RespiraMFM: A Multimodal Foundation Model with Contrastive Audio-Language Alignment for Respiratory Disease Identification

Anonymous ACL submission

Abstract

Respiratory diseases remain a leading cause 002 of global mortality, where timely and accurate diagnosis is critical to improving patient out-004 comes and reducing healthcare burdens. While prior work has explored audio-based models for 006 respiratory disease detection, such unimodal approaches often suffer from limited gener-007 alizability and diagnostic precision. In this paper, we propose RespiraMFM, a Multimodal Foundation Model that integrates res-011 piratory sounds with patient medical history and symptoms to enhance diagnostic accuracy 012 and disease detection capabilities. We introduce an effective contrastive alignment strategy for audio-text multimodal integration, allowing the model to learn better cross-modal representations between respiratory sounds and 017 corresponding textual clinical information. We 019 evaluate RespiraMFM across five major respiratory diseases using seven real-world datasets in both supervised fine-tuning and zero-shot settings, achieving a 14.4% improvement in AU-ROC on supervised tasks and a 13.4% gain on zero-shot tasks over existing baselines. These findings underscore the potential of our framework to advance early diagnosis and improve 027 clinical decision-making in respiratory disease management.

1 Introduction

037

041

Respiratory diseases, such as COVID-19, tuberculosis (TB), chronic obstructive pulmonary disease (COPD), asthma, and pneumonia, remain a leading cause of morbidity and mortality worldwide (Weinberger et al., 2020). Most existing works (Baur et al., 2024; Zhang et al., 2024a) on detecting those respiratory diseases rely solely on audio inputs such as coughing sounds or stethoscope recordings. However, their performance is often constrained by the limited information that audio data alone can provide.

To mitigate the constraints of relying solely

on audio inputs, researchers have proposed multimodal methods (Kim et al., 2024; Zhang et al., 2024b) that combine respiratory audio with relevant clinical information, such as symptoms (e.g., fever, fatigue, chest pain) and lifestyle factors such as smoking history. For instance, BTS (Kim et al., 2024) proposes an audio-text model that uses metadata associated with respiratory sounds. This model concatenates learned representation from audio and text encoders to improve respiratory disease identification. Moreover, RespLLM (Zhang et al., 2024b) utilizes a large language model (LLM) as the text encoder alongside a separate audio encoder, with a linear projector for matching dimensions. Despite these efforts, existing methods still face two major limitations: (1) they lack effective fusion mechanisms for deep cross-modal interaction, as they often rely on simple concatenation or linear projection to combine modalities, which fails to capture complex semantic relationships between audio and text features. and (2) their zero-shot performance reveals limited generalization capability, leading to suboptimal performance on novel or unseen diseases. Since the audio and text encoders are typically trained independently, their resulting representations may not be well-aligned, hindering the model's ability to exploit complementary cross-modal information.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

To address these challenges, we introduce RespiraMFM, a Multimodal Foundation Model that effectively bridges audio and textual representations for respiratory disease identification. Specifically, we improve existing multimodal systems with a two-stage learning framework. First, we use a contrastive learning-based module to align the representation from both audio and text modalities. This process facilitates the second stage, where LLMs can effectively leverage the aligned audiotext representation, instead of modality-specific features, to accurately identify the respiratory disease.

We evaluate RespiraMFM using seven real-

world datasets that cover five of the most common respiratory diseases: COVID-19, TB, COPD, 084 asthma, and pneumonia. We highlight four of our findings: (1) RespiraMFM consistently outperforms state-of-the-art multimodal baselines on respiratory disease identification by achieving a 14.4% improvement in AUROC on supervised tasks and a 13.4% improvement on zero-shot tasks. (2) RespiraMFM achieves superior generalization ca-091 pabilities and effectively detects unseen respiratory diseases without requiring any training samples of those diseases. (3) RespiraMFM significantly reduces the training data requirement, achieving comparable performance with an order of magnitude less training data compared to the baselines. (4) Our contrastive alignment module effectively unifies audio and text modalities, leading to consistent AUROC improvements across all tasks compared 100 to models without this module. 101

2 Related Work

102

103

104

105

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

128

129

130

131

132

2.1 Single-Modal Models

The majority of existing respiratory disease identification methods (Yang et al., 2020; Ma et al., 2020; Chang et al., 2022) rely solely on audio inputs such as cough sounds or stethoscope recordings. Bae et al. (2023) introduce a contrastive learning framework to enhance respiratory sound classification using Audio Spectrogram Transformer (AST). By mixing spectrogram patches generated from raw audio data and applying contrastive loss, the model learns robust and discriminative features, which are subsequently passed to a linear classifier for respiratory disease identification. OPERA (Zhang et al., 2024a) curates large-scale unlabeled respiratory audio datasets and pretrains three foundational models using self-supervised learning. Among them, OPERA-CT, the best-performing model, is a contrastive learning-based transformer model, which is used as general-purpose feature extractors for respiratory disease classification tasks. HeAR (Health Acoustic Representations) (Baur et al., 2024) introduces a self-supervised generative learning-based framework trained on a large corpus of healthrelated audio data. By leveraging generative objectives during pretraining, HeAR learns generalizable audio representations which is utilized for downstream disease diagnosis tasks via simple linear probes. Despite the promising results of singlemodal models, their performance is limited by the information available from audio data alone.

2.2 Multimodal Models

Unlike single-modal models, multimodal models 134 combine audio data with textual information such 135 as patient symptoms and medical history, leading 136 to more accurate diagnoses. BTS (Kim et al., 2024) 137 introduces a text-audio model that combines res-138 piratory sounds with metadata transformed into 139 descriptive text. It uses the Contrastive Language-140 Audio Pretraining (CLAP) (Elizalde et al., 2023) 141 model to extract features from both modalities, fol-142 lowed by a linear classifier for respiratory disease 143 classification. However, the use of a basic linear 144 classifier limits its ability to generalize in zero-shot 145 scenarios or when encountering new diseases. To 146 date, RespLLM (Zhang et al., 2024b) is one of the 147 early efforts that applies a multimodal LLM frame-148 work integrating text and audio representations for 149 respiratory disease prediction. Their approach uti-150 lizes a pretrained encoder to extract audio and text features and a trainable linear projector to align the 152 feature dimensions with LLM. However, since each 153 modality encoder is trained separately, the result-154 ing representations are often distinct and may not 155 be directly compatible across modalities. While a 156 linear projector can align the encoder output dimen-157 sions with those expected by the LLM, it does not 158 ensure semantic alignment between modalities. To 159 address these limitations, we propose a contrastive 160 alignment module that facilitates more effective in-161 tegration by aligning audio and text representations 162 in a shared semantic space. Our approach goes 163 beyond mere dimensional alignment, aiming to es-164 tablish a shared representation space that enables 165 effective integration of multimodal information. 166

133

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

3 RespiraMFM

3.1 Overview

Figure 1 provides an overview of the proposed RespiraMFM framework. In the data curation stage shown in Figure 1(a), given the multimodal respiratory datasets, we extract and pre-process the raw audio data and the corresponding patient symptoms to construct the instruction tuning data for respiratory disease identification. As shown in Figure 1(b), the curated multimodal data are first processed by the audio and text encoders to get audio and text representations, respectively. One key component of our framework is the alignment module that reduces the domain mismatches between audio features and the language model embeddings. The alignment module is trained separately via con-



(a) Data Curation

(b) Multimodal foundation model for respiratory disease identification.



183trastive learning. Upon completion of training, the184alignment module is frozen and incorporated into185the instruction tuning stage. During instruction186tuning, the curated data are passed through each187encoder to obtain modality-specific representations,188which are then fused by concatenating them. The189resulting multimodal representation is subsequently190fed into the LLM to generate predictions for respi-191ratory disease classification.

3.2 Data Curation

The multimodal respiratory datasets consist of both respiratory audio recordings and the correspond-194 ing patient-reported symptoms in either JSON or 195 tabular format. The objective of data curation is to 196 create the instruction tuning data by pre-processing 197 the respiratory audio recordings, generating in-198 struction prompts, and converting patient-reported symptoms into structured textual representations. For audio recordings, each recording was normalized to 8 seconds in length by either truncating longer recordings or padding shorter ones through repetition. The audio signals are then processed with a 64ms Hann window with a 32ms step size, and subsequently converted into mel spectrograms denoted as x_a using the features extracted from 207 a pre-trained OPERA-CT encoder (Zhang et al., 2024a). Patient metadata varies across datasets in terms of structure and format. As shown in Fig-210 ure 1(a), we select relevant symptoms (Table 5) 211 from the tabular data and apply a standardized template to generate a textual representation x_c . We 213 utilize task-specific prompts x_p such as - "Classify 214 whether the participant has COVID-19 given the 215 following information. The 2 classes are: healthy, 216 COVID19. Please output 0 for healthy and 1 for 217



Figure 2: Illustration of contrastive learning-based audio-text alignment.

COVID19" to guide the LLM in producing disease classification outputs.

218

219

220

221

222

223

224

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

3.3 Contrastive Learning-based Audio-Text Aligning

We introduce a contrastive learning-based audiotext alignment module to align audio features with language model embeddings effectively. Specifically, we employ a pre-trained audio encoder that generates 768-dimensional embeddings from the input audio. In contrast, large language models typically operate with higher-dimensional input embeddings (e.g., 2048 or 4096, depending on the model). Therefore, it is necessary to match the audio dimension to use as input into the LLM. Prior work (Zhang et al., 2024b) addresses this by introducing a trainable linear projector to map audio embeddings into the higher-dimensional space required by the LLM. However, given the fundamental differences between the audio and text encoders in both architecture and representational semantics, simple dimensional alignment may be insufficient to achieve effective multimodal fusion (Lyu et al., 2023, 2024). To enable more effective cross-modal alignment between audio and text, we adopt a con-

trastive learning strategy—an approach shown to 242 yield powerful multimodal representations in mod-243 els like CLIP (Radford et al., 2021). As shown in 244 Figure 2, we compute text embeddings $\mathbf{e}_t \in \mathbb{R}^d$ using a frozen LLaMA model, where $\mathbf{e}_t = f_T(x_c)$, 246 f_T is the text encoder (LLM), x_c is the textual con-247 text, and d is the embedding dimension of LLM. Similarly, audio embeddings $\mathbf{e}_a \in \mathbb{R}^{768}$ are extracted using a frozen OPERA encoder, where $\mathbf{e}_a = f_O(x_a), f_O$ is the pre-trained Opera-CT au-251 dio encoder, and x_a is the mel-spectrogram of the raw audio data.

> A lightweight projection head $f_{\theta} : \mathbb{R}^{768} \to \mathbb{R}^d$ is trained to map audio embeddings into the same semantic space as the text embeddings. The training objective minimizes a contrastive loss that encourages matched audio-text pairs to be close while pushing unmatched pairs apart.

Formally, for a batch of N paired samples, we define the normalized embeddings as:

$$\mathbf{z}_i^a = \frac{f_{\theta}(\mathbf{e}_i^a)}{\|f_{\theta}(\mathbf{e}_i^a)\|}, \quad \mathbf{z}_i^t = \frac{\mathbf{e}_i^t}{\|\mathbf{e}_i^t\|}.$$

The contrastive loss (Chen et al., 2020) is given by:

$$\mathcal{L}_{\text{contrast}} = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(\mathbf{z}_{i}^{a} \cdot \mathbf{z}_{i}^{t}/\tau)}{\sum_{j=1}^{N} \exp(\mathbf{z}_{i}^{a} \cdot \mathbf{z}_{j}^{t}/\tau)},$$

where τ is a temperature scaling factor. By training this projection head with contrastive supervision, we achieve better semantic alignment across modalities while keeping the audio and text encoders frozen. The model architecture and additional training details about the aligner module are presented in Appendix D.

3.4 Instruction Tuning

260

261

262

263

264

269

271

272

275

277

We employ instruction tuning to guide the LLM in understanding and following task-specific prompts that connect the multi-modal input and the corresponding diagnostic outcomes. The core components of instruction tuning are described below.

278Multimodal Fusion: The audio and text features279are fused at the embedding level by concatena-280tion. During this stage, we utilize the contrastively281trained alignment module (f_{θ}) from the previous282step (§3.3), keeping its weights frozen to preserve283the learned representations. Similarly, the text en-284coder (f_T) and the audio encoder (f_O) are also285kept frozen during this stage. The inputs include286mel-spectrogram (x_a) , curated patient symptom

descriptions as contextual information (x_c) , and a task-specific prompt (x_p) . Audio embeddings $z_a \in \mathbb{R}^d$ are extracted via the audio encoder and subsequently projected to match the input dimensionality of the LLM.

$$z_a = f_\theta(f_O(x_a)) \tag{29}$$

287

289

290

293

294

297

299

300

301

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

325

326

327

328

329

Simultaneously, the prompt and contextual text are processed through the LLM's encoder to obtain their respective representations, denoted as $z_p \in \mathbb{R}^d$ and $z_c \in \mathbb{R}^d$, corresponding to the prompt and context embeddings.

$$z_p = f_T(x_p), \quad z_c = f_T(x_c).$$
 290

Finally, we concatenate the audio (z_a) , prompt (z_p) , and context (z_c) embeddings to get a combined embedding of a longer sequence:

$$z_{fusion} = z_a \parallel z_p \parallel z_c \tag{30}$$

where \parallel denotes concatenation operation and $z_{fusion} \in \mathbb{R}^d$.

Large Language Model: We utilize Bio-Medical-LLaMA-8B¹, a domain-adapted version of LLaMA-3-8B-Instruct fine-tuned on a specialized biomedical dataset, as the backbone LLM. To adapt it for our classification task, we extend the model by appending a linear classification head atop the transformer architecture. We first form the multimodal fusion embeddings by concatenating audio and text representations. These fused embeddings are then fed into the LLM to produce a sequence of hidden states. A pooling layer is then applied to obtain the latent representation z_h . Specifically, we adopt the default pooling strategy used in LlamaForSequenceClassification, which selects the hidden state corresponding to the final token in the sequence. Finally, a linear classification head is applied to the pooled representation to produce prediction scores for different respiratory disease identification tasks.

$$z_h = Pool_{final}(f_{LLM}(z_{fusion}))$$
³²⁴

This vector z_h is then passed through a classification head comprising fully connected layers, followed by a softmax function to produce class probability distributions. The model is trained using cross-entropy loss:

¹https://huggingface.co/ContactDoctor/Bio-Medical-Llama-3-8B

Table 1: Summary of the datasets and tasks.

Task ID	Dataset	Disease	#Train/Test
T1	UK COVID-19 (Coppock et al., 2024)	COVID-19	20717/11121
T2	Coughvid (Orlandic et al., 2021)	COVID-19	7958/2464
T3	TBscreen (Sharma et al., 2024)	TB	20302/8051
T4	ICBHI (Rocha et al., 2019)	COPD	462/366
T5	Coswara (Bhattacharya et al., 2023)	COVID-19	-/1747
T6	CodaTB (Huddart et al., 2024)	TB	-/2053
T7	KAUH (Fraiwan et al., 2022)	COPD	-/132
T8	KAUH (Fraiwan et al., 2022)	Asthma	-/201
Т9	KAUH (Fraiwan et al., 2022)	Pneumonia	-/120

$$\mathcal{L}_{\text{CE}} = -\sum_{i=1}^{C} y_i \log(\hat{y}_i)$$

where y_i and \hat{y}_i are the true and predicted probabilities for class *i*, respectively.

Training Details: The instruction tuning process combines task-specific instructions x_p with multimodal audio (x_a) and text (x_c) inputs to ensure the model generates outputs that align with the desired response format. Additionally, we use LoRA (Low-Rank Adaptation) (Hu et al., 2021), a parameterefficient fine-tuning (PEFT) technique designed to preserve the inherent knowledge of a pre-trained LLM. The model was trained for 40 epochs, and the training configuration further optimizes LoRA with parameters like a rank (r) of 16, scaling factor (α) of 32, and a dropout of 0.1.

4 Experimental Setup

4.1 Datasets and Tasks

We evaluate the performance of RespiraMFM using seven real-world datasets, covering five of the most common respiratory diseases: COVID-19, TB, COPD, asthma, and pneumonia. These datasets include both respiratory audio recordings (e.g., coughing sound, stethoscope sound) and the associated metadata, such as patient-reported symptoms and medical history. Based on these datasets, we construct nine respiratory disease identification tasks as summarized in Table 1. Datasets associated with tasks T1 through T4 are used for training and in-domain evaluation using held-out test sets, while datasets associated with tasks T5 through T9 are reserved for zero-shot evaluation. For each task, the model is trained on the combined training data from T1 to T4. For example, in T5, the model is trained using all data from T1 to T4 and evaluated on the T5 test set. Notably, T8 and T9 involve entirely new diseases (asthma and pneumonia) not seen during training, allowing us to assess the model's generalization ability to previously unseen conditions in a zero-shot setting. Details of each dataset and task are provided in Appendix A. 366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

382

385

386

387

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

4.2 Baselines and Evaluation Metrics

Baselines: We compare RespiraMFM against two state-of-the-art multimodal baselines: BTS (Kim et al., 2024) and RespLLM (Zhang et al., 2024b). More details about the baselines are included in Appendix C.

Evaluation Metrics: To ensure fair comparison, we follow prior works on respiratory disease detection to use the Area Under the Receiver Operating Characteristic Curve (AUROC) (Janssens and Martens, 2020) as the evaluation metric for all the tasks. To ensure robust evaluation, each result was obtained through three independent runs. The mean and standard deviation of the AUROC scores across these runs are reported.

4.3 Implementation Details

We utilized PyTorch 2.3.0, transformers 4.47.1 (Wolf et al., 2020), and accelerated on four NVIDIA A100-80G GPUs. The training process uses a batch size of 16.

5 Results

5.1 Overall Performance

First, we compare the performance of RespiraMFM with the baselines under the supervised learning setting on the held-out test sets of the training datasets on tasks T1 through T4. The results are summarized in Table 2. As shown, RespiraMFM consistently outperforms both BTS and RespLLM across all four tasks. Overall, the average AUROC RespiraMFM has achieved over tasks T1 through T4 is 0.895, representing 12.3% improvement over BTS (average AUROC: 0.797) and 14.7% gain (average AUROC: 0.780) over RespLLM. These results demonstrate the strong performance of RespiraMFM in identifying a wide range of respiratory diseases, advancing the state of the arts.

5.2 Zero-Shot Performance

Next, we evaluate the zero-shot performance of RespiraMFM under the following two scenarios.

Unseen Datasets:Regarding the unseen datasets410condition, we compare RespiraMFM with BTS411and RespLLM in performing tasks T5-T7. In these412

331

333

338

339

340

341

342

343

345

347

361

Table 2: AUROC comparison for respiratory disease recognition task. Results are shown in $mean \pm std$ format of three individual runs. The light teal color indicates the second highest results, and heavy teal color indicates the highest results. The values in parentheses represent the relative improvement (%) of RespiraMFM over the strongest baseline for each task.

Task ID	Dataset	Disease	BTS	RespLLM	RespiraMFM (ours)
T1	UK COVID-19	COVID-19	0.909 ± 0.012	0.903 ± 0.002	0.914 ± 0.002 († 1.14 %)
T2	Coughvid	COVID-19	0.617 ± 0.014	0.627 ± 0.008	$0.722 \pm 0.013~(\uparrow 15.15~\%)$
Т3	TBscreen	TB	0.670 ± 0.015	0.593 ± 0.015	$0.946 \pm 0.013~(\uparrow 41.19~\%)$
T4	ICBHI	COPD	0.993 ± 0.004	0.997 ± 0.001	$1.000 \pm 0.000 (\uparrow 0.3 \%)$

tasks, the datasets used for evaluation are not seen 413 414 during training, though the target diseases remain the same. Specifically, the models are trained on 415 one or more datasets for a given disease and eval-416 417 uated on a different, unseen dataset for the same condition. For example, in task T5, the training 418 data includes other COVID-19 datasets such as 419 UKCOVID-19 and CoughVid, and is evaluated on 420 the unseen Coswara dataset. As shown in Table 3, 421 the proposed RespiraMFM consistently outper-422 forms both multi-modal baselines on these unseen 423 datasets. Specifically, our average AUROC over 494 these tasks is 0.827, outperforming BTS by 14.3% 425 (average 0.723) and RespLLM by 12.5% (aver-426 age 0.735) on average AUROC. For example, our 427 model trained on the UKCovid-19 and CoughVid 428 datasets also show strong performance in classify-429 ing COVID-19 disease within the Coswara dataset. 430 Moreover, RespiraMFM demonstrates a 38% rel-431 ative performance improvement in COPD detection 432 on the KAUH dataset compared to the other multi-433 modal baselines. 434

435 **Unobserved Respiratory Diseases:** Regarding the unobserved respiratory diseases, we further 436 compare RespiraMFM with BTS and RespLLM 437 on the prediction of asthma (T8) and pneumonia 438 (T9). In both tasks, the models are trained on 439 datasets from T1 to T4, none of which include 440 instances of asthma or pneumonia. As shown in 441 Table 3, despite having no disease-specific training 442 data for these conditions, RespiraMFM consis-443 tently outperforms both baselines. Specifically, it 444 achieves an 19.8% relative improvement in pneu-445 monia prediction over BTS and RespLLM. Overall, 446 these results suggest that RespiraMFM general-447 448 izes effectively across datasets and to previously unseen respiratory diseases. 449

450 5.3 Effects of Data and Model Scaling

451

452

Effect of Data Scaling: To assess how the training dataset size impacts the model performance, we



Figure 3: Effect of dataset scaling.

conducted experiments on Task T1 by systemati-453 cally varying the number of training examples. In 454 this experiment, the model was trained on a com-455 bined set of UKCOVID-19, Coughvid, TBscreen, 456 and ICBHI datasets and evaluated on the held-out 457 test of the UKCOVID-19 dataset. Starting with a 458 full training set of 49,439 samples, we randomly 459 sampled subsets at varying fractions and compared 460 our model with the baselines on the same test set. 461 We explored two configurations for this experiment: 462 a single-modal setup using only audio features as 463 input, and a multi-modal setup that integrates both 464 audio and textual features as input. The results are 465 shown in Figure 3. Figure 3a, which corresponds 466 to the single-modal setting using only audio input, 467 shows a clear trend of improved performance with 468 increasing training samples, indicating that larger 469 datasets lead to better performance. Our model 470 consistently outperforms both BTS and RespLLM 471 across all data fractions, with notably strong per-472 formance even at low data availability. While all 473 models benefit from more data, ours maintains a 474 consistent lead. In contrast, Figure 3b illustrates 475 the multi-modal configuration, where both audio 476 and text features are used as input. Here, our model 477 rapidly approaches peak performance with mini-478 mal training data and significantly outperforms the 479 baselines across nearly all data scales. These re-480 sults highlight the strength of multi-modal integra-481 Table 3: AUROC comparison for respiratory disease recognition task of zero-shot prediction on new dataset. Results are shown in $mean \pm std$ format of three individual runs. The light teal color indicates the second highest results,

and heavy teal color indicates the highest results. The values in parentheses represent the relative improvement (%) of RespiraMFM over the strongest baseline for each task.

ID	Dataset	Task	BTS	RespLLM	RespiraMFM (ours)
T5	Coswara	Covid	0.905 ± 0.008	0.925 ± 0.008	0.927 ± 0.006 († 0.22 %)
T6	CodaTB	TB	0.645 ± 0.016	0.649 ± 0.018	0.681 ± 0.013 († 4.93 %)
T7	KAUH	COPD	0.619 ± 0.013	0.633 ± 0.012	$0.874 \pm 0.005~(\uparrow 38.07~\%)$
T8	KAUH	Asthma	0.632 ± 0.015	0.596 ± 0.011	0.658 ± 0.011 (↑ 4.11 %)
Т9	KAUH	pneumonia	0.542 ± 0.025	0.604 ± 0.015	0.724 ± 0.010 († 19.85 %)



Mild or No Moderate Healthy Total symptoms symptoms 0.3576 Audio 0.3571 0.7266 0.6102 Text 0.3294 0.619 0.9766 0.7934 Audio+Text 0.4047 0.6587 0.9849 0.8203

Figure 4: Performance comparison of BiomedLLaMA models with different scales (1B vs. 8B) across all tasks. The 8B model consistently outperforms the 1B model, with larger gains observed on tasks involving unseen diseases.

Table 4: Performance comparison of audio-only, textonly, and multimodal (audio+text) models across different patient groups in the Coswara dataset. **Bold** indicates the best performance and <u>underlined</u> indicates the second-best.

tion, especially in clinical contexts where labeled data is often limited. The findings suggest that multi-modal models are particularly well-suited for deployment in resource-constrained healthcare settings, offering high diagnostic performance even with sparse training data.

482

483

484

486

487

488

489

490

491

492

493

494

495

496

497

499 500

501

504

Scaling Law of Model Size: To investigate whether respiratory instruction-tuning on largerscale models yields better results, we validate 1B and 8B versions of BiomedLLaMA across all tasks. As shown in Figure 4, the 8B model matches or outperforms the 1B model on nearly all tasks, demonstrating the benefits of scaling model size in respiratory disease recognition. Notably, Performance gains are more substantial in tasks involving new and unseen diseases (T6-T9), suggesting that larger models possess stronger generalization capabilities and are better equipped to handle distribution shifts in real-world clinical settings. However, the 1B model performs competitively compared to the 8B model, suggesting it remains a viable option for deployment on resource-constrained devices such as mobile platforms.

5.4 Ablation Study

Uni-Modality vs. Multi-Modality: To assess the effectiveness of multimodal integration compared to unimodal inputs, we conducted experiments on Task T5, aiming to understand whether combining audio and textual information offers complementary benefits that improve diagnostic performance beyond what a single modality can achieve alone. In this experiment, the model is trained on the combined data from all available training datasets and evaluated in a zero-shot setting on the Coswara dataset. We select the Coswara dataset for this experiment because it provides both disease labels and additional metadata describing patient health status, including severity levels such as asymptomatic (no symptoms), mild, moderate, and healthy. We group these into three broad categories-mild or no symptoms, moderate symptoms, and healthy-and evaluate models in three configurations: audio-only input (uni-modal), text-only input (uni-modal), and multimodal input combining both audio and text. Accuracy is used as the evaluation metric for all configurations. As shown in Table 4, for cases with mild or no symptoms, the audio-only model outperforms the text-only model based on the symptom information. Conversely,

507

508

509

510

511

512

513

514

515

516

517

518

519

520

522

523

524

525

526

528

529



Figure 5: Performance comparison between single modal and multimodal models on different Tasks.

the text-only model performs better compared to the audio-only model for symptomatic and healthy individuals. On the other hand, the multimodal model, which integrates both audio and text information, consistently outperforms both unimodal models across all severity levels and on the overall dataset. In summary, these results demonstrate the clear advantage of combining multiple modalities for improved disease prediction.

Contrastive Alignment: To assess the effectiveness of the contrastive alignment module introduced in §3.3, we conducted a comparative study on tasks T1-T6 by training the model with and without this component. We evaluated two configurations: (a) using audio-only input to isolate the impact on unimodal audio data, and (b) combining audio and text to evaluate performance in a multimodal setting. For the baseline (w/o alignment), we employed a standard linear projector commonly used in prior work (Zhang et al., 2024b). Results are presented in Figure 5. As shown in the left plot (Audio Only), the alignment module consistently yields higher AUC scores across all tasks, indicating more informative audio representations for instruction tuning. Similarly, in the multimodal configuration (right plot), the aligned model matches or outperforms the baseline in every case. These findings suggest that contrastive alignment not only strengthens unimodal audio features but also contributes positively to overall representation quality in multimodal scenarios.

562Generic vs. In-Domain LLM: In this experiment,563we evaluate the contribution of specialized medi-564cal domain knowledge in the in-domain BiomedL-565LaMA model compared to the general-purpose566base LLaMA model for disease detection. We use567both the 1B and 8B variants of BiomedLLaMA to568compare against the corresponding base LLaMA569models across all tasks in multimodal settings. The



Figure 6: Performance comparison of general-purpose (LLaMA-3.2) and domain-specific (BiomedLLaMA) LLMs across various tasks. (a) shows results for 1B models, while (b) shows results for 8B models.

results are presented in Figure 6, where subfigure (a) corresponds to the 1B variant and subfigure (b) corresponds to the 8B variant. In both cases, indomain BiomedLLaMA models consistently match or outperform their general-purpose counterparts, with more pronounced gains observed in larger 8B model variants. This suggests the effectiveness of using in-domain LLMs, particularly in complex multimodal tasks where domain-specific knowledge plays a critical role in disease detection. 570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

590

591

592

593

594

596

597

598

599

600

601

602

6 Conclusion

In this paper, we introduced RespiraMFM, a multimodal foundation model designed to detect respiratory diseases by integrating respiratory sound recordings with patient-reported symptoms and medical history. We also proposed an effective method for multimodal alignment of text and audio input, demonstrating strong performance across nine tasks involving five major respiratory diseases using diverse real-world datasets. We also showed that the model can maintain high diagnostic accuracy even with limited training data, making it suitable for deployment in data-scarce healthcare environments. Overall, RespiraMFM offers a scalable, non-invasive, and clinically relevant solution for early and accurate respiratory disease detection, with the potential to support medical professionals and improve decision-making across a variety of healthcare settings.

7 Limitation

While our proposed multimodal foundation model shows strong performance across various respiratory disease detection tasks, it has some limitations.

561

531

532

533

534

535

536

705

706

707

708

709

653

The model's effectiveness depends on the qual-603 ity and consistency of symptom metadata, which can differ significantly between datasets and clin-605 ical environments. For instance, in Task T2, the model performs relatively lower compared to other COVID-19 detection tasks (T1 and T5), likely due to the limited or less informative symptom data available in the Coswara dataset, making accurate 610 diagnosis more challenging. Additionally, although the model integrates audio and symptom data, in-612 corporating additional modalities such as medical imaging or wearable sensor data could further im-614 prove its diagnostic accuracy and robustness.

8 Ethics Statement

We foresee no ethical concerns with our work. All the datasets used in this study were anonymized and excluded any participant identity information.

References

616

617

618

619

621

623

627

633

634

641

647

652

- Sangmin Bae, June-Woo Kim, Won-Yang Cho, Hyerim Baek, Soyoun Son, Byungjo Lee, Changwan Ha, Kyongpil Tae, Sungnyun Kim, and Se-Young Yun. 2023.
 Patch-mix contrastive learning with audio spectrogram transformer on respiratory sound classification. *arXiv preprint arXiv:2305.14032*.
- Sebastien Baur, Zaid Nabulsi, Wei-Hung Weng, Jake Garrison, Louis Blankemeier, Sam Fishman, Christina Chen, Sujay Kakarmath, Minyoi Maimbolwa, Nsala Sanjase, et al. 2024. Hearhealth acoustic representations. *arXiv preprint arXiv:2403.02522.*
 - Debarpan Bhattacharya, Neeraj Kumar Sharma, Debottam Dutta, Srikanth Raj Chetupalli, Pravin Mote, Sriram Ganapathy, C Chandrakiran, Sahiti Nori, KK Suhail, Sadhana Gonuguntla, et al. 2023. Coswara: A respiratory sounds and symptoms dataset for remote screening of sars-cov-2 infection. *Scientific Data*, 10(1):397.
 - Yi Chang, Zhao Ren, Thanh Tam Nguyen, Wolfgang Nejdl, and Björn W Schuller. 2022. Example-based explanations with adversarial attacks for respiratory sound analysis. *arXiv preprint arXiv:2203.16141*.
- Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. 2022. Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection. In *ICASSP 2022-2022 IEEE International Conference on Acoustics*, Speech and Signal Processing (ICASSP), pages 646– 650. IEEE.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for

contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR.

- Harry Coppock, George Nicholson, Ivan Kiskin, Vasiliki Koutra, Kieran Baker, Jobie Budd, Richard Payne, Emma Karoune, David Hurley, Alexander Titcomb, et al. 2024. Audio-based ai classifiers show no evidence of improved covid-19 screening over simple symptoms checkers. *Nature Machine Intelligence*, 6(2):229–242.
- Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. 2023. Clap learning audio concepts from natural language supervision. In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5. IEEE.
- Mohammad Fraiwan, Luay Fraiwan, Mohanad Alkhodari, and Omnia Hassanin. 2022. Recognition of pulmonary diseases from lung sounds using convolutional neural networks and long short-term memory. *Journal of Ambient Intelligence and Humanized Computing*, pages 1–13.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Sophie Huddart, Vijay Yadav, Solveig K Sieberts, Larson Omberg, Mihaja Raberahona, Rivo Rakotoarivelo, Issa N Lyimo, Omar Lweno, Devasahayam J Christopher, Nguyen Viet Nhung, et al. 2024. a dataset of solicited cough sound for tuberculosis triage testing. *Scientific Data*, 11(1):1149.
- A Cecile JW Janssens and Forike K Martens. 2020. Reflection on modern methods: Revisiting the area under the roc curve. *International journal of epidemiology*, 49(4):1397–1403.
- June-Woo Kim, Miika Toikkanen, Yera Choi, Seoung-Eun Moon, and Ho-Young Jung. 2024. Bts: Bridging text and sound modalities for metadata-aided respiratory sound classification. *arXiv preprint arXiv:2406.06786*.
- Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-Ilm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*.
- Yuanhuiyi Lyu, Xu Zheng, Jiazhou Zhou, and Lin Wang. 2024. Unibind: Llm-augmented unified and balanced representation space to bind them all. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 26752–26762.
- Yi Ma, Xinzi Xu, and Yongfu Li. 2020. Lungrn+ nl: An improved adventitious lung sound classification using non-local block resnet neural network with mixup data augmentation. In *Interspeech*, pages 2902–2906.

Lara Orlandic, Tomas Teijeiro, and David Atienza. 2021. The coughvid crowdsourcing dataset, a corpus for the study of large-scale cough analysis algorithms. *Scientific Data*, 8(1):156.

710

712

714 715

716

717

719

720

721

722 723

724

725

726

727 728

729

730

731

733

734

736

737

738

739 740

741

742

743

744 745

746

747

753

754 755

757

758

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.
- Bruno M Rocha, Dimitris Filos, Luís Mendes, Gorkem Serbes, Sezer Ulukaya, Yasemin P Kahya, Nikša Jakovljevic, Tatjana L Turukalo, Ioannis M Vogiatzis, Eleni Perantoni, et al. 2019. An open access database for the evaluation of respiratory sound classification algorithms. *Physiological measurement*, 40(3):035001.
 - Manuja Sharma, Videlis Nduba, Lilian N Njagi, Wilfred Murithi, Zipporah Mwongera, Thomas R Hawn, Shwetak N Patel, and David J Horne. 2024. Tbscreen: A passive cough classifier for tuberculosis screening with a controlled dataset. *Science Advances*, 10(1):eadi0282.
 - Daniel M Weinberger, Jenny Chen, Ted Cohen, Forrest W Crawford, Farzad Mostashari, Don Olson, Virginia E Pitzer, Nicholas G Reich, Marcus Russi, Lone Simonsen, et al. 2020. Estimation of excess deaths associated with the covid-19 pandemic in the united states, march to may 2020. JAMA internal medicine, 180(10):1336–1344.
 - Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.
 - Zijiang Yang, Shuo Liu, Meishu Song, Emilia Parada-Cabaleiro, and Björn W Schuller. 2020. Adventitious respiratory classification using attentive residual neural networks.
 - Yuwei Zhang, Tong Xia, Jing Han, Yu Wu, Georgios Rizos, Yang Liu, Mohammed Mosuily, Jagmohan Chauhan, and Cecilia Mascolo. 2024a. Towards open respiratory acoustic foundation models: Pretraining and benchmarking. *arXiv preprint arXiv:2406.16148*.
- Yuwei Zhang, Tong Xia, Aaqib Saeed, and Cecilia Mascolo. 2024b. Respllm: Unifying audio and text with multimodal llms for generalized respiratory health prediction. *arXiv preprint arXiv:2410.05361*.

A Additional Details on Datasets

761

763

764

765

767

771

773

774

775

776

779

782

786

In this study, we used the following datasets: **UK COVID-19**: The UK COVID-19 Vocal Audio Dataset (Coppock et al., 2024) represents the largest collection of SARS-CoV-2 PCR-referenced audio recordings to date, compiled in the United Kingdom. The dataset features PCR test results linked to 70,794 out of 72,999 participants, with 24,155 of the 25,776 confirmed positive cases accurately documented. Notably, respiratory symptoms were reported by 45.62% of the participants, providing valuable symptomatic metadata for analysis. All the audio recordings were captured in the .wav format. In our study, we adopt the official train-test split released with the dataset.



(a) Class Distribution in covid datasets (UK covid-19, coughvid and coswara)



(b) Class Distribution in TB datasets (TBscreen and Coda TB)



(c) Class Distribution in ICBHI and KAUH datasets



Coswara: The Coswara dataset (Bhattacharya et al., 2023) is a diverse collection of respiratory sounds and detailed metadata, recorded between April 2020 and February 2022 from 2,635 individuals, including 1,819 SARS-CoV-2 negative, 674 positive, and 142 recovered cases. It features nine categories of respiratory sounds, covering variations of breathing, coughing, and speech, providing a rich dataset for analyzing respiratory health. In addition to audio recordings, the dataset includes comprehensive metadata, capturing demographic



Figure 8: Symptom Occurrence Distribution by COVID-19 Test Result in UK COVID-19 and Coswara Datasets.

details such as age, gender, and geographic location, along with health-related information like symptoms, pre-existing respiratory conditions, comorbidities, and COVID-19 test status. We follow the official data split, which contains 70% samples for training, 15% for validation, and 15% for testing. 787

788

789

791

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

COUGHVID: The COUGHVID dataset (Orlandic et al., 2021) is a large-scale, publicly available collection of over 25,000 crowdsourced cough recordings, covering a diverse range of ages, genders, geographic locations, and COVID-19 statuses. The database contains approximately 35 hours of audio recordings, comprising around 37,000 segmented cough samples. An automatic cough classifier was used to filter recordings, retaining only those with a minimum probability of 0.8 of containing cough sounds. The final distribution of labeled recordings was as follows: 25% COVID-positive cases, 35% symptomatic cases, 25% healthy individuals, and 15% with no reported health status.

TBscreen: The TBscreen dataset (Sharma et al., 2024) was collected in Nairobi and comprises cough recordings from 149 subjects diagnosed with pulmonary tuberculosis (TB) and 46 control subjects with other respiratory illnesses. The dataset includes a total of 33,000 passive coughs and 1,600 forced coughs, all recorded in a controlled setting to ensure consistency across subjects with similar demographics. To standardize the data for applica-

909

910

911

912

913

914

868

869

tions, each cough recording was processed to have
a fixed duration of one second. Longer recordings
were segmented into multiple one-second audio
files, while shorter recordings were centered and
padded with zeros to maintain uniformity.

CodaTB: The CodaTB dataset (Huddart et al., 822 2024) is a large, multi-country collection of cough sounds from individuals undergoing evaluation for tuberculosis (TB). It comprises over 700,000 cough 825 recordings from 2,143 participants, along with de-826 tailed demographic, clinical, and microbiological 827 diagnostic information. The dataset was collected as part of broader TB research studies, where participants underwent a baseline questionnaire, clinical 831 examination, and sputum collection for TB testing at the time of enrollment. Comprehensive metadata accompanies the cough recordings, including 833 age, gender, height, weight, smoking status, and duration of cough. Additionally, HIV status was determined either through self-reported diagnosis or confirmed positive test results. The dataset was split into training (n = 1,105) and validation (n = 1,105)838 839 1,038) subsets.

ICBHI: The ICBHI Respiratory Sound Database (Rocha et al., 2019) was originally compiled to support the International Conference on Biomedical Health Informatics (ICBHI) 2017 scientific challenge and is now publicly available for research. It consists of a combination of public and private datasets collected independently by two research teams across two different countries over several years. The dataset contains 5.5 hours of respiratory sound recordings, comprising 6,898 respiratory cycles from 126 subjects. The 920 audio samples in the dataset have been manually annotated by respiratory experts, classifying them based on the presence of crackles, wheezes, both, or no adventitious respiratory sounds. Additionally, the dataset provides diagnostic labels for chronic obstructive pulmonary disease (COPD), pneumonia, and asthma, enabling the development of machinelearning models for disease classification.

841

842

844

849

852

853

857

KAUH: The KAUH (King Abdulaziz University Hospital) dataset (Fraiwan et al., 2022) is a collection of respiratory sound recordings from 112 subjects, including 35 healthy individuals and 77 patients with pulmonary conditions. Lung sounds were recorded using a Electronic Stethoscope, which was placed at multiple points on the chest wall to capture respiratory sounds while avoiding heart sounds. The recordings were processed and extracted using Heart and Lung Sound Visualization software, which allows exporting data with three different filter settings (Bell, Diaphragm, and Extended) to emphasize different frequency ranges relevant to lung sounds.

B More Details of Audio Encoder

We utilized the Opera-CT encoder (Zhang et al., 2024a), to extract audio features from raw audio signals. Opera-CT is a contrastive learningbased hierarchical token-semantic audio transformer (Chen et al., 2022). It operates by dividing the mel-spectrogram into patches, which are embedded as input tokens for the transformer. The model leverages a hierarchical architecture with window attention, optimizing both computational efficiency and memory usage by restricting attention to localized windows. The transformer has 31 million parameters and produces output features of size $D_a = 768$.

C Baselines

We compared RespiraMFM with the following sate-of-the-art multimodal baselines:

BTS: BTS (Kim et al., 2024) proposes a module called Bridging Text and Sound (BTS), which aligns respiratory audio and text metadata by utilizing CLAP (Elizalde et al., 2023) as a dual-purpose encoder for both modalities. In this approach, CLAP independently processes text and audio data through separate encoders. The resulting embeddings are then concatenated and passed through a linear classifier to perform the disease prediction.

RespLLM: RespLLM (Zhang et al., 2024b) introduces a multimodal approach using a pre-trained audio encoder and a Large Language Model for diagnosing respiratory diseases using audio recordings and patient metadata. RespLLM employs a trainable linear projector to align audio embeddings with the language model's input space. In contrast, our method adopts a contrastively trained projection head, which enables more effective alignment between audio and text modalities.

D Additional Details on Contrastive Aligner

D.1 Model Architecture

The contrastive alignment module is implemented as a multi-layer perceptron (MLP) with normalization and regularization components. Specifically,

Table 5: Dataset-wise	e patient symptoms	and medical history	selection
-----------------------	--------------------	---------------------	-----------

Dataset	Patient information
Uk covid-19	Age, sex, cough, new continuous cough, runny or blocked nose, shortness of breath, sore throat, abdominal pain, diarrhea, fatigue, fever, headache, changes to sense of smell or taste, loss of taste, asthma, other symptoms
Coughvid	Age, sex, fever and muscle pain, other respiratory symptoms
TBscreen	Age, sex, fever, cough, night sweats, cough with blood, smoking status, previous TB history, HIV status, cough duration
ICBHI	Age, sex, BMI, child weight, child height, recording device placement
Coswara	Age, sex, cold, cough, diarrhea, fever, loss of smell and taste, muscle pain, breathing difficulties, fatigue, sore throat
CodaTB	Age, sex, fever, weight loss, night sweats, cough with blood, previous TB history, HIV status, cough duration
KAUH	Age, sex, recording device placement, sound type

915the projection head maps an input embedding of di-916mension 768 into a higher-dimensional contrastive917space of 2048 through an intermediate hidden layer918of size 1024. The architecture consists of a linear919transformation followed by Layer Normalization,920ReLU activation, and dropout (rate = 0.1). A final921linear layer produces the output embeddings used922for contrastive supervision.

D.2 Training

923

924

925

926

927

929

931

932

933

935

937

938

939

We trained the alignment module using the same dataset employed during instruction-tuning. The model was optimized for 100 epochs with learning rate 0.001.

D.3 Embeddings Visualization

Figure 9 presents t-SNE visualizations of audio embeddings from the UK-COVID-19 and Coswara datasets, both before and after applying contrastive alignment with text. The post-alignment visualization (on the right) shows significantly improved class-wise clustering, indicating that the contrastive alignment strategy effectively enhances the discriminative power of the audio features with respect to the respiratory disease categories.

E Additional Details on Instruction-Tuning

Value
40
32
16
0.1
16
256
2e-4
AdamW
linear
0.1

Table 6: Training hyperparameters

Models	Embedding Dimension
LLaMA-3.2 (1B)	2048
LLaMA-3.2 (8B)	4096
BioMedLLaMA (1B)	2048
BioMedLLaMA (8B)	4096

Table 7: Embedding dimension of different language models used in RespiraMFM.



(b) Coswara dataset

Figure 9: t-SNE visualization of audio features across different datasets. The left panel represents the raw output from the Opera-CT encoder, while the right panel shows the embeddings generated by the contrastively trained projector.