z_j\|_2 = 1$.

Class dependent **triplet Loss** formulation is used to maximise the distance between distinct class
 centroids and minimise intra-class separation, following Kumar et al. (2020); Hermans et al. (2017).
 Audio embeddings obtained from the encoder network were used as an anchor point X_a . Let Z_a be
 the centroid vector of the class corresponding to true label y_a , while Z_j indicates remaining centroid
 vectors such that $\{j \in \mathcal{J} \forall j \neq a\}$, The loss with margin $\epsilon \in (0.1 - 0.5)$ is given by

$$\mathcal{L}_{\text{triplet}} = \sum_{a,j} \max(\|WX_a - Z_a\| - \|WX_a - Z_j\| + \epsilon, 0)$$

169 During the training process, this loss is averaged over a mini-batch of data points, the class centroids
 170 are updated to new locations as per predicted labels and stochastic gradient descent (SGD) is
 171 performed for θ and W_j . Audio encoder output latent vectors usually have high dimensions and the
 172 above loss may suffer poorly due to involved distance computation. Low rank nature of W ensures
 173 that distance computation in above loss function is sensible.

174 **Feature Regularisation** High dimensional feature space embedding suffer from feature collapse
 175 and feature redundancy in latent space which can adversely affect uncertainty prediction Liu et al.
 176 (2020); van Amersfoort et al. (2020). These problems can be alleviated by encouraging latent space
 177 smoothness and sensitivity, or alternatively by regularising the the weights W to follow bi-Lipschitz
 178 condition Liu et al. (2020)

$$L_1 * \|x_1 - x_2\|_X \leq \|f_W(x_1) - f_W(x_2)\|_H \leq L_2 * \|x_1 - x_2\|_X$$

179 This ensures the mapping $\|f_W(x_1) - f_W(x_2)\|_H$ has meaningful correspondence in input space with
 180 respect to a well defined distance measure $\|x_1 - x_2\|_X$ Liu et al. (2020). This condition also ensures
 181 smoothness in latent space such that the audio embeddings are not too sensitive to small variations in
 182 input.

183 We use spectral normalisation to enforce bi-Lipschitz condition during UQ training, following
 184 the analysis Smith et al. (2021); Liu et al. (2020) that adding spectral normalisation before each
 185 convolution layer leads to bi-Lipschitz condition. Apart from being simpler in implementation (with
 186 minor changes to encoder architecture such as replacing L2 norm layer by spectral norm), spectral
 187 normalisation is significantly faster Smith et al. (2021) and is more stable during training compared
 188 to Jacobian Gradient penalty implemented in van Amersfoort et al. (2020).

189 3 Experiments

190 We will now demonstrate the utility of proposed framework in quantifying uncertainties of audio
191 driven disease diagnosis. We first start with a brief description of datasets, evaluation criterion and
192 implementation details. (detailed description and data histograms are covered in Appendix)

193 3.1 Datasets

194 We conduct extensive experiments using two popular audio-driven healthcare diagnosis datasets.

195 The **ICBHI** Rocha et al. (2018) dataset is the largest publicly available respiratory audio repository
196 recorded from 128 patients with a total of 6898 labelled breathing cycles (Label distribution 3642
197 normal, 1864 crackle, 886 wheeze, and 506 cycles as both). The highly unbalanced dataset constitutes
198 a 4-class audio classification task.

199 **COSWARA** Sharma et al. (2020) consists of a diverse set of manually curated audio records from
200 2635 individuals, of which 1819 are SARS-CoV-2 negative, 674 are positive subjects, and the
201 remaining unlabelled or noisy samples are filtered out. Speech recordings of numbers (1-20) counted
202 at a fast pace were used for this 2-class classification and disease detection task.

203 3.2 Self-supervised Audio Encoder

204 Self-supervised learning (SSL) is an attractive approach for healthcare audio datasets where the
205 data size is limited and manual annotation is expensive Sharma et al. (2020); Rocha et al. (2018).
206 Three different SSL models are employed as audio encoders for the empirical evaluation. First, an
207 image-based **ResNet-50** is used as the backbone with a residual block of two 3×3 convolution layers
208 and a skip connection between each block. The network is trained on the self-supervised task of
209 spectral feature prediction and reconstruction of the log-Mel spectrogram. Further, **Wav2Vec** Baevski
210 et al. (2020) and **PASE** Ravanelli et al. (2020) are used as direct waveform feature encoders. Each
211 encoder is pre-trained on the respective SSL pretext task and used to obtain latent representations
212 from raw audio.

213 Let $a(t) \in \mathcal{A}$ be an input audio waveform and $y = 1, \dots, J$ be its corresponding label. The feature
214 encoder gives embedding vectors $X_w(a) \in \mathcal{R}^d$, where $d = 256$ is the fixed latent dimension.

215 3.3 Preprocessing

216 All audio files were resampled to a fixed rate of 22.05kHz. The ICBHI respiratory sounds were
217 cropped/padded to max a length of 7s Gairola et al. (2021); Kulkarni et al. (2023), while COSWARA
218 speech were fixed to 10s length Sharma et al. (2020). In the case of ResNet, each audio was
219 transformed to log Mel-spectrogram using 128 frequency bins. An input size of (128, 350) was used
220 for ICBHI, whereas, for COSWARA, the input size was (128, 500). For both cases, the dataset was
221 divided into three non-overlapping portions such that the test set (20%) and validation set (20%)
222 contained audio records from different patients than that of the train set (60%).

223 3.4 Evaluation

224 For measuring accuracy of model, **sensitivity** ($\frac{TP}{TP+FN}$), and **specificity** ($\frac{FP}{FP+TN}$) scores were
225 used. Each score measures class-wise prediction accuracy in the case of the unbalanced dataset. The
226 notations TN, FN denote true and false negative rates and TP, FP denote true and false positive
227 rates, respectively. Average of these two scores ($\frac{SP+SN}{2}$) was used for comparison with SoTA
228 models Rocha et al. (2018). The area under the receiver operating curve (**AUROC**) was used as an
229 indicative probability of correctly classifying a randomly selected unseen sample.

230 Most common measure predictive uncertainty is Expected Calibration error (**ECE**). Low ECE
231 indicates model accuracy closely follows predicted uncertainty estimates, i.e. low model accuracy in
232 high-uncertainty regions and vice versa. To calculate ECE on a test set, all test samples are grouped
233 in $k = 10$ equal bins according to uncertainty scores. ECE was calculated as the absolute sum of
234 differences between expected model confidence and accuracy for each bin. A small ECE indicates
235 better performance as the model accurately quantifies uncertainties in its prediction. Experiments

236 show that ECE values drastically reduce with the proposed UQ implementation while maintaining
 237 the model’s accuracy.

238 4 Results and Discussion

239 The first goal of the experiments is to answer the question ‘whether model uncertainty score follows
 240 model accuracy’. Figure 7 show reliability diagrams of ICBHI 4-class classification model using
 241 PASE encoder as backbone. The model output is divided in equally spaced bins according to estimated
 242 confidence score for each bin. Reliability plots show the average accuracy of the examples in each
 243 corresponding confidence bin. We also visualise the confidence scores (1- uncertainty) with class
 244 conditional histograms of correctly and incorrectly classified outputs. The proposed model reliably
 245 predicts high uncertainty misclassified examples while producing high uncertainty for accurately
 classified examples.

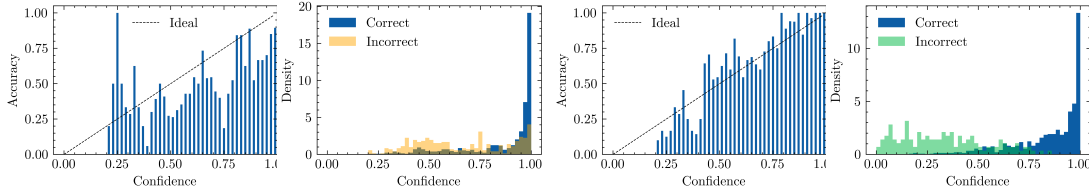


Figure 3: Reliability diagrams before and after feature distance based uncertainty calibration. Plots show that proposed models predicts UQ scores that closely follow the model accuracy. (low confidence scores for low accuracy data regions and vice versa)

246

Table 1: Performance comparison of different base encoder models with and without uncertainty estimation (ICBHI)

Model	Base (Dirichlet)		Base+UQ	
	AUROC	ECE	AUROC	ECE
PASE	0.835±0.01	0.121±0.01	0.905±0.01	0.055±0.01
Wav2Vec	0.778±0.02	0.148±0.01	0.812±0.02	0.069±0.01
ResNet	0.746±0.01	0.106±0.01	0.862±0.01	0.041±0.01

247 A similar analysis is conducted for different choices of base feature encoder (Table 1) by considering
 248 the ECE (error) and AUROC (accuracy) of ICBHI respiratory classification task using different audio
 249 encoders (ResNet, PASE Ravanelli et al. (2020) and Wav2Vec Baeovski et al. (2020)) with and without
 250 UQ estimation. A significant reduction in ECE values is observed among all three feature models.
 251 This means the model is more uncertain for false predictions and more confident for correct outputs.
 252 ResNet achieves higher relative improvement compared to direct waveform-based audio encoders.
 253 This is due to ResNet having higher embedding dimension compared to SSL encoders and thus
 254 adversely affecting the class conditional density estimation in latent space Ren et al. (2021). The
 255 low rank class-wise linear transformation enables distribution aware low dimensional transformation,
 256 improving both AUROC and ECE score.

257 **Classification Accuracy** and dataset variability of uncertainty aware models are compared in Table
 258 2. Bootstrapping is used to compute the maximum confidence interval. The proposed model shows
 259 a significant advantage in ECE prediction over other UQ methods with marginal improvements in
 260 model accuracy.

261 **An ablation study** was conducted to study the incremental effects of various loss functions by
 262 fixing the feature encoder of the proposed UQ model. Table 3a displays the ECE and accuracy
 263 improvements with each additional loss term. A significant reduction in ECE error is observed upon
 264 the inclusion of triplet loss term for both datasets.

265 **Compute efficiency** of the proposed method, in terms of the number of parameters (in Millions) and
 266 inference time (in milliseconds), is compared with those of popular UQ models in Table 3b. The
 267 scores show the expected inference time for a single sample averaged over the test set compared

Table 2: Evaluation of the UQ framework for two different datasets with fixed feature encoder (PASE)

Model	ECE	SN(%)	SP(%)	AUROC
ICBHI₄-class				
Base	0.161 \pm 0.01	79.8 \pm 4.71	50.5 \pm 6.21	0.782 \pm 0.01
MC Drop.	0.064 \pm 0.01	79.6 \pm 5.31	42.6 \pm 5.91	0.732 \pm 0.02
Ensemble	0.051 \pm 0.01	83.1 \pm 3.71	57.7 \pm 1.91	0.888 \pm 0.01
Our (UQ)	0.045\pm0.01	82.1\pm4.07	55.1\pm3.75	0.823\pm0.01
COSWARA₂-class				
Base	0.191 \pm 0.02	96 \pm 3.32	72.9 \pm 2.21	0.781 \pm 0.01
MC Drop.	0.074 \pm 0.01	96 \pm 5.59	70 \pm 4.19	0.951 \pm 0.01
Ensemble	0.060 \pm 0.01	96.6 \pm 3.15	77.9 \pm 4.98	0.964 \pm 0.01
Our (UQ)	0.058\pm0.01	95.9\pm4.81	74.6\pm2.91	0.961\pm0.01

Table 3: Ablation study (a) of proposed UQ framework to study effects of modification terms, along with network size (Millions) and inference time (sec) of different UQ models (Results on non-intersecting splits of ICBHI dataset with PASE as feature encoder)

(a) Ablation study			(b) Network size			
Model	ECE	AUROC	Method	AUROC	Params	Inference
Old (Softmax)	0.158 \pm 0.01	0.741 \pm 0.02	Base (logits)	0.782	26M	1.8 ms
Base (Dirichlet)	0.149 \pm 0.01	0.876 \pm 0.01	MC Dropout	0.732	26M	4.3 ms
+ KL Divergence	0.104 \pm 0.01	0.921 \pm 0.02	Ensemble - 5	0.888	132M	9.8 ms
+ Triplet loss	0.086 \pm 0.01	0.923 \pm 0.02	Ensemble - 15	0.891	395M	29 ms
+ Regularisation	0.065 \pm 0.01	0.918 \pm 0.01	Mahalanobis	0.823	26M	2.1 ms

268 against AUROC scores. In this case, PASE is used as the base model. The ensemble model performed
 269 well but was extremely slow at inference time with a large number of parameters, increasing the
 270 storage and compute overhead. The Mahalanobis distance-based uncertainty estimation enables
 271 lightweight and fast inference while improving model accuracy.

Table 4: Comparison with SoTA models and recent studies on four-class respiratory anomaly detection (ICBHI dataset)

Method	Performance			
	SN(%)	SP(%)	Acc.	
ResNet Gairola et al. (2021)	40.1	72.3	56.2	
ResNeST Wang and Wang (2022)	70.4	40.2	55.3	
CNN8-Pt Ren et al. (2022)	72.9	27.8	50.4	
ResNet Chang et al. (2022)	69.9	35.8	52.9	
CVAE-Tr Bae et al. (2023)	81.7	43.1	62.4	
Our (UQ) ECE-	0.058	82.1\pm4.07	55.1\pm3.75	68.5\pm3.92

272 **Comparison** with state-of-the-art (SoTA) models for ICBHI 4+class respiratory sound classification
 273 task is presented in Table 4. The proposed model improved the accuracy scores over the current
 274 SoTA by 6.1%. A validation set sensitivity score of 82.1% indicates the ability to correctly identify
 275 true positives from unseen patient samples recorded using different digital stethoscopes. Accounting
 276 for the uncertainties not only provides a nuanced understanding of output but also improves model
 277 performances for audio-driven disease diagnosis.

278 5 Uncertainty Visualisation and Decision Making

279 The outputs produced by sampling from Dirichlet distribution (output of probabilistic classifier for a
 280 single input) satisfy the property that $\sum_j(p_j) = 1$, where p_j is probability $P[y = j|X]$. For a three
 281 class problem (ICBHI - wheeze, crackle, healthy), each of these samples fall on the 2D plane defined
 282 by $\sum_j(p_j) = 1$. Figure 4 shows uncertainty visualisations on the simplex plane. This uncertainty

283 is sample specific (data/ aleatoric uncertainty) indicating inherent label noise or ambiguity in the samples.

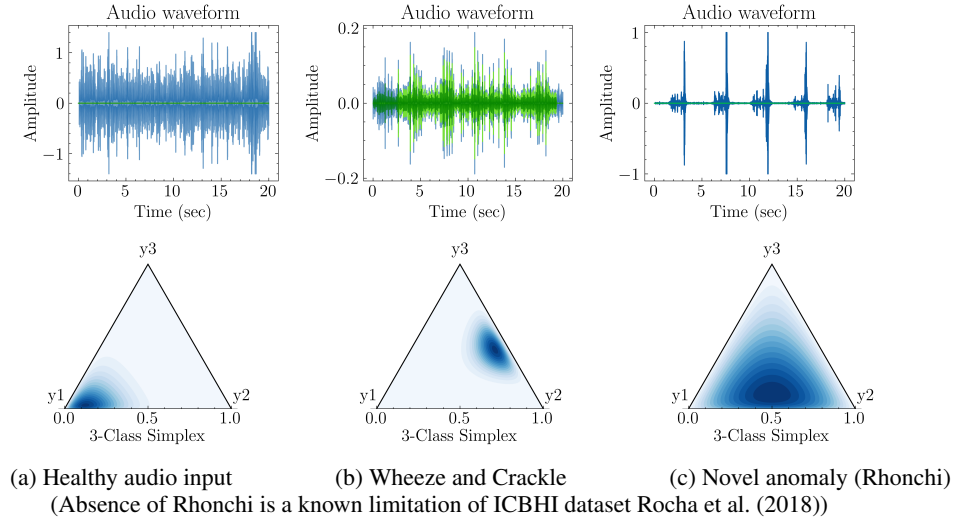


Figure 4: Plots visualising data uncertainty corresponding to each input audio sample. The network predicts Dirichlet distribution parameters (α) in single forward pass which are then used to plot probability density over the simplex.

284

285 At the same time, the model uncertainty (predictive) is given by mixture of the Gaussian models fitted
 286 at each class centroid $d_{UQ} = \sum_j \mathcal{N}(D_j | C_j, \sigma_j)$. This estimate is independent model prediction
 287 at a given sample. This is a measure of learning capacity of model for current input and can also
 288 be used as OOD indicator (epistemic uncertainty). A threshold on the UQ score can be used as a
 289 decision factor for audio-driven medical diagnosis. If the predicted UQ value is higher than this
 290 threshold, the model is not sufficiently confident in its prediction; thus, the disease diagnosis output
 291 is rejected. In such cases, second or multiple evaluations using re-recording of input audio samples
 292 are recommended. If the resulting uncertainty, after multiple empirical evaluations, is still higher
 293 than the threshold, then the particular sample is selected for clinical or manual diagnosis. This avoids
 294 the risk of erroneous predictions via uncertainty quantification. As a result, the proposed framework
 295 improves the performance of audio-driven disease detection system along with patient safety. (Such
 296 threshold based rejection was not used during experiments and results, however it can be a useful
 297 tool for medical decision making)

298 6 Conclusion

299 In this work, a framework for uncertainty-aware disease diagnosis was proposed using speech and
 300 non-speech inputs. The UQ framework enables confidence scoring to improve the reliability of
 301 model outputs. Evaluations of the popular COSWARA and ICBHI datasets illustrate the superiority of
 302 the proposed model over the popular ensemble and Monte Carlo dropout method. Using the same
 303 ResNet backbone, the UQ aware model outperformed softmax-based SoTA models for respiratory
 304 disease Bae et al. (2023) without using data driven oversampling techniques. Using the UQ model for
 305 the ICBHI dataset, an improvement of 6.1% was observed over the SoTA models. Furthermore, for
 306 speech-driven COVID detection, quantifying data uncertainty improves AUROC scores by 18.1%.
 307 The UQ model performs well on unseen datasets, as seen from results on non-intersecting inter-patient
 308 data splits, and is equally applicable to more general datasets. Results also show the effectiveness and
 309 applicability of the Mahalanobis distance-based metric for different general-purpose audio encoders.
 310 Finally, the proposed framework enables fast and lightweight UQ estimation, making it more suitable
 311 for implementation in mobile and IoT devices for continuous health monitoring owing to its small
 312 size and lower number of trainable parameters.

313 **References**

- 314 Simon Bachstein, Gutachter Prof, Markus Pauly, and Florian Wilhelm. 2019. Uncertainty Quantifica-
315 tion in Deep Learning.
- 316 Sangmin Bae, June-Woo Kim, Won-Yang Cho, Hyerim Baek, Soyoun Son, Byungjo Lee, Changwan
317 Ha, Kyongpil Tae, Sungnyun Kim, and Se-Young Yun. 2023. Patch-Mix Contrastive Learning with
318 Audio Spectrogram Transformer on Respiratory Sound Classification. In *Proc. INTERSPEECH*
319 *2023*. 5436–5440. <https://doi.org/10.21437/Interspeech.2023-1426>
- 320 Alexei Baevski, Steffen Schneider, and Michael Auli. 2020. vq-wav2vec: Self-Supervised Learning
321 of Discrete Speech Representations. <http://arxiv.org/abs/1910.05453> arXiv:1910.05453
322 [cs].
- 323 Viktor Bengs, Willem Waegeman, and Willem Waegeman@ugent Be. [n. d.]. Pitfalls of Epistemic
324 Uncertainty Quantification through Loss Minimisation.
- 325 Yi Chang, Zhao Ren, Thanh Tam Nguyen, Wolfgang Nejdl, and Björn W. Schuller. 2022. Example-
326 based Explanations with Adversarial Attacks for Respiratory Sound Analysis. In *Proc. INTER-*
327 *SPEECH 2022*. arXiv:2203.16141 [cs.SD]
- 328 Siddhartha Dalal and Vishal Misra. 2024. The Matrix: A Bayesian learning model for LLMs. (2
329 2024). <http://arxiv.org/abs/2402.03175>
- 330 Fatih Demir, Abdulkadir Sengur, and Varun Bajaj. 2020. Convolutional neural networks based
331 efficient approach for classification of lung diseases. *Health Information Science and Systems* 8, 1
332 (Dec. 2020), 4. <https://doi.org/10.1007/s13755-019-0091-3>
- 333 Terrance DeVries and Graham W. Taylor. 2018. Learning Confidence for Out-of-Distribution
334 Detection in Neural Networks. (2 2018). <http://arxiv.org/abs/1802.04865>
- 335 Sina Däubener, Lea Schönherr, Asja Fischer, and Dorothea Kolossa. 2020. Detecting Adversarial
336 Examples for Speech Recognition via Uncertainty Quantification. arXiv:2005.14611 [eess.AS]
- 337 Siddhartha Gairola, Francis Tom, Nipun Kwatra, and Mohit Jain. 2021. RespireNet: A Deep Neural
338 Network for Accurately Detecting Abnormal Lung Sounds in Limited Data Setting. <http://arxiv.org/abs/2011.00196> arXiv:2011.00196 [cs, eess].
- 340 Yarín Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian Approximation: Representing
341 Model Uncertainty in Deep Learning. In *Proceedings of The 33rd International Conference on*
342 *Machine Learning (Proceedings of Machine Learning Research, Vol. 48)*, Maria Florina Balcan
343 and Kilian Q. Weinberger (Eds.). PMLR, New York, New York, USA, 1050–1059. <https://proceedings.mlr.press/v48/gal16.html>
- 345 Yarín Gal and Zg201@cam Ac Uk. 2016. Dropout as a Bayesian Approximation: Representing
346 Model Uncertainty in Deep Learning Zoubin Ghahramani. <http://yarin.co>.
- 347 Alexander Hermans, Lucas Beyer, and Bastian Leibe. 2017. In Defense of the Triplet Loss for Person
348 Re-Identification. *ArXiv (2017)*. arXiv:1703.07737 [cs.CV]
- 349 José Miguel Hernández-Lobato and Ryan P. Adams. 2015. Probabilistic Backpropagation for Scalable
350 Learning of Bayesian Neural Networks. (2 2015). <http://arxiv.org/abs/1502.05336>
- 351 Dae Y. Kang, Pamela N. DeYoung, Justin Tantiengloc, Todd P. Coleman, and Robert L. Owens. 2021.
352 Statistical uncertainty quantification to augment clinical decision support: a first implementation
353 in sleep medicine. *npj Digital Medicine* 4 (12 2021). Issue 1. <https://doi.org/10.1038/s41746-021-00515-3>
- 355 Shubham Kulkarni, Hideaki Watanabe, and Fuminori Homma. 2023. Self-Supervised Audio Encoder
356 with Contrastive Pretraining for Respiratory Anomaly Detection. In *2023 IEEE International*
357 *Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)*. 1–5. <https://doi.org/10.1109/ICASSPW59220.2023.10193030>
- 358

- 359 Puneet Kumar, Sidharth Jain, Balasubramanian Raman, Partha Pratim Roy, and Masakazu Iwamura.
360 2020. End-to-end Triplet Loss based Emotion Embedding System for Speech Emotion Recognition.
361 In *2020 25th International Conference on Pattern Recognition (ICPR)*. arXiv:2010.06200 [cs.SD]
- 362 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2016. Simple and Scalable
363 Predictive Uncertainty Estimation using Deep Ensembles. (12 2016). [http://arxiv.org/abs/
364 1612.01474](http://arxiv.org/abs/1612.01474)
- 365 Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. [n. d.]. A Simple Unified Framework for
366 Detecting Out-of-Distribution Samples and Adversarial Attacks.
- 367 Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with Confidence: Uncertainty
368 Quantification for Black-box Large Language Models. arXiv:2305.19187 [cs.CL]
- 369 Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshmi-
370 narayanan. 2020. Simple and Principled Uncertainty Estimation with Deterministic Deep Learning
371 via Distance Awareness. (6 2020). <http://arxiv.org/abs/2006.10108>
- 372 Simon W. McKnight, Aidan O. T. Hogg, Vincent W. Neo, and Patrick A. Naylor. 2023. Uncertainty
373 Quantification in Machine Learning for Joint Speaker Diarization and Identification. (12 2023).
374 <http://arxiv.org/abs/2312.16763>
- 375 Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial
376 Intelligence* 267 (2019), 1–38. <https://doi.org/10.1016/j.artint.2018.07.007>
- 377 Jishnu Mukhoti, Andreas Kirsch, Joost van Amersfoort, Philip H. S. Torr, and Yarin Gal. 2021. Deep
378 Deterministic Uncertainty: A Simple Baseline. (2 2021). <http://arxiv.org/abs/2102.11582>
- 379 Radford M. Neal. 1995. Bayesian Learning for Neural Networks. [https://api.
380 semanticscholar.org/CorpusID:60809283](https://api.semanticscholar.org/CorpusID:60809283)
- 381 Jekaterina Novikova and Aparna Balagopalan. [n. d.]. On Speech Datasets in Machine Learning for
382 Healthcare. <https://www.researchgate.net/publication/338593562>
- 383 Mirco Ravanelli, Jianyuan Zhong, Santiago Pascual, Pawel Swietojanski, Joao Monteiro, Jan Trmal,
384 and Yoshua Bengio. 2020. Multi-task self-supervised learning for Robust Speech Recognition.
385 <http://arxiv.org/abs/2001.09239> arXiv:2001.09239 [cs, eess].
- 386 Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshmi-
387 narayanan. 2021. A Simple Fix to Mahalanobis Distance for Improving Near-OOD Detection.
388 arXiv:2106.09022 [cs.LG]
- 389 Zhao Ren, Thanh Tam Nguyen, Wolfgang Nejdl, and Wolfgang Nejdl. 2022. Prototype Learning for
390 Interpretable Respiratory Sound Analysis. In *ICASSP 2022 - 2022 IEEE International Conference
391 on Acoustics, Speech and Signal Processing (ICASSP)*. 9087–9091. [https://doi.org/10.
392 1109/ICASSP43922.2022.9747014](https://doi.org/10.1109/ICASSP43922.2022.9747014)
- 393 B. M. Rocha, D. Filos, P. Carvalho, and N. Maglaveras. 2018. Respiratory Sound Database for the
394 Development of Automated Classification. In *Springer Precision Medicine*. Vol. 66. Springer
395 Singapore, Singapore, 33–37. https://doi.org/10.1007/978-981-10-7419-6_6
- 396 Murat Sensoy, Lance Kaplan, and Melih Kandemir. 2018. Evidential Deep Learning to Quantify
397 Classification Uncertainty. In *Advances in Neural Information Processing Systems (Neurips 2018,
398 Vol. 31)*. Curran Associates, Inc. [https://proceedings.neurips.cc/paper_files/paper/
399 2018/file/a981f2b708044d6fb4a71a1463242520-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/a981f2b708044d6fb4a71a1463242520-Paper.pdf)
- 400 Neeraj Sharma, Prashant Krishnan, Rohit Kumar, Shreyas Ramoji, Srikanth Raj Chetupalli, Nirmala
401 R., Prasanta Kumar Ghosh, and Sriram Ganapathy. 2020. Coswara - A Database of Breathing,
402 Cough, and Voice Sounds for COVID-19 Diagnosis. In *Interspeech 2020 (interspeech 2020)*. ISCA.
403 <https://doi.org/10.21437/interspeech.2020-2768>
- 404 Lewis Smith, Joost van Amersfoort, Haiwen Huang, Stephen Roberts, and Yarin Gal. 2021. Can
405 convolutional ResNets approximately preserve input distances? A frequency analysis perspective.
406 arXiv:2106.02469 [cs.LG] <https://arxiv.org/abs/2106.02469>

- 407 Meishu Song, Andreas Triantafyllopoulos, Zhonghao Zhao, Kun Qian, Zijiang Yang, Hiroki Takeuchi,
408 Toru Nakamura, Akifumi Kishi, Tetsuro Ishizawa, Kazuhiro Yoshiuchi, Xin Jing, Vincent Karas,
409 Bin Hu, Björn W Schuller, and Yoshiharu Yamamoto. 2023. Daily Mental Health Monitoring
410 From Speech: A Real-World Japanese Dataset and Multitask Learning Analysis. <https://www.researchgate.net/publication/370561148>
411
- 412 Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. 2020. Uncertainty Estimation
413 Using a Single Deep Deterministic Neural Network. (3 2020). <http://arxiv.org/abs/2003.02037>
414
- 415 Aishwarya Venkataramanan, Assia Benbihi, Martin Laviale, and Cedric Pradalier. 2023. Gaussian
416 Latent Representations for Uncertainty Estimation using Mahalanobis Distance in Deep Classifiers.
417 [arXiv:2305.13849](https://arxiv.org/abs/2305.13849) [cs.CV]
- 418 Zijie Wang and Zhao Wang. 2022. A Domain Transfer Based Data Augmentation Method for
419 Automated Respiratory Classification. In *ICASSP 2022 - 2022 IEEE International Conference*
420 *on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, Singapore, Singapore, 9017–9021.
421 <https://doi.org/10.1109/ICASSP43922.2022.9746941>
- 422 Tong Xia, Jing Han, Lorena Qendro, Ting Dang, and Cecilia Mascolo. 2021. Uncertainty-Aware
423 COVID-19 Detection from Imbalanced Sound Data. *ArXiv* abs/2104.02005 (2021). <https://api.semanticscholar.org/CorpusID:233025205>
424
- 425 Yijun Xiao and William Yang Wang. 2019. Quantifying uncertainties in natural language processing
426 tasks. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence* (Honolulu,
427 Hawaii, USA) (*AAAI'19/IAAI'19/EAAI'19*). AAAI Press, Article 899, 8 pages. <https://doi.org/10.1609/aaai.v33i01.33017322>
428

429 **Appendix / supplemental material**

430 **A Algorithm**

431 We consider the probabilistic function learning problem between input audio space $X \in \mathcal{X}$ and
 432 corresponding discrete label space $Y \in \mathcal{Y}$, where X and Y are random variables. We denote $x \in \mathcal{X}$
 433 and $y \in \mathcal{Y}$ as data samples from joint space $\mu_{XY} = (X, Y)$, with a joint distribution function denoted
 434 by μ_{XY} . A model trained with uncertainty aware classification tries to approximate the conditional
 435 probability distribution $\mu_{Y|X} = \mathbf{P}[Y = y|X]$.

436 The training dataset $D_n = \{(x_i, y_i) \forall i = 1 \dots N\}$ (a subset of joint space μ_{XY}) is used to train an
 437 estimator $\hat{Y} = g(X; W)$, where g denotes a neural network with parameter W . Further, the collected
 438 dataset itself can be inherently noisy or error-prone, which results in data uncertainty (also known as
 439 Aleatoric uncertainty). The noise in the dataset is indicated by $x_i = \hat{x} + \eta_i$ and $y_i = \hat{y} + \epsilon_j$ where the
 440 the observed noisy dataset is given by, $D_n = \{(x_i, y_i)\}$. This aleatoric uncertainty is irreducible and
 441 can only be estimated as expected variance in the output for a fixed input X , and a given estimator
 442 $f(x|W)$. However is not the only source of uncertainty in the estimator, the output variance does not
 443 capture the uncertainty in estimator or the learning process itself.

444 The total predictive probability, can be expanded as follows -

$$\mu_{Y|X} = \mathbf{P}[Y = y|X] \tag{2}$$

$$= \mathbf{P}[Y = y|X, D_n]P[X \in D_n] \tag{3}$$

$$.. + \mathbf{P}[Y = y|X, D_n]P[X \notin D_n] \tag{4}$$

$$\tag{5}$$

445 The second term in above equation signifies distribution uncertainty, i.e. uncertainty associated with
 446 limitations of training data. This can be reduced by obtaining more training data i.e. by minimising
 447 $P[X \notin D_n]$. Assuming the dataset D is used to learn a function $y = f(x; W)$, parametrized by the
 448 weights W , the first term can further be expanded as follows -

$$\mu_{Y|X,D} = \mathbf{P}[Y = y|X, D_n] \tag{6}$$

$$= \int P(y|X, w)dP(w, D) \tag{7}$$

$$= \int P(y|X, w)P(w|D)dw \tag{8}$$

$$\tag{9}$$

449 This integral is called as inference using posterior density $P[w|D]$, this computation involves test
 450 time optimisation, by formulating closed form of posterior density. The second term, posterior in
 451 above equation can be decomposed as

$$P(w|D) = \frac{\mathbf{P}[D|W]P[W]}{P[W]} \tag{10}$$

$$P(D) = \int P(D|w)dP(w) \tag{11}$$

$$= \int P(D|w)P(w)dw \tag{12}$$

$$\tag{13}$$

452 This integral is called as marginal integration to compute a form for posterior density from a presumed
 453 prior $p[w]$. Often this marginal is intractable for most non trivial forms of likelihood functions $p[D|W]$.
 454 The inference integral is often approximated using multiple forward pass via Dropout or Ensemble
 455 modelling. The goal of the proposed distance based model is to provide an efficient single forward
 456 pass alternative to approximate the marginal and inference integrals.

457 Given that we can view an ensemble member as a single deterministic model and vice versa, this
 458 provides an intuitive explanation for why single deterministic models report inconsistent and widely
 459 varying predictive entropies and confidence scores for OoD samples for which a Deep Ensemble
 460 would report high epistemic uncertainty (expected information gain) and high predictive entropy.

461 Assuming that $p(y|x, \omega)$ only depends on $p(y|x)$ and $\mathbb{I}[Y; w|x]$, we model the distribution of $p(y|x, \omega)$
 462 (as a function of ω) using a Dirichlet distribution $Dir(\alpha)$ which satisfies:

$$p(y|x) = \frac{\alpha_i}{\alpha_0} \quad (14)$$

$$H[Y|x] - \mathbb{I}[Y; w|x] = \psi(\alpha_0 + 1) \quad (15)$$

$$(16)$$

463 Then, we can model the softmax distribution using a random variable $\mathbf{p} \sim Dir(\alpha)$ as:

$$P(y|x; w) \approx Cat(\mathbf{p}). \quad (17)$$

464 The variance $VarH[Y|x; w]$ of the softmax entropy for different samples x given $p(y|x)$ and
 465 $\mathbb{I}[Y; w|x]$ is then approximated by $VarY|\mathbf{p}$: This is the estimate of Aleatoric Uncertainty in the
 466 model. For the random variable, $\mathbf{p} \sim Dir(\alpha)$, the expected entropy $\mathbb{E}_{\mathbf{p} \sim Dir(\alpha)} \mathbb{H}_{Y \sim Cat(\mathbf{p})}[Y]$ of the
 467 categorical distribution $Y \sim Cat(\mathbf{p})$ is given by

$$\mathbb{E}_{\mathbf{p}(\mathbf{p}|\alpha)} \mathbb{H}[Y | \mathbf{p}] = \psi(\alpha_0 + 1) - \sum_{y=1}^K \frac{\alpha_i}{\alpha_0} \psi(\alpha_i + 1)$$

468 Proof. Applying the sum rule of expectations and 3 from 1.1 we can write

$$\begin{aligned} \mathbb{E} \mathbb{H}[Y | \mathbf{p}] &= \mathbb{E} \left[- \sum_{i=1}^K \mathbf{p}_i \log \mathbf{p}_i \right] = - \sum_i \mathbb{E} [\mathbf{p}_i \log \mathbf{p}_i] \\ &= - \sum_i \frac{\alpha_i}{\alpha_0} (\psi(\alpha_i + 1) - \psi(\alpha_0 + 1)) \end{aligned}$$

469 The result follows after rearranging and making use of $\sum_i \frac{\alpha_i}{\alpha_0} = 1$.

470 B Base Model

471 B.1 Feature Encoder

472 A feature encoder serves as base model (backbone) of the framework. The feature encoder serves as
 473 an indicative audio classification backbone. The proposed framework can accommodate any state of
 474 the art audio encoder and does not require any modification in the feature encoder training process
 475 and architecture. Specifically, when $a(t) \in \mathcal{A}$ be an input audio waveform, the feature encoder gives
 476 embedding vectors $X_w(a) \in \mathcal{R}^d$, where $d = 256$ is the fixed latent dimension. Next we explain the
 base encoders used for experiments

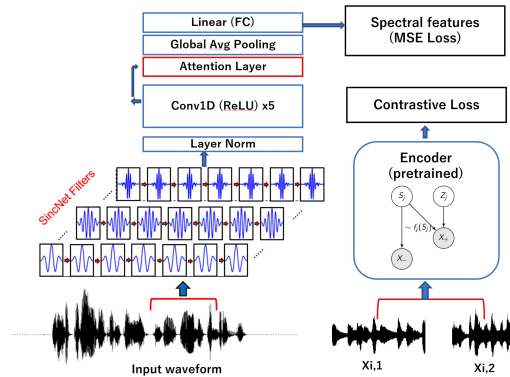


Figure 5: Self supervised feature encoder architecture for PASE+ Ravanelli et al. (2020)

478 B.2 Wav2vec2.0

479 wav2vec is a self-supervised learning model trained to learn representations of raw audio waveforms
480 directly, without relying on manual transcriptions or labels. Wav2vec employs contrastive learning
481 to learn powerful representations from raw audio inputs. A more recent version, wav2vec 2.0
482 introduces a more sophisticated approach by masking portions of the latent space rather than the
483 raw audio. Wav2Vec 2.0 significantly improves upon the quality of learned representations and
484 demonstrates exceptional performance in downstream speech tasks. We use Wav2vec2.0 as one of
485 the backbone feature extractor in the proposed framework. The Wav2vec model predicts the masked
486 latent representations, encouraging it to capture rich contextual information. The output embedding
487 dimension of the Wav2vec encoder is fixed to be 256.

488 B.3 PASE +

489 Problem Agnostic Speech Encoder (PASE) is another self supervised audio feature encoder which
490 employs multiple neural networks, termed "workers," to tackle various self-supervised tasks. These
491 workers contribute to learning rich and discriminative representations. To ensure robust feature
492 vectors with respect to small variations in input audio, PASE+ introduces an online speech distortion
493 module that artificially corrupts the input audio, forcing the encoder to learn more invariant features.
494 As shown in Figure 5 PASE+ also uses bidirectional attention layers to combine convolution outputs
495 to better capture both short-term and long-term speech dynamics.

496 B.4 ResNet

497 An image-based **ResNet-50** is used as the backbone with a residual block of two 3×3 convolution
498 layers and a skip connection between each block. The network is trained on the self-supervised task
499 of spectral feature prediction and reconstruction of the log-Mel spectrogram. The network consists of
500 a series of convolution layers. Each of these layers is defined with 64 channels, kernel strides (5, 2, 2,
501 2, 2, 2, 2), and kernel widths (7, 3, 3, 3, 3, 3, 2, 2), respectively, followed by batch normalization
502 and ReLU activation. The interval between two sequential samples in the feature encoder output Z is
503 15ms, and the receptive audio field is 20 ms. The output from convolution layers is concatenated and
504 passed to a multi-head attention layer and a fully connected layer with an embedding size of 256.
505 Like PASE Ravanelli et al. (2020), the final linear layer is used to predict speech features such as log
506 power spectrum (LPS), MFCCs, prosody, 40 FBANKS and 40 Gammatone features. The architecture
507 is pre-trained on an open source audio dataset called Audioset ?, consisting of a wide variety of input
508 sounds ranging such as birds, coughs, speech and machine sounds. During pretraining, the model
509 predicts a set of 12 supervised tasks consisting of regression and binary feature banks such as log
510 power spectrum (LPS), MFCCs, prosody and Gammatone features. This pretraining ensures that the
511 ResNet representations are tuned capture short and long-range audio dynamics over a wide variety of
512 input sounds. These representations are proven to outperform spectrogram-based large CNN models
513 and standard acoustic features for different classification and speech recognition tasks Ravanelli et al.
514 (2020). These representations are then frozen to compute encoding for respiratory cycle datasets.
515 Experiments show that no significant improvement are observed with additional complete finetuning
516 on the ICBHI dataset during the training phase compared to the frozen representation.

517 C Criteria

518 C.1 Probabilistic Classifier

519 We train the classification using unweighted combination of negative log likelihood term \mathcal{L}^{NLL} and
520 a KL-divergence term, following the Sensoy et al. (2018); Bachstein et al. (2019). Appendix covers
521 Loss function derivations and final expressions.

522 The loss function expressions of \mathcal{L}^{NLL} and \mathcal{L}^{KL} are respectively

$$\mathcal{L}^{NLL} = \sum_{c=1}^C y_c (\log(\alpha_0) - \log(\alpha_c)) \quad (18)$$

$$\begin{aligned} \mathcal{L}^{KL} &= \log \left(\frac{\Gamma(\sum_{c=1}^C \tilde{\alpha}_c)}{\Gamma(C) \prod_{c=1}^C \Gamma(\tilde{\alpha}_c)} \right) \\ &+ \sum_{c=1}^C (\tilde{\alpha}_c - 1) \left(\psi(\tilde{\alpha}_c) - \psi \left(\sum_{c=1}^C \tilde{\alpha}_c \right) \right) \end{aligned} \quad (19)$$

523 in which $\tilde{\alpha}_c = y_c + (1 - y_c)\alpha_c$ and $\psi(\cdot)$ is *Digamma* function.

524 These two losses can viewed intuitively as a union of **Bayes Risk Approximation** losses, which is
525 defined with respect to class conditional density prediction. We use Bayes risk formulation from PAC
526 learning nomenclature as given below,

$$\mathcal{L}_i(\Theta) = \sum_{j=1}^K (y_{ij} - \mathbf{E}[p_{ij}])^2 + \text{Var}(p_{ij}) \quad (20)$$

$$= \sum_{j=1}^K \underbrace{(y_{ij} - \alpha_{ij}/S_i)^2}_{\mathcal{L}_{ij}^{\text{err}}} + \underbrace{\frac{\alpha_{ij}(S_i - \alpha_{ij})}{S_i^2(S_i + 1)}}_{\mathcal{L}_{ij}^{\text{var}}} \quad (21)$$

$$= \sum_{j=1}^K (y_{ij} - \hat{p}_{ij})^2 + \frac{\hat{p}_{ij}(1 - \hat{p}_{ij})}{(S_i + 1)}. \quad (22)$$

527 C.2 Uncertainty Calibration Network

528 During the forward pass, a class label for each sample is given by softmax of distance scores
529 $y_i = \mathbf{Argmin} Z_j(X_i)$ as the maximum correlation (minimum distance) between data point X_i and
530 class centroids Z_j . For the UQ estimate, the set of Mahalanobis distances is normalised through the
531 division of maximum class distance. The model uncertainty is given by mixture of the Gaussian
532 models fitted at each class centroid $d_{\text{UQ}} = \sum_j \mathcal{N}(D_j | z_j, \sigma_j)$.

The class centroids, Z_j , are updated for every mini-batch of training using an exponential moving average of the feature vectors of data points corresponding to class j :

$$Z_{t+1,j} = \gamma Z_{t,j} + \frac{1}{n_j} (1 - \gamma) \sum_i (W_j X_i)$$

533 where n_j is number of samples in the j^{th} class, and γ is a hyper-parameter similar to momentum
534 gradient descent. After each update, the class vectors are normalised such that $\|Z_j\|_2 = 1$.

Class dependent **triplet Loss** formulation is used to maximise the distance between distinct class centroids and minimise intra-class separation, following Kumar et al. (2020); Hermans et al. (2017). Audio embeddings obtained from the encoder network were used as an anchor point X_a . Let Z_a be the centroid vector of the class corresponding to true label y_a , while Z_j indicates remaining centroid vectors such that $\{j \in \mathcal{J} \forall j \neq a\}$. The loss with margin $\epsilon \in (0.1 - 0.5)$ is given by

$$\mathcal{L}_{\text{triplet}} = \sum_{a,j} \mathbf{max} (\|W X_a - Z_a\| - \|W X_a - Z_j\| + \epsilon, 0)$$

535 During the training process, this loss is averaged over a mini-batch of data points, the class centroids
536 are updated to new locations as per predicted labels and stochastic gradient descent (SGD) is
537 performed for θ and W_j .

538 **C.3 Evaluation**

539 For measuring accuracy of model, **sensitivity** ($\frac{TP}{TP+FN}$), and **specificity** ($\frac{FP}{FP+TN}$) scores were
 540 used. Each score measures class-wise prediction accuracy in the case of the unbalanced dataset. The
 541 notations TN, FN denote true and false negative rates and TP, FP denote true and false positive
 542 rates, respectively. Average of these two scores ($\frac{SP+SN}{2}$) was used for comparison with SoTA
 543 models Rocha et al. (2018). The area under the receiver operating curve (**AUROC**) was used as an
 544 indicative probability of correctly classifying a randomly selected unseen sample.

Most common measure predictive uncertainty is Expected Calibration error (**ECE**). Low ECE indicates model accuracy closely follows predicted uncertainty estimates, i.e. low model accuracy in high-uncertainty regions and vice versa. At high thresholds, the model is tolerant of low confidence predictions, and thus, the model accuracy should decrease. At low uncertainty thresholds, the model should have high accuracy and confidence scores. To calculate ECE on a test set, all test samples are grouped in $k = 10$ equal bins according to uncertainty scores. ECE was calculated as the absolute sum of differences between expected model confidence and accuracy for each bin. A small ECE indicates better performance as the model accurately quantifies uncertainties in its prediction. Experiments show that ECE values drastically reduce with the proposed UQ implementation while maintaining the model’s accuracy. The expected difference between

$$ECE = \sum_{k=1}^{10} \frac{n_k}{n} |ascore(B_k) - uscore(B_k)|$$

545 where n_k is number of samples in k^{th} bin, *ascore* and *uscore* are average accuracy and uncertainty
 546 estimates for each bin B_k . Experiments show that calibration error drastically reduces with the
 547 addition of UQ models while maintaining the model accuracy in AUROC scores. It can be interpreted
 548 as the probability that a positive example (in-distribution) will have a higher detection score than a
 549 negative example (out-of-distribution).

550 **D Datasets**

551 We conduct extensive experiments using two popular audio-driven healthcare diagnosis datasets.

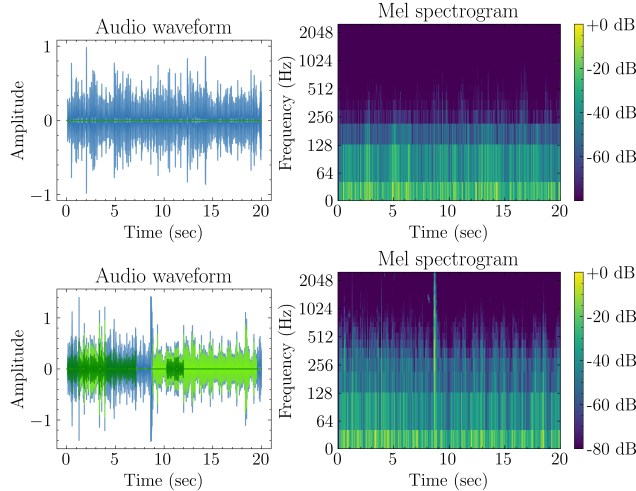


Figure 6: Audio samples showing varying degrees of anomalous (green) and healthy (blue) classes illustrating the necessity of uncertainty quantification

552 The **ICBHI** Rocha et al. (2018) dataset is the largest publicly available respiratory audio repository
 553 recorded from 128 patients with a total of 6898 labelled breathing cycles (Label distribution 3642
 554 normal, 1864 crackle, 886 wheeze, and 506 cycles as both). The highly unbalanced dataset constitutes
 555 a 4-class audio classification task. Figure 6 shows audio samples showing varying degrees of
 556 anomalous (orange) and healthy (blue) classes. The input sample contains illustrating the necessity of
 557 uncertainty quantification. We share training and validation sets of this dataset for SoTA comparison.

558 **COSWARA** Sharma et al. (2020) consists of a diverse set of manually curated audio records from
 559 2635 individuals, of which 1819 are SARS-CoV-2 negative, 674 are positive subjects, and the
 560 remaining unlabelled or noisy samples are filtered out. Speech recordings of numbers (1-20) counted
 561 at a fast pace were used for this 2-class classification and disease detection task. The dataset is
 562 manually curated and has approximately 10% noisy audio samples.

563 All audio files were resampled to a fixed rate of 22.05kHz. The ICBHI respiratory sounds were
 564 cropped/padded to max a length of 7s Gairola et al. (2021); Kulkarni et al. (2023), while COSWARA
 565 speech were fixed to 10s length Sharma et al. (2020). In the case of ResNet, each audio was
 566 transformed to log Mel-spectrogram using 128 frequency bins. An input size of (128, 350) was used
 567 for ICBHI, whereas, for COSWARA, the input size was (128, 500). For both cases, the dataset was
 568 divided into three non-overlapping portions such that the test set (20%) and validation set (20%)
 569 contained audio records from different patients than that of the train set (60%).¹

570 E Experiments

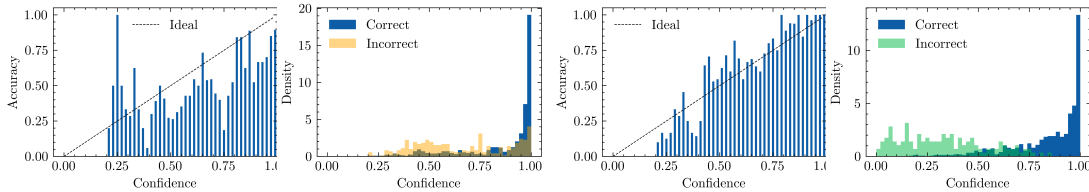


Figure 7: Reliability diagrams before and after feature distance based uncertainty calibration. Plots show that proposed models predicts UQ scores that closely follow the model accuracy. (low confidence scores for low accuracy data regions and vice versa)

571 The proposed framework is trained independently in two stages. The distance transformation matrix
 572 W and audio feature encoders were optimised during the first stage of training process. It is important
 573 to note that the goal of feature encoder training not to represent state-of-the-art for any particular
 574 task – the goal is to demonstrate value of quantifying model uncertainty independent of the model
 575 prediction. We will show that across various of of-the-shelf audio feature encoders, the addition of
 576 UQ framework enables significant gains in model utility by not only quantifying model confidence
 577 but also reducing the calibration error of the model. This point is reinforced here using 2D synthetic
 578 dataset. In second stage of training the probabilistic classifier is optimised using KL divergence loss.
 579 In this second stage we show that, using off-the-shelf encoders it is possible to achieve and state of
 580 the art performance on popular disease diagnosis task.

581 E.1 UQ on 2D dataset

In the proposed uncertainty quantification framework weighted feature distance D_j between the model output and centroids is computed as:

$$D_j(X_t, Z_j) = \sqrt{\frac{\|W_j X_t - Z_j\|^2}{2m\sigma_j^2}}$$

582 where length scale σ_j is a trainable parameter and acts as class dependent normalising hyper-
 583 parameter.

584 If the matrix W is assumed to be Identity Matrix the above formulation computes Mahalanobis
 585 distance (MD) from the centroids. The learnable nature of W acts as an adaptive dimensionality
 586 reduction on the latent space X and the output WX can be expected to represent global distributions
 587 as well as class dependent local distributions.

588 Figure 8 shows comparison of uncertainty estimates obtained using distance based metric and
 589 ensemble based model. In contrast to multiple feed forward evaluation models, a single shot
 590 estimation of distance function in feature space gives an approximation of class conditional density.

¹Training and validation set labels shared as supplementary material.

591 A Mahalanobis distance metric Lee et al. ([n. d.]); Venkataramanan et al. (2023) between output and
 592 the class centroids is has already been shown to act as an approximation of class conditional density
 593 and outperform empirical ensemble models for the task of OOD detection.

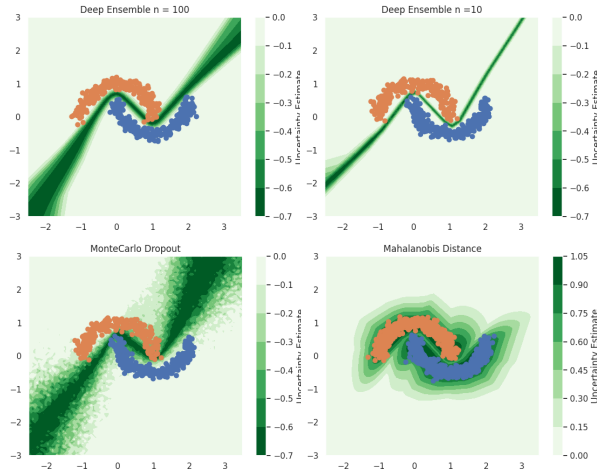


Figure 8: Comparison of proposed UQ model with popular Bayesian methods using confidence heat maps (Green) for a 2D synthetic dataset

594 F Related Work

595 Much research has been devoted to UQ in deep learning models for computer vision (CV) DeVries
 596 and Taylor (2018); Mukhoti et al. (2021) speech and language processing Lin et al. (2023); Däubener
 597 et al. (2020). Traditional UQ models, such as Bayesian neural networks Dalal and Misra (2024); Neal
 598 (1995), Monte Carlo dropout Gal and Ghahramani (2016); Xiao and Wang (2019), and deep ensemble
 599 Lakshminarayanan et al. (2016) models have been popular for speech classification and automatic
 600 speech recognition (ASR) Däubener et al. (2020). Kalman filtering with Monte Carlo dropout has
 601 been used to quantify data and model uncertainties in speaker identification McKnight et al. (2023).
 602 In medical diagnosis, a nonparametric model for UQ has been used as a noise metric for Diffusion
 603 MRI. Statistical UQ using ensemble models has also been used to augment clinical decision support
 604 in medicine Kang et al. (2021). These models are resource-intensive due to multiple training runs
 605 to form an ensemble and/or several feed-forward evaluations for a single inference. However, the
 606 explainability via uncertainty quantification of speech and audio-driven disease classification remains
 607 under-explored and is paramount for system reliability and patient safety Xia et al. (2021). In this
 608 work we treat model prediction and uncertainty quantification as independent tasks. We emphasise
 609 quantification of both reducible epistemic and aleatoric (irreducible) uncertainties using a single
 610 inference model with minimal modifications to backbone architecture.

611 G Further Comments

612 The proposed framework introduced a new way to measure how confident the model is in its
 613 predictions. This is called model uncertainty. Instead of just giving a single answer, the model
 614 provides a range of possible outcomes along with how likely each outcome is. This helps in making
 615 more reliable decisions, especially in medical applications.

616 When the model is very uncertain about its prediction, it can be considered as a "flag" to say, "I'm not
 617 sure about this." In this case, the system can ask for more information, like re-recording the audio,
 618 before making a final decision. This helps prevent incorrect diagnoses.

619 The framework was tested on two well-known datasets: COSWARA and ICBHI. It performed better
 620 than other popular methods, especially in detecting respiratory diseases. The model was able to
 621 correctly identify diseases more often and was less likely to make mistakes. One of the strengths of
 622 this framework is its ability to handle different types of audio data. It can work well with various

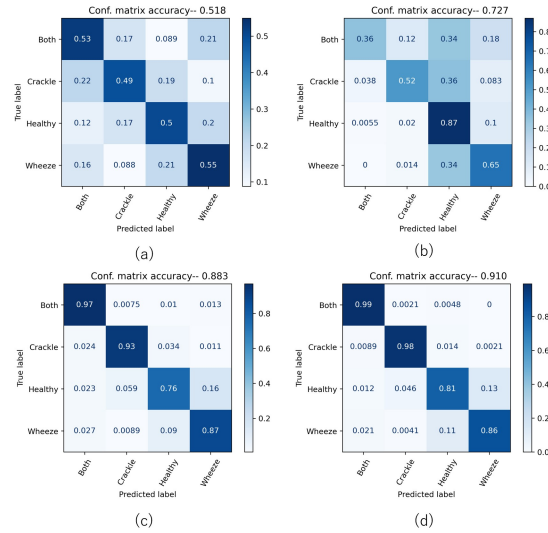


Figure 9: Evaluation of individual anomalous class performance

623 audio encoders and can be used for different diseases. This makes it a versatile tool for many medical
 624 applications. Practical Implications

625 In contrast to ensemble based SoTA alternatives, the proposed UQ model is relatively small and fast,
 626 it can be used on mobile devices or other devices with limited computing power. This means it could
 627 be used for real-time monitoring of people’s health.

628 In summary, the model offers a significant improvement in audio-based disease detection. It is more
 629 accurate, reliable, and practical than existing methods. By considering the model’s uncertainty, it
 630 helps to reduce the risk of incorrect diagnoses and improve patient safety.

631 Finally, Our results show that: there is a necessity of quantification and identification of prediction
 632 uncertainties in deep learning models for audio-driven disease estimation. Further it is necessary to
 633 distinguish between learning aleatoric and epistemic uncertainty, which is unexpected and violates
 634 assumptions on simple distribution based uncertainty quantification methods. We expect that our
 635 formulation and results help practitioners and researchers choose uncertainty methods and expand
 636 the use of disentangled uncertainties, as well as motivate additional research into this topic.

637 **NeurIPS Paper Checklist**

638 **1. Claims**

639 Question: Do the main claims made in the abstract and introduction accurately reflect the
640 paper's contributions and scope?

641 Answer: [\[Yes\]](#)

642 Justification:

643 Guidelines:

- 644 • The answer NA means that the abstract and introduction do not include the claims
645 made in the paper.
- 646 • The abstract and/or introduction should clearly state the claims made, including the
647 contributions made in the paper and important assumptions and limitations. A No or
648 NA answer to this question will not be perceived well by the reviewers.
- 649 • The claims made should match theoretical and experimental results, and reflect how
650 much the results can be expected to generalize to other settings.
- 651 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
652 are not attained by the paper.

653 **2. Limitations**

654 Question: Does the paper discuss the limitations of the work performed by the authors?

655 Answer: [\[Yes\]](#)

656 Justification: We discuss detailed limitations and related works in appendix.

657 Guidelines:

- 658 • The answer NA means that the paper has no limitation while the answer No means that
659 the paper has limitations, but those are not discussed in the paper.
- 660 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 661 • The paper should point out any strong assumptions and how robust the results are to
662 violations of these assumptions (e.g., independence assumptions, noiseless settings,
663 model well-specification, asymptotic approximations only holding locally). The authors
664 should reflect on how these assumptions might be violated in practice and what the
665 implications would be.
- 666 • The authors should reflect on the scope of the claims made, e.g., if the approach was
667 only tested on a few datasets or with a few runs. In general, empirical results often
668 depend on implicit assumptions, which should be articulated.
- 669 • The authors should reflect on the factors that influence the performance of the approach.
670 For example, a facial recognition algorithm may perform poorly when image resolution
671 is low or images are taken in low lighting. Or a speech-to-text system might not be
672 used reliably to provide closed captions for online lectures because it fails to handle
673 technical jargon.
- 674 • The authors should discuss the computational efficiency of the proposed algorithms
675 and how they scale with dataset size.
- 676 • If applicable, the authors should discuss possible limitations of their approach to
677 address problems of privacy and fairness.
- 678 • While the authors might fear that complete honesty about limitations might be used by
679 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
680 limitations that aren't acknowledged in the paper. The authors should use their best
681 judgment and recognize that individual actions in favor of transparency play an impor-
682 tant role in developing norms that preserve the integrity of the community. Reviewers
683 will be specifically instructed to not penalize honesty concerning limitations.

684 **3. Theory Assumptions and Proofs**

685 Question: For each theoretical result, does the paper provide the full set of assumptions and
686 a complete (and correct) proof?

687 Answer: [\[Yes\]](#)

688 Justification: Supplementary material covers complete proof for Dirichlet Uncertainty
689 Quantification

690 Guidelines:

- 691 • The answer NA means that the paper does not include theoretical results.
- 692 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
693 referenced.
- 694 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 695 • The proofs can either appear in the main paper or the supplemental material, but if
696 they appear in the supplemental material, the authors are encouraged to provide a short
697 proof sketch to provide intuition.
- 698 • Inversely, any informal proof provided in the core of the paper should be complemented
699 by formal proofs provided in appendix or supplemental material.
- 700 • Theorems and Lemmas that the proof relies upon should be properly referenced.

701 4. Experimental Result Reproducibility

702 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
703 perimental results of the paper to the extent that it affects the main claims and/or conclusions
704 of the paper (regardless of whether the code and data are provided or not)?

705 Answer: [Yes]

706 Justification:

707 Guidelines:

- 708 • The answer NA means that the paper does not include experiments.
- 709 • If the paper includes experiments, a No answer to this question will not be perceived
710 well by the reviewers: Making the paper reproducible is important, regardless of
711 whether the code and data are provided or not.
- 712 • If the contribution is a dataset and/or model, the authors should describe the steps taken
713 to make their results reproducible or verifiable.
- 714 • Depending on the contribution, reproducibility can be accomplished in various ways.
715 For example, if the contribution is a novel architecture, describing the architecture fully
716 might suffice, or if the contribution is a specific model and empirical evaluation, it may
717 be necessary to either make it possible for others to replicate the model with the same
718 dataset, or provide access to the model. In general, releasing code and data is often
719 one good way to accomplish this, but reproducibility can also be provided via detailed
720 instructions for how to replicate the results, access to a hosted model (e.g., in the case
721 of a large language model), releasing of a model checkpoint, or other means that are
722 appropriate to the research performed.
- 723 • While NeurIPS does not require releasing code, the conference does require all submis-
724 sions to provide some reasonable avenue for reproducibility, which may depend on the
725 nature of the contribution. For example
 - 726 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
727 to reproduce that algorithm.
 - 728 (b) If the contribution is primarily a new model architecture, the paper should describe
729 the architecture clearly and fully.
 - 730 (c) If the contribution is a new model (e.g., a large language model), then there should
731 either be a way to access this model for reproducing the results or a way to reproduce
732 the model (e.g., with an open-source dataset or instructions for how to construct
733 the dataset).
 - 734 (d) We recognize that reproducibility may be tricky in some cases, in which case
735 authors are welcome to describe the particular way they provide for reproducibility.
736 In the case of closed-source models, it may be that access to the model is limited in
737 some way (e.g., to registered users), but it should be possible for other researchers
738 to have some path to reproducing or verifying the results.

739 5. Open access to data and code

740 Question: Does the paper provide open access to the data and code, with sufficient instruc-
741 tions to faithfully reproduce the main experimental results, as described in supplemental
742 material?

743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794

Answer: [No]

Justification: We use open source healthcare datasets and provide appropriate links for data download, no data is released along with this paper. Further, we provide step-by-step guidance to reproduce experiments conducted in this study.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Supplementary material further describes validation splits in details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We use 10-fold CV for reporting all numerical figures in the results

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 795
- 796
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

805 8. Experiments Compute Resources

806 Question: For each experiment, does the paper provide sufficient information on the com-
807 puter resources (type of compute workers, memory, time of execution) needed to reproduce
808 the experiments?

809 Answer: [Yes]

810 Justification: Covered in Supplementary Material

811 Guidelines:

- 812
- 813
- 814
- 815
- 816
- 817
- 818
- 819
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

820 9. Code Of Ethics

821 Question: Does the research conducted in the paper conform, in every respect, with the
822 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

823 Answer: [Yes]

824 Justification:

825 Guidelines:

- 826
- 827
- 828
- 829
- 830
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

831 10. Broader Impacts

832 Question: Does the paper discuss both potential positive societal impacts and negative
833 societal impacts of the work performed?

834 Answer: [Yes]

835 Justification:

836 Guidelines:

- 837
- 838
- 839
- 840
- 841
- 842
- 843
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- 844
- 845
- 846
- 847
- 848
- 849
- 850
- 851
- 852
- 853
- 854
- 855
- 856
- 857
- 858
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

859 11. Safeguards

860 Question: Does the paper describe safeguards that have been put in place for responsible
861 release of data or models that have a high risk for misuse (e.g., pretrained language models,
862 image generators, or scraped datasets)?

863 Answer: [NA]

864 Justification:

865 Guidelines:

- 866
- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 875
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

876 12. Licenses for existing assets

877 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
878 the paper, properly credited and are the license and terms of use explicitly mentioned and
879 properly respected?

880 Answer: [NA]

881 Justification:

882 Guidelines:

- 883
- 884
- 885
- 886
- 887
- 888
- 889
- 890
- 891
- 892
- 893
- 894
- 895
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

896 • If this information is not available online, the authors are encouraged to reach out to
897 the asset’s creators.

898 **13. New Assets**

899 Question: Are new assets introduced in the paper well documented and is the documentation
900 provided alongside the assets?

901 Answer: [NA]

902 Justification:

903 Guidelines:

- 904 • The answer NA means that the paper does not release new assets.
- 905 • Researchers should communicate the details of the dataset/code/model as part of their
906 submissions via structured templates. This includes details about training, license,
907 limitations, etc.
- 908 • The paper should discuss whether and how consent was obtained from people whose
909 asset is used.
- 910 • At submission time, remember to anonymize your assets (if applicable). You can either
911 create an anonymized URL or include an anonymized zip file.

912 **14. Crowdsourcing and Research with Human Subjects**

913 Question: For crowdsourcing experiments and research with human subjects, does the paper
914 include the full text of instructions given to participants and screenshots, if applicable, as
915 well as details about compensation (if any)?

916 Answer: [NA]

917 Justification:

918 Guidelines:

- 919 • The answer NA means that the paper does not involve crowdsourcing nor research with
920 human subjects.
- 921 • Including this information in the supplemental material is fine, but if the main contribu-
922 tion of the paper involves human subjects, then as much detail as possible should be
923 included in the main paper.
- 924 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
925 or other labor should be paid at least the minimum wage in the country of the data
926 collector.

927 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
928 Subjects**

929 Question: Does the paper describe potential risks incurred by study participants, whether
930 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
931 approvals (or an equivalent approval/review based on the requirements of your country or
932 institution) were obtained?

933 Answer: [NA]

934 Justification:

935 Guidelines:

- 936 • The answer NA means that the paper does not involve crowdsourcing nor research with
937 human subjects.
- 938 • Depending on the country in which research is conducted, IRB approval (or equivalent)
939 may be required for any human subjects research. If you obtained IRB approval, you
940 should clearly state this in the paper.
- 941 • We recognize that the procedures for this may vary significantly between institutions
942 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
943 guidelines for their institution.
- 944 • For initial submissions, do not include any information that would break anonymity (if
945 applicable), such as the institution conducting the review.