# PLOS ONE

# Benchmarking of eight recurrent neural network variants for breath phase and adventitious sound detection on a self-developed open-access lung sound database—HF_Lung_V1

## --Manuscript Draft--

| | |
|---|---|
| **Manuscript Number:** | PONE-D-21-10240 |
| **Article Type:** | Research Article |
| **Full Title:** | Benchmarking of eight recurrent neural network variants for breath phase and adventitious sound detection on a self-developed open-access lung sound database—HF_Lung_V1 |
| **Short Title:** | Automated lung sound analysis database |
| **Corresponding Author:** | Shang-Ran Huang<br>Heroic Faith Medical Science Co Ltd<br>Taipei, TAIWAN |
| **Keywords:** | adventitious sound; auscultation; convolutional neural networks; lung sound; recurrent neural networks; respiratory monitor |
| **Abstract:** | A reliable, remote, and continuous real-time respiratory sound monitor with automated respiratory sound analysis ability is urgently required in many clinical scenarios—such as in monitoring disease progression of coronavirus disease 2019—to replace conventional auscultation with a handheld stethoscope. However, a robust computerized respiratory sound analysis algorithm has not yet been validated in practical applications. In this study, we developed a lung sound database (HF_Lung_V1) comprising 9,765 audio files of lung sounds (duration of 15 s each), 34,095 inhalation labels, 18,349 exhalation labels, 13,883 continuous adventitious sound (CAS) labels (comprising 8,457 wheeze labels, 686 stridor labels, and 4,740 rhonchi labels), and 15,606 discontinuous adventitious sound labels (all crackles). We conducted benchmark tests for long short-term memory (LSTM), gated recurrent unit (GRU), bidirectional LSTM (BiLSTM), bidirectional GRU (BiGRU), convolutional neural network (CNN)-LSTM, CNN-GRU, CNN-BiLSTM, and CNN-BiGRU models for breath phase detection and adventitious sound detection. We also conducted a performance comparison between the LSTM-based and GRU-based models, between unidirectional and bidirectional models, and between models with and without a CNN. The results revealed that these models exhibited adequate performance in lung sound analysis. The GRU-based models outperformed, in terms of F1 scores and areas under the receiver operating characteristic curves, the LSTM-based models in most of the defined tasks. Furthermore, all bidirectional models outperformed their unidirectional counterparts. Finally, the addition of a CNN improved the accuracy of lung sound analysis, especially in the CAS detection tasks. |
| **Order of Authors:** | Fu-Shun Hsu |
| | Shang-Ran Huang |
| | Chien-Wen Huang |
| | Chao-Jung Huang |
| | Yuan-Ren Cheng |
| | Chun-Chieh Chen |
| | Jack Hsiao |
| | Chung-Wei Chen |
| | Li-Chin Chen |
| | Yen-Chun Lai |
| | Bi-Fang Hsu |
| | Nian-Jhen Lin |

| | |
|---|---|
| | Wan-Ling Tsai |
| | Yi-Lin Wu |
| | Tzu-Ling Tseng |
| | Ching-Ting Tseng |
| | Yi-Tsun Chen |
| | Feipei Lai |

| Additional Information: | |
|---|---|
| **Question** | **Response** |
| **Financial Disclosure**<br><br>Enter a financial disclosure statement that describes the sources of funding for the work included in this submission. Review the submission guidelines for detailed requirements. View published research articles from *PLOS ONE* for specific examples.<br><br>This statement is required for submission and **will appear in the published article** if the submission is accepted. Please make sure it is accurate.<br><br>**Unfunded studies**<br>Enter: *The author(s) received no specific funding for this work.*<br><br>**Funded studies**<br>Enter a statement with the following details:<br>• Initials of the authors who received each award<br>• Grant numbers awarded to each author<br>• The full name of each funder<br>• URL of each funder website<br>• Did the sponsors or funders play any role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript?<br>• **NO** - Include this sentence at the end of your statement: *The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*<br>• **YES** - Specify the role(s) played.<br><br>\* typeset | This study was partially funded by the Raising Children Medical Foundation, Taiwan (http://http://www.raising.org.tw). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. |
| **Competing Interests**<br><br>Use the instructions below to enter a | The authors have declared that no competing interests exist. |

competing interest statement for this submission. On behalf of all authors, disclose any competing interests that could be perceived to bias this work—acknowledging all financial support and any other relevant financial or non-financial competing interests.

This statement is required for submission and **will appear in the published article** if the submission is accepted. Please make sure it is accurate and that any funding sources listed in your Funding Information later in the submission form are also declared in your Financial Disclosure statement.

View published research articles from *PLOS ONE* for specific examples.

**NO authors have competing interests**

Enter: *The authors have declared that no competing interests exist.*

**Authors with competing interests**

Enter competing interest details beginning with this statement:

*I have read the journal's policy and the authors of this manuscript have the following competing interests: [insert competing interests here]*

\* typeset

---

**Ethics Statement**

Enter an ethics statement for this submission. This statement is required if the study involved:

• Human participants
• Human specimens or tissue
• Vertebrate animals or cephalopods
• Vertebrate embryos or tissues
• Field research

The recordings were approved by the Research Ethics Review Committee of Far Eastern Memorial Hospital (case number: 107052-F). Written informed consent was obtained from the 18 patients. This study was conducted in accordance with the 1964 Helsinki Declaration and its later amendments or comparable ethical standards.

Write "N/A" if the submission does not require an ethics statement.

General guidance is provided below. Consult the submission guidelines for detailed instructions. **Make sure that all information entered here is included in the Methods section of the manuscript.**

**Format for specific study types**

**Human Subject Research (involving human participants and/or tissue)**
- Give the name of the institutional review board or ethics committee that approved the study
- Include the approval number and/or a statement indicating approval of this research
- Indicate the form of consent obtained (written/oral) or the reason that consent was not obtained (e.g. the data were analyzed anonymously)

**Animal Research (involving vertebrate animals, embryos or tissues)**
- Provide the name of the Institutional Animal Care and Use Committee (IACUC) or other relevant ethics board that reviewed the study protocol, and indicate whether they approved this research or granted a formal waiver of ethical approval
- Include an approval number if one was obtained
- If the study involved *non-human primates*, add *additional details* about animal welfare and steps taken to ameliorate suffering
- If anesthesia, euthanasia, or any kind of animal sacrifice is part of the study, include briefly which substances and/or methods were applied

**Field Research**

Include the following details if this study involves the collection of plant, animal, or other materials from a natural setting:
- Field permit number
- Name of the institution or relevant body that granted permission

| | |
|---|---|
| **Data Availability**<br><br>Authors are required to make all data underlying the findings described fully available, without restriction, and from the time of publication. PLOS allows rare exceptions to address legal and ethical concerns. See the PLOS Data Policy and FAQ for detailed information.<br><br>A Data Availability Statement describing where the data can be found is required at submission. Your answers to this question constitute the Data Availability Statement and **will be published in the article**, if accepted.<br><br>**Important:** Stating 'data available on request from the author' is not sufficient. If your data are only available upon request, select 'No' for the first question and explain your exceptional situation in the text box.<br><br>Do the authors confirm that all data underlying the findings described in their manuscript are fully available without restriction? | Yes - all data are fully available without restriction |
| **Describe where the data may be found in full sentences. If you are copying our sample text, replace any instances of XXX with the appropriate details.**<br><br>• If the data are **held or will be held in a public repository**, include URLs, accession numbers or DOIs. If this information will only be available after acceptance, indicate this by ticking the box below. For example: *All XXX files are available from the XXX database (accession number(s) XXX, XXX.).*<br>• If the data are all contained **within the manuscript and/or Supporting Information files**, enter the following: *All relevant data are within the manuscript and its Supporting Information files.*<br>• If neither of these applies but you are able to provide **details of access elsewhere**, with or without limitations, please do so. For example: | All relevant data are within the manuscript and its Supporting Information files. |

| | |
|---|---|
| *Data cannot be shared publicly because of [*XXX*]. Data are available from the* XXX *Institutional Data Access / Ethics Committee (contact via* XXX*) for researchers who meet the criteria for access to confidential data.*<br><br>*The data underlying the results presented in the study are available from (include the name of the third party and contact information or URL).*<br>• This text is appropriate if the data are owned by a third party and authors do not have permission to share the data.<br><br>* typeset | |
| Additional data availability information: | |

1 **Benchmarking of eight recurrent neural network variants** for breath phase and
2 **adventitious sound detection on a self-developed open-access lung sound**
3 **database—HF_Lung_V1**

4 Fu-Shun Hsu[1,2,3], Shang-Ran Huang[3], Chien-Wen Huang[4], Chao-Jung Huang[3], Yuan-Ren Cheng[3,5,6],

5 Chun-Chieh Chen[4], Jack Hsiao[7], Chung-Wei Chen[2], Li-Chin Chen[8], Yen-Chun Lai[3], Bi-Fang Hsu[3],

6 Nian-Jhen Lin[3,9], Wan-Ling Tsai[3], Yi-Lin Wu[3], Tzu-Ling Tseng[3], Ching-Ting Tseng[3], Yi-Tsun Chen[3],

7 Feipei Lai[1,*]

8

9 [1] Graduate Institute of Biomedical Electronics and Bioinformatics, National Taiwan University, Taipei,

10 Taiwan

11 [2] Department of Critical Care Medicine, Far Eastern Memorial Hospital, New Taipei, Taiwan

12 [3] Heroic Faith Medical Science Co., Ltd., Taipei, Taiwan

13 [4] Avalanche Computing Inc., Taipei, Taiwan

14 [5] Department of Life Science, College of Life Science, National Taiwan University, Taipei, Taiwan

15 [6] Institute of Biomedical Sciences, Academia Sinica, Taipei, Taiwan

16 [7] HCC Healthcare Group, New Taipei, Taiwan

17 [8] Research Center for Information Technology Innovation, Academia Sinica, Taipei, Taiwan

18 [9] Division of Pulmonary Medicine, Far Eastern Memorial Hospital, New Taipei, Taiwan

19

20 Short Title: Automated lung sound analysis database

21

22    *Corresponding Author

23    E-mail: flai@csie.ntu.edu.tw

24

## ABSTRACT

A reliable, remote, and continuous real-time respiratory sound monitor with automated respiratory sound analysis ability is urgently required in many clinical scenarios—such as in monitoring disease progression of coronavirus disease 2019—to replace conventional auscultation with a handheld stethoscope. However, a robust computerized respiratory sound analysis algorithm has not yet been validated in practical applications. In this study, we developed a lung sound database (HF_Lung_V1) comprising 9,765 audio files of lung sounds (duration of 15 s each), 34,095 inhalation labels, 18,349 exhalation labels, 13,883 continuous adventitious sound (CAS) labels (comprising 8,457 wheeze labels, 686 stridor labels, and 4,740 rhonchi labels), and 15,606 discontinuous adventitious sound labels (all crackles). We conducted benchmark tests for long short-term memory (LSTM), gated recurrent unit (GRU), bidirectional LSTM (BiLSTM), bidirectional GRU (BiGRU), convolutional neural network (CNN)-LSTM, CNN-GRU, CNN-BiLSTM, and CNN-BiGRU models for breath phase detection and adventitious sound detection. We also conducted a performance comparison between the LSTM-based and GRU-based models, between unidirectional and bidirectional models, and between models with and without a CNN. The results revealed that these models exhibited adequate performance in lung sound analysis. The GRU-based models outperformed, in terms of *F1* scores and areas under the receiver operating characteristic curves, the LSTM-based models in most of the defined tasks. Furthermore, all bidirectional models outperformed their unidirectional counterparts. Finally, the addition of a CNN improved the accuracy of lung sound analysis, especially in the CAS

44    detection tasks.

45

## 1.  Introduction

Respiration is vital for the normal functioning of the human body. Therefore, clinical physicians are frequently required to examine respiratory conditions. Respiratory auscultation [1-3] using a stethoscope has long been a crucial first-line physical examination. The chestpiece of a stethoscope is usually placed on a patient's chest or back for lung sound auscultation or over the patient's tracheal region for tracheal sound auscultation. During auscultation, breath cycles can be inferred, which help clinical physicians evaluate the patient's respiratory rate. In addition, pulmonary pathologies are suspected when the frequency or intensity of respiratory sounds changes or when adventitious sounds, including continuous adventitious sounds (CASs) and discontinuous adventitious sounds (DASs), are identified [1, 2, 4]. Patients with coronavirus disease 2019 exhibit adventitious sounds [5]; hence, auscultation may be a useful approach for disease diagnosis [6] and disease progression tracking. However, auscultation performed using a conventional handheld stethoscope involves some limitations [7]. First, the interpretation of auscultation results substantially depends on the subjectivity of the practitioners. Even experienced clinicians might not have high consensus rates in their interpretations of auscultatory manifestations [8, 9]. Second, auscultation is a qualitative analysis method. Comparing auscultation results between individuals and quantifying the sound change by reviewing historical records are difficult tasks. Third, prolonged continuous monitoring of respiratory sound is almost impractical.

5

66    To overcome the aforementioned limitations, computerized methods for respiratory sound

67    recording and analyses based on traditional signal processing and machine learning have been

68    proposed and reviewed [4, 10-13]. With the advent of the deep learning era, studies have developed

69    novel deep learning–based methods for respiratory sound analysis. However, many of such studies

70    have focused on only distinguishing healthy participants from participants with respiratory disorders

71    [14-18] and distinguishing various types of normal breathing sounds from adventitious sounds

72    [19-25]. Only a few studies [26-29] have explored the use of deep learning for detecting breath

73    phases and adventitious sounds. Moreover, most previous studies on computerized lung sound

74    analysis have been limited by insufficient data. As of writing this paper, the largest reported

75    respiratory sound database is ICBHI 2017 Challenge [30], which comprises 6,898 breath cycles and

76    10,775 events of wheezes and crackles acquired from 126 individuals.

77    Data size plays a major role in the creation of a robust and accurate deep learning–based respiratory

78    sound analysis algorithm [31, 32]. Accordingly, the first aim of the present study was to establish a

79    large and open-access respiratory sound database for training such algorithms for the detection of

80    breath phase and adventitious sounds, mainly focusing on lung sounds. The second aim was to conduct

81    a benchmark test on the established lung sound database by using eight recurrent neural network

82    (RNN)-based models. RNNs [33] are effective for time-series analysis; long short-term memory

83    (LSTM) [34] and gated recurrent unit (GRU) [35] networks, which are two RNN variants, exhibit

84    superior performance to the original RNN model. However, whether LSTM models are superior to

6

85 GRU models (and vice versa) in many applications, particularly in respiratory sound analysis, is

86 inconclusive. Bidirectional RNN models [36, 37] can transfer not only past information to the future

87 but also future information to the past; these models consistently exhibit superior performance to

88 unidirectional RNN models in many applications [38-40] as well as in breath phase and crackle

89 detection [29]. However, whether bidirectional RNN models outperform unidirectional RNN models in

90 CAS detection has yet to be determined. Furthermore, the convolutional neural network (CNN)–RNN

91 structure has been proven to be suitable for heart sound analysis [41], lung sound analysis [19], and

92 other tasks [39, 42]. Nevertheless, the application of the CNN–RNN structure in respiratory sound

93 detection has yet to be fully investigated. Benchmarking can enable demonstrating the reliability and

94 goodness of a database; it can also be applied to investigate the performance of the RNN variants in

95 respiratory analysis.

96    In summary, the aims of this study are outlined as follows:

97 ■ Establish the largest open-access lung sound database as of writing this paper—HF_Lung_V1

98    (https://gitlab.com/techsupportHF/HF_Lung_V1).

99 ■ Conduct a performance comparison between LSTM and GRU models, between unidirectional and

100    bidirectional models, and between models with and without a CNN in breath phase and

101    adventitious sound detection based on lung sound data.

102 ■ Discuss factors influencing model performance.

103

104 **2   Establishment of the lung sound database**

105 *2.1   Data sources and patients*

106     The lung sound database was established using two sources. The first source was a database

107 used in a datathon in Taiwan Smart Emergency and Critical Care (TSECC), 2020, under the license

108 of Creative Commons Attribution 4.0 (CC BY 4.0), provided by the Taiwan Society of Emergency

109 and Critical Care Medicine. Lung sound recordings in the TSECC database were acquired from 261

110 patients.

111     The second source was sound recordings acquired from 18 residents of a respiratory care ward

112 (RCW) or a respiratory care center (RCC) in Northern Taiwan between August 2018 and October

113 2019. The recordings were approved by the Research Ethics Review Committee of Far Eastern

114 Memorial Hospital (case number: 107052-F). Written informed consent was obtained from the 18

115 patients. This study was conducted in accordance with the 1964 Helsinki Declaration and its later

116 amendments or comparable ethical standards.

117     All patients were Taiwanese and aged older than 20 years. Descriptive statistics regarding the

118 patients' demographic data, major diagnosis, and comorbidities are presented in Table 1; however,

119 information on the patients in the TSECC database is missing. Moreover, all 18 RCW/RCC residents

120 were under mechanical ventilation.

121

122 **Table 1. Demographic data of patients.**

|  | Subjects from RCW/RCC | Subjects in TSECC Database |
|---|---|---|
| Number (n) | 18 | 261 |
| Gender (M/F) | 11/7 | NA |
| Age | 67.5 (36.7, 98.3) | NA |
| Height (cm) | 163.6 (147.2, 180.0) | NA |
| Weight (kg) | 62.1 (38.2, 86.1) | NA |
| BMI (kg/m$^2$) | 23.1 (15.6, 30.7) | NA |
| Respiratory Diseases |  |  |
| ARF | 4 (22.2%) | NA |
| CRF | 8 (44.4%) | NA |
| COPD AE | 1 (5.6%) | NA |
| COPD | 2 (11.1%) | NA |
| Pneumonia | 4 (22.2%) | NA |
| ARDS | 1 (5.6%) | NA |
| Emphysema | 1 (5.6%) | NA |
| Comorbidity |  |  |
| CKD | 1 (5.6%) | NA |
| AKI | 3 (16.7%) | NA |
| CHF | 2 (11.1%) | NA |
| DM | 7 (38.9%) | NA |
| HTN | 6 (33.3%) | NA |
| Malignancy | 1 (5.6%) | NA |
| Arrythmia | 1 (5.6%) | NA |
| CAD | 1 (5.6%) | NA |

123 RCW: respiratory care ward, RCC: respiratory care center, ARF: acute respiratory failure, CRF: chronic respiratory failure, COPD AE: chronic

124 obstructive pulmonary disease acute exacerbation, COPD: chronic obstructive pulmonary disease, ARDS: acute respiratory distress syndrome, CKD:

125 chronic kidney disease, AKI: acute kidney injury, CHF: chronic heart failure, DM: diabetes, HTN: hypertension, CAD: cardiovascular disease. The mean

126 values of the age, height, weight, and BMI are presented, with the corresponding 95% CI in parentheses.

127

128

129 *2.2 Sound recording*

130 Breathing lung sounds were recorded using two devices: (1) a commercial electronic

131     stethoscope (Littmann 3200, 3M, Saint Paul, Minnesota, USA) and (2) a customized multichannel

132     acoustic recording device (HF-Type-1) that supports the connection of eight electret microphones.

133     The signals collected by the HF-Type-1 device were transmitted to a tablet (Surface Pro 6, Microsoft,

134     Redmond, Washington, USA; Fig 1). Breathing lung sounds were collected at the eight locations

135     (denoted by L1–L8) indicated in Fig 2a. The auscultation locations are described in detail in the

136     caption of Fig 2. The two devices had a sampling rate of 4,000 Hz and a bit depth of 16 bits. The

137     audio files were recorded in the WAVE (.wav) format.

138

139

140     **Fig 1. Customized multichannel acoustic recording device (HF-Type-1) connected to a tablet.**

141

142     **Fig 2. Auscultation locations and lung sound recording protocol.** (a) Auscultation locations (L1–
143     L8): L1: second intercostal space (ICS) on the right midclavicular line (MCL); L2: fifth ICS on the
144     right MCL; L3: fourth ICS on the right midaxillary line (MAL); L4: tenth ICS on the right MAL; L5:
145     second ICS on the left MCL; L6: fifth ICS on the left MCL; L7: fourth ICS on the left MAL; and L8:
146     tenth ICS on the left MAL. (b) A standard round of breathing lung sound recording with Littmann
147     3200 and HF-Type-1 devices. The white arrows represent a continuous recording, and the small red
148     blocks represent 15-s recordings. When the Littmann 3200 device was used, 15.8-s signals were
149     recorded sequentially from L1 to L8. Subsequently, all recordings were truncated to 15 s. When the
150     HF-Type-1 device was used, sounds at L1, L2, L4, L5, L6, and L8 were recorded simultaneously.
151     Subsequently, each 2-min signal was truncated to generate new 15-s audio files.

152

153         All lung sounds in the TSECC database were collected using the Littmann 3200 device only,

154     where 15.8-s recordings were obtained sequentially from L1 to L8 (Fig 2b; Littmann 3200). One

155 round of recording with the Littmann 3200 device entails a recording of lung sounds from L1 to L8.

156 The TSECC database was composed of data obtained from one to three rounds of recording with the

157 Littmann 3200 device for each patient.

158 We recorded the lung sounds of the 18 RCW/RCC residents by using both the Littmann 3200

159 device and the HF-Type-1 device. The Littmann 3200 recording protocol was in accordance with that

160 used in the TSECC database, except that data from four to five rounds of lung sound recording were

161 collected instead. The HF-Type-1 device was used to record breath sounds at L1, L2, L4, L5, L6, and

162 L8. One round of recording with the HF-Type-1 device entails a synchronous and continuous

163 recording of breath sounds for 30 min (Fig 2b; HF-Type-1). However, the recording with the

164 HF-Type-1 device was occasionally interrupted; in this case, the recording duration was <30 min.

165 Voluntary deep breathing was not mandated during the recording of lung sounds. The statistics

166 of the recordings are listed in Table 2.

167

168 **Table 2. Statistics of recordings and labels of HF_Lung_V1 database.**

|  | Littmann 3200 | HF-Type-1 | Total |
|---|---|---|---|
| Subjects |  |  |  |
| n | 261 | 18 | 261 |
| Recordings |  |  |  |
| Filename prefix | steth_ | trunc_ | NA |
| Rounds of recording | 748 | 70 | NA |
| No of 15-sec recordings | 4504 | 5261 | 9765 |
| Total duration (min) | 1126 | 1315.25 | 2441.25 |
| Labels |  |  |  |
| No of I | 16535 | 17560 | 34095 |

| | | | |
|---|---|---|---|
| Total duration of I (min) | 257.17 | 271.02 | 528.19 |
| Mean duration of I (sec) | 0.93 | 0.93 | 0.93 |
| No of E | 9107 | 9242 | 18349 |
| Total duration of E (min) | 160.25 | 132.60 | 292.85 |
| Mean duration of E (sec) | 1.06 | 0.86 | 0.96 |
| No of C/W/S/R | 6984/3974/152/2858 | 6899/4483/534/1882 | 13883/8457/686/4740 |
| Total duration of C/W/S/R (min) | 105.90/63.92/1.94/40.04 | 85.26/55.80/7.52/21.94 | 191.16/119.73/9.46/61.98 |
| Mean duration of C/W/S/R (sec) | 0.91/0.97/0.76/0.84 | 0.74/0.75/0.85/0.70 | 0.83/0.85/0.83/0.78 |
| No of D | 7266 | 8340 | 15606 |
| Total duration of D (min) | 111.75 | 55.80 | 230.87 |
| Mean duration of D (sec) | 0.92 | 0.87 | 0.89 |

169 I: inhalation, E: exhalation, W: wheeze, S: stridor, R: rhonchus, C: continuous adventitious sound, D: discontinuous adventitious sound. W, S, and R were

170 combined to form C.

171

172

173 *2.3 Audio file truncation*

174    In this study, the standard duration of an audio signal used for inhalation, exhalation, and

175 adventitious sound detection was 15 s. This duration was selected because a 15-s signal contains at

176 least three complete breath cycles, which are adequate for a clinician to reach a clinical conclusion.

177 Furthermore, a 15-s breath sound was be used previously for verification and validation [43] .

178    Because each audio file generated by the Littmann 3200 device had a length of 15.8 s, we

179 cropped out the final 0.8-s signal from the files (Fig 2b; Littmann 3200). Moreover, we used only the

180 first 15 s of each 2-min signal of the audio files (Fig 2b; HF-Type-1) generated by the HF-Type-1

181 device. Table 2 presents the number of truncated 15-s recordings and the total duration.

182

183 *2.4 Data labeling*

184     Because the data in the TSECC database contains only classification labels indicating whether a

185     CAS or DAS exists in a recording, we attempted to label the event level of all sound recordings. Two

186     board-certified respiratory therapists (NJL and YLW) and one board-certified nurse (WLT), with 8, 3,

187     and 13 years of clinical experience, respectively, were recruited to label the start and end points of

188     inhalation (I), exhalation (E), wheeze (W), stridor (S), rhonchus (R), and DAS (D) events in the

189     recordings. They labeled the sound events by listening to the recorded breath sounds while

190     simultaneously observing the corresponding patterns on a spectrogram by using customized labeling

191     software [44]. The labelers were asked not to label sound events if they could not clearly identify the

192     corresponding sound or if an incomplete event at the beginning or end of an audio file caused

193     difficulty in identification. BFH held regular meetings to ensure that the labelers had good agreement

194     on labeling criteria based on a few samples by judging the mean pseudo-$\kappa$ value [27]. When

195     developing artificial intelligence (AI) detection models, we combined the W, S, and R labels to form

196     CAS labels (C). Moreover, the D labels comprised only crackles, which were not differentiated into

197     coarse or fine crackles. The labelers were asked to label the period containing crackles but not a

198     single explosive sound (generally less than 25 ms) of a crackle. Each recording was annotated by

199     only one labeler; thus, the labels did not represent perfect ground truth. However, we used the labels

200     as ground-truth labels for model training, validation, and testing. The statistics of the labels are listed

201     in Table 2.

202

203 **3.   Inhalation, exhalation, CAS, and DAS detection**

204 *3.1 Framework*

205    The inhalation, exhalation, CAS, and DAS detection framework developed in this study is

206 displayed in Fig 3. The prominent advantage of the research framework is its modular design.

207 Specifically, each unit of the framework can be tested separately, and the algorithms in different parts

208 of the framework can be modified to achieve optimal overall performance. Moreover, the output of

209 some blocks can be used for multiple purposes. For instance, the spectrogram generated by the

210 preprocessing block can be used as the input of a model or for visualization in the user interface for

211 real-time monitoring.

212

213 **Fig. 3. Pipeline of detection framework.**

214

215    The framework comprises three parts: preprocessing, deep learning–based modeling, and

216 postprocessing. The preprocessing part involves signal processing and feature engineering

217 techniques. The deep learning–based modeling part entails the use of a well-designed neural network

218 for obtaining a sequence of classification predictions rather than a single prediction. The

219 postprocessing part involves merging the segment prediction results and eliminating the burst event.

220

221 *3.2  Preprocessing*

14

222    We processed the lung sound recordings at a sampling frequency of 4 kHz. First, to eliminate

223    the 60-Hz electrical interference and a part of the heart sound noise, we applied a high-pass filter to

224    the recordings by setting a filter order of 10 and cut-off frequency of 80 Hz. The filtered signals were

225    then processed using the short-time Fourier transform (STFT). In the STFT, we set a Hanning

226    window size of 256 and hop length of 64; no additional zero-padding was applied. Thus, a 15-s

227    sound signal could be transformed into a corresponding spectrogram with a size of $938 \times 129$. To

228    obtain the spectral information regarding the lung sounds, we extracted the following features [29,

229    45]:

230    ■    Spectrogram: We extracted 129-bin log-magnitude spectrograms.

231    ■    Mel frequency cepstral coefficients (MFCCs): We extracted 20 static coefficients, 20 delta

232         coefficients ($\Delta$), and 20 acceleration coefficients ($\Delta^2$). We used 40 mel bands within a frequency

233         range of 0–4,000 Hz. The frame width used to calculate the delta and acceleration coefficients

234         was set to 9, which resulted in a 60-bin vector per frame.

235    ■    Energy summation: We computed the energy summation of four frequency bands, namely 0–

236         250, 250–500, 500–1,000, and 0–2,000 Hz, and obtained four values per time frame.

237    After extracting the aforementioned features, we concatenated them to form a $938 \times 193$ feature

238    matrix. Subsequently, we conducted min–max normalization on each feature. The values of the

239    normalized features ranged between 0 and 1.

240

*3.3 Deep learning models*

242     We investigated the performance of eight RNN models, namely LSTM, GRU, bidirectional

243     LSTM (BiLSTM), bidirectional GRU (BiGRU), CNN-LSTM, CNN-GRU, CNN-BiLSTM, and

244     CNN-BiGRU, in terms of inhalation, exhalation, and adventitious sound detection. Fig 4 illustrates

245     the detailed model structures. The outputs of the LSTM, GRU, BiLSTM, and BiGRU models were

246     $938 \times 1$ vectors, and those of the CNN-LSTM, CNN-GRU, CNN-BiLSTM, and CNN-BiGRU

247     models were $469 \times 1$ vectors. An element in these vectors was set to 1 if an inhalation, exhalation,

248     CAS, or DAS occurred within a time segment in which the output value passed the thresholding

249     criterion; otherwise, the element was set to 0.

250

251     **Fig. 4. Model architectures and postprocessing for inhalation, exhalation, CAS, and DAS**
252     **segment and event detection.** (a) LSTM and GRU models; (b) BiLSTM and BiGRU models; and (c)
253     CNN-LSTM, CNN-GRU, CNN-BiLSTM, and CNN-BiGRU models.

254

255     For a fairer comparison of the performance of the unidirectional and bidirectional models, we

256     trained additional simplified (SIMP) BiLSTM, SIMP BiGRU, SIMP CNN-BiLSTM, and SIMP

257     CNN-BiGRU models by adjusting the number of trainable parameters. Parameter adjustment was

258     conducted by halving the number of cells of the LSTM and GRU layers.

259     We used Adam as the optimizer in the benchmark model, and we set the initial learning rate to

260     0.0001 with a step decay ($0.2\times$) when the validation loss did not decrease for 10 epochs. The learning

261    process stopped when no improvement occurred over 50 consecutive epochs.

262

263    *3.4   Postprocessing*

264        The prediction vectors obtained using the adopted models can be further processed for different

265    purposes. For example, we can transform the prediction result from frames to time for real-time

266    monitoring. The breathing duration of most humans lies within a certain range; we considered this

267    fact in our study. Accordingly, when the prediction results obtained using the models indicated that

268    two consecutive inhalation events occurred within a very small interval, we checked the continuity of

269    these two events and decided whether to merge them, as illustrated in the bottom panel of Fig 4a. For

270    example, when the interval between the *j*th and *i*th events was smaller than $T$ s, we computed the

271    difference in frequency between their energy peaks ($|\boldsymbol{p_j} - \boldsymbol{p_i}|$). Subsequently, if the difference was

272    below a given threshold $P$, the two events were merged into a single event. In the experiment, $T$ was

273    set to 0.5 s, and $P$ was set to 25 Hz. After the merging process, we further assessed whether a burst

274    event existed. If the duration of an event was shorter than 0.05 s, the event was deleted.

275

276    *3.5   Dataset arrangement and cross-validation*

277        We adopted fivefold cross-validation in the training dataset to train and validate the models.

278    Moreover, we used an independent testing dataset to test the performance of the trained models.

279    According to our preliminary experience, the acoustic patterns of the breath sounds collected from

280    one patient at different auscultation locations or between short intervals had many similarities. To

281    avoid potential data leakage caused by our methods of collecting and truncating the breath sound

282    signals, we assigned all truncated recordings collected on the same day to only one of the training,

283    validation, or testing datasets; this is because these recordings might have been collected from the

284    same patient within a short period. The statistics of the datasets are listed in Table 3. We used only

285    audio files containing CASs and DASs to train and test their corresponding detection models.

286

287 **Table 3. Statistics of the datasets and labels of the HF_Lung_V1 database.**

|  | Training Dataset | Testing Dataset | Total |
|---|---|---|---|
| Recordings | | | |
| No of 15-sec recordings | 7809 | 1956 | 9765 |
| Total duration (min) | 1952.25 | 489 | 2441.25 |
| Labels | | | |
| No of I | 27223 | 6872 | 34095 |
| Total duration of I (min) | 422.17 | 105.97 | 528.14 |
| Mean duration of I (sec) | 0.93 | 0.93 | 0.93 |
| | | | |
| No of E | 15601 | 2748 | 18349 |
| Total duration of E (min) | 248.05 | 44.81 | 292.85 |
| Mean duration of E (sec) | 0.95 | 0.98 | 0.96 |
| | | | |
| No of C/W/S/R | 11464/7027/657/3780 | 2419/1430/29/960 | 13883/8457/686/4740 |
| Total duration of C/W/S/R (min) | 160.16/100.71/9.10/50.35 | 31.01/19.02/0.36/11.63 | 191.16/119.73/9.46/61.98 |
| Mean duration of C/W/S/R (sec) | 0.84/0.86/0.83/0.80 | 0.77/0.80/0.74/0.73 | 0.83/0.85/0.83/0.78 |
| | | | |
| No of D | 13794 | 1812 | 15606 |
| Total duration of D (min) | 203.59 | 27.29 | 230.87 |
| Mean duration of D (sec) | 0.89 | 0.90 | 0.89 |

288  I: inhalation, E: exhalation, W: wheeze, S: stridor, R: rhonchus, C: continuous adventitious sound, D: discontinuous adventitious sound. W, S, and R were

289  combined to form C.

290

291

292  *3.6  Task definition and evaluation metrics*

293  [4] clearly defined classification and detection at the segment, event, and recording levels. In

294  this study, we performed two tasks. The first task involved performing detection at the segment level.

295  The acoustic signal of each lung sound recording was transformed into a spectrogram. The temporal

296  resolution of the spectrogram depended on the window size and overlap ratio of the STFT. The

297  aforementioned parameters were fixed such that each spectrogram was a matrix of size $938 \times 129$.

298  Thus, each recording contained 938 time segments (time frames), and each time segment was

299  automatically labeled (Fig 5b) according to the ground-truth event labels (Fig 5a) assigned by the

300  labelers. The output of the prediction process was a sequential prediction matrix (Fig 5c) of size 938

301  $\times 1$ in the LSTM, GRU, BiLSTM, and BiGRU models and size $469 \times 1$ in the CNN-LSTM,

302  CNN-GRU, CNN-BiLSTM, and CNN-BiGRU models. By comparing the sequential prediction with

303  the ground-truth time segments, we could define true positive (TP; orange vertical bars in Fig 5d),

304  true negative (TN; green vertical bars in Fig 5d), false positive (FP; black vertical bars in Fig 5d),

305  and false negative (FN; yellow vertical bars in Fig 5d) time segments. Subsequently, the models'

306  sensitivity and specificity in classifying the segments in each recording were computed.

307

19

308 **Fig 5. Task definition and evaluation metrics.** (a) Ground-truth event labels, (b) ground-truth time
309 segments, (c) AI inference results, (d) segment classification, (e) event detection, and (f) legend. JI:
310 Jaccard index.

311

312     The second task entailed event detection at the recording level. After completing the sequential

313 prediction (Fig 5c), we assembled the time segments associated with the same label into a

314 corresponding event (Fig 5e). We also derived the start and end times of each assembled event. The

315 Jaccard index (JI; [27] was used to determine whether an AI inference result correctly matched the

316 ground-truth event. For an assembled event to be designated as a TP event (orange horizontal bars in

317 Fig 5e), the corresponding JI value must be greater than 0.5. If the JI was between 0 and 0.5, the

318 assembled event was designated as an FN event (yellow horizontal bars in Fig 5e), and if it was 0,

319 the assembled event was designated as an FP event (black horizontal bars in Fig 5e). A TN event

320 cannot be defined in the task of event detection.

321     The performance of the models was evaluated using the *F1* score, and that of segment detection

322 was evaluated using the receiver operating characteristic (ROC) curve and area under the ROC curve

323 (AUC). In addition, the mean absolute percentage error (MAPE) of event detection was derived. The

324 accuracy, positive predictive value (PPV), sensitivity, specificity, and *F1* score of the models are

325 presented in the section of Supporting information.

326

327 *3.7  Hardware and software*

20

328        We trained the baseline models on an Ubuntu 18.04 server that was provided by the National

329   Center for High-Performance Computing in Taiwan [Taiwan Computing Cloud (TWCC)] and was

330   equipped with an Intel(R) Xeon(R) Gold 6154 @3.00 GHz CPU with 90 GB RAM. To manage the

331   intensive computation involved in RNN training, we implemented the training module by using the

332   TensorFlow 2.10, CUDA 10, and CuDNN 7 programs to run the NVIDIA Titan V100 card on the

333   TWCC server for GPU acceleration.

**4 Results**

*4.1 LSTM versus GRU models*

336    Table 4 presents the *F1* scores used to compare the eight LSTM- and GRU-based models. When

337 a CNN was not added, the GRU models outperformed the LSTM models by 0.7%–9.5% in terms of

338 the *F1* scores. However, the CNN-GRU and CNN-BiGRU models did not outperform the

339 CNN-LSTM and CNN-BiLSTM models in terms of the *F1* scores (and vice versa).

340

341 **Table 4. Comparison of *F1* scores between LSTM-based models and GRU-based models.**

| Models | n of trainable parameters | Inhalation | | Exhalation | | CASs | | DASs | |
|---|---|---|---|---|---|---|---|---|---|
| | | *F1* score | | *F1* score | | *F1* score | | *F1* score | |
| | | Segment Detection | Event Detection | Segment Detection | Event Detection | Segment Detection | Event Detection | Segment Detection | Event Detection |
| LSTM | 300,609 | 73.9% | 76.1% | 51.8% | 57.0% | 15.1% | 12.2% | 62.6% | 59.1% |
| GRU | 227,265 | **76.2%** | **78.9%** | **59.8%** | **65.6%** | **24.6%** | **20.1%** | **65.9%** | **62.5%** |
| BiLSTM | 732,225 | 78.1% | 84.0% | 57.3% | 63.9% | 19.8% | 19.1% | 69.6% | 70.0% |
| BiGRU | 552,769 | **80.3%** | **86.2%** | **64.1%** | **70.9%** | **26.9%** | **25.6%** | **70.3%** | **71.4%** |
| CNN-LSTM | 3,448,513 | 77.6% | 81.1% | **57.7%** | **62.1%** | 45.3% | 42.5% | **68.8%** | 64.4% |
| CNN-GRU | 2,605,249 | **78.4%** | **82.0%** | 57.2% | 62.0% | **51.5%** | **49.8%** | 68.0% | **64.6%** |
| CNN-BiLSTM | 6,959,809 | **80.6%** | **86.3%** | 60.4% | 65.6% | 47.9% | 46.4% | **71.2%** | **70.8%** |
| CNN-BiGRU | 5,240,513 | **80.6%** | 86.2% | **62.2%** | **68.5%** | **53.3%** | **51.6%** | 70.6% | 70.0% |

342 The bold values indicate the higher *F1* score between the compared pairs of models.

343

344    According to the ROC curves presented in Fig 6a–d, the GRU-based models outperformed the

345 LSTM-based models in all compared pairs, except for one pair, in terms of DAS segment detection

346 (AUC of 0.891 for CNN-BiLSTM vs 0.889 for CNN-BiGRU).

**Fig. 6. ROC curves for (a) inhalation, (b) exhalation, (c) CAS, and (d) DAS segment detection.**

The corresponding AUC values are presented.

*4.2 Unidirectional versus bidirectional models*

As presented in Table 5, the bidirectional models outperformed their unidirectional counterparts

in all the defined tasks by 0.4%−9.8% in terms of the *F1* scores, even when the bidirectional models

had fewer trainable parameters after model adjustment.

**Table 5. Comparison of *F1* scores between the unidirectional and bidirectional models.**

| Models | n of trainable parameters | Inhalation | | Exhalation | | CASs | | DASs | |
|---|---|---|---|---|---|---|---|---|---|
| | | *F1* score | | *F1* score | | *F1* score | | *F1* score | |
| | | Segment Detection | Event Detection | Segment Detection | Event Detection | Segment Detection | Event Detection | Segment Detection | Event Detection |
| LSTM | 300,609 | 73.9% | 76.1% | 51.8% | 57.0% | 15.1% | 12.2% | 62.6% | 59.1% |
| SIMP BiLSTM | 235,073 | **77.8%** | **84.1%** | **55.8%** | **62.4%** | **19.8%** | **17.9%** | **68.8%** | **68.9%** |
| GRU | 227,265 | 76.2% | 78.9% | 59.8% | 65.6% | 24.6% | 20.1% | 65.9% | 62.5% |
| SIMP BiGRU | 178,113 | **80.1%** | **86.1%** | **63.7%** | **70.0%** | **25.0%** | **22.2%** | **70.3%** | **71.3%** |
| CNN-LSTM | 3,448,513 | 77.6% | 81.1% | 57.7% | 62.1% | 45.3% | 42.5% | 68.8% | 64.4% |
| SIMP CNN-BiLSTM | 3,382,977 | **80.0%** | **85.8%** | **60.4%** | **66.2%** | **50.8%** | **50.2%** | **70.2%** | **70.2%** |
| CNN-GRU | 2,605,249 | 78.4% | 82.0% | 57.2% | 62.0% | 51.5% | 49.8% | 68.0% | 64.6% |
| SIMP CNN-BiGRU | 2,556,097 | **80.1%** | **85.9%** | **62.4%** | **68.4%** | **52.6%** | **51.5%** | **69.9%** | **69.5%** |

The bold values indicate the higher *F1* score between the compared pairs of models. SIMP means the number of trainable parameters is adjusted.

*4.3 Models with CNN versus those without CNN*

According to Table 6, the models with a CNN outperformed those without a CNN in 26 of the

361    32 compared pairs.

362

363

**Table 6. Comparison of *F1* scores between models without and with a CNN.**

| | | Inhalation | | Exhalation | | CASs | | DASs | |
|---|---|---|---|---|---|---|---|---|---|
| | | *F1* score | | *F1* score | | *F1* score | | *F1* score | |
| Models | n of trainable parameters | Segment Detection | Event Detection | Segment Detection | Event Detection | Segment Detection | Event Detection | Segment Detection | Event Detection |
| LSTM | 300,609 | 73.9% | 76.1% | 51.8% | 57.0% | 15.10% | 12.20% | 62.60% | 59.10% |
| CNN-LSTM | 3,448,513 | **77.6%** | **81.1%** | **57.7%** | **62.1%** | **45.30%** | **42.50%** | **68.80%** | **64.40%** |
| BiLSTM | 732,225 | 76.2% | 78.9% | **59.8%** | **65.6%** | 19.80% | 17.90% | 68.80% | 68.90% |
| CNN-BiLSTM | 6,959,809 | **78.4%** | **82.0%** | 57.2% | 62.0% | **50.80%** | **50.20%** | **70.20%** | **70.20%** |
| GRU | 227,265 | 78.1% | 84.0% | 57.3% | 63.9% | 24.60% | 20.10% | 65.90% | 62.50% |
| CNN-GRU | 2,605,249 | **80.6%** | **86.3%** | **60.4%** | **65.6%** | **51.50%** | **49.80%** | **68.00%** | **64.60%** |
| BiGRU | 178,113 | 80.3% | **86.2%** | **64.1%** | **70.9%** | 25.00% | 22.20% | **70.30%** | **71.30%** |
| CNN-BiGRU | 2,556,097 | **80.6%** | **86.2%** | 62.2% | 68.5% | **52.60%** | **51.50%** | 69.90% | 69.50% |

The bold values indicate the higher *F1* score between the compared pairs of models.

366

367    The models with a CNN exhibited higher AUC values than did those without a CNN (Fig 6a–d),

368    except that BiGRU had a higher AUC value than did CNN-BiGRU in terms of inhalation detection

369    (0.963 vs 0.961), GRU had a higher AUC value than did CNN-GRU in terms of exhalation detection

370    (0.886 vs 0.883), and BiGRU had a higher AUC value than did CNN-BiGRU in terms of exhalation

371    detection (0.911 vs 0.899).

372    Moreover, compared with the LSTM, GRU, BiLSTM, and BiGRU models, the CNN-LSTM,

373    CNN-GRU, CNN-BiLSTM, and CNN-BiGRU models exhibited flatter and lower MAPE curves

374 over a wide range of threshold values in all event detection tasks (Fig 7a–d).

375

376

377 **Fig 7. MAPE curves for (a) inhalation, (b) exhalation, (c) CAS, and (d) DAS event detection.**

378 **5 Discussion**

379 *5.1 Benchmark results*

380   According to the *F1* scores presented in Table 4, among models without a CNN, the GRU and

381 BiGRU models consistently outperformed the LSTM and BiLSTM models in all defined tasks.

382 However, the GRU-based models did not have superior *F1* scores among models with a CNN.

383 Regarding the ROC curves and AUC values (Fig 6a–d), the GRU-based models consistently

384 outperformed the other models in all but one task. Accordingly, we can conclude that GRU-based

385 models perform slightly better than LSTM-based models in lung sound analysis. Previous studies

386 have also compared LSTM- and GRU-based models [38, 46, 47]. Although a concrete conclusion

387 cannot be drawn regarding whether LSTM-based models are superior to the GRU-based models (and

388 vice versa), GRU-based models have been reported to outperform LSTM-based models in terms of

389 computation time [38, 47].

390   As presented in Table 5, the bidirectional models outperformed their unidirectional counterparts

391 in all defined tasks, a finding that is consistent with several previously obtained results [29, 36, 38,

392 40].

393    A CNN can facilitate the extraction of useful features and enhance the prediction accuracy of

394    RNN-based models. The benefits engendered by a CNN are particularly vital in CAS detection. For

395    the models with a CNN, the *F1* score improvement ranged from 26.0% to 30.3% and the AUC

396    improvement ranged from 0.067 to 0.089 in the CAS detection tasks. Accordingly, we can infer that

397    considerable information used in CAS detection resides in the local positional arrangement of the

398    features. Thus, a two-dimensional CNN facilitates the extraction of the associated information.

399    Notably, CNN-induced improvements in model performance in the inhalation, exhalation, and DAS

400    detection tasks were not as high as those observed in the CAS detection tasks. The MAPE curves

401    (Fig 7a–d) reveal that a model with a CNN has more consistent predictions over various threshold

402    values.

403    In our previous study [26], an attention-based encoder–decoder architecture based on ResNet

404    and LSTM exhibited favorable performance in inhalation (*F1* score of 90.4%) and exhalation (*F1*

405    score of 93.2%) segment detection tasks. However, the model was established on the basis of a very

406    small dataset (489 recordings of 15-s-long lung sounds). Moreover, the model involves a

407    complicated architecture; hence, it is impossible to implement real-time respiratory monitoring in

408    devices with limited computing power, such as smartphones or medical-grade tablets.

409    Few studies have performed event detection at the recording level by using a comparatively

410    simple deep learning model. [29] used the BiGRU model and one-dimensional labels (similar to

411    those used in the present study) for breath phase and crackle detection. Their BiGRU model

26

412  exhibited comparable performance to our models in terms of inhalation event detection (*F1* scores,

413  87.0% vs 86.2%) and in terms of DAS event detection (*F1* scores, 72.1% vs 71.4%). However, the

414  performance of the BiGRU model differed considerably from that of our models in terms of

415  exhalation detection (*F1* scores: 84.6% vs 70.9%). One of the reasons for this discrepancy is that [29]

416  established their ground-truth labels on the basis of the gold-standard signals of a pneumotachograph.

417  Another reason is that an exhalation label is not always available following an inhalation label in our

418  data. Finally, we did not specifically control the sounds we recorded; for example, we did not ask

419  patients to perform voluntary deep breathing or keep ambient noise down. The factors influencing

420  the model performance are further discussed in the next section.

421

422  *5.2  Factors influencing model performance*

423      The benchmark performance of the proposed models may have been influenced by the

424  following factors: (1) unusual breathing patterns; (2) imbalanced data; (3) low signal-to-noise ratio

425  (SNR); (4) noisy labels, including class and attribute noise, in the database; and (5) sound

426  overlapping.

427      Fig 8 displays most of the breath patterns present in the HF_Lung_V1 database. Fig 8a

428  illustrates the general pattern of a breath cycle in the lung sounds when the ratio of inhalation to

429  exhalation durations is approximately 2:1 and an expiratory pause is noted [3, 4]. Fig 8b presents a

430  frequent condition under which an exhalation is not completely heard by the labelers. However,

27

431    because we did not ask the subjects to breath voluntarily when recording the sound, many unusual

432    breath patterns might have been recorded, such as patterns caused by shallow breathing, fast

433    breathing, and apnea as well as those caused by double triggering of the ventilator [48] and air

434    trapping [49, 50]. These unusual breathing patterns might confuse the labeling and learning

435    processes and result in poor testing results.

436

437    **Fig 8. Patterns of normal breathing lung sounds.** (a) General lung sound patterns and (b) general
438    lung sound patterns with unidentifiable exhalations. "I" represents an identifiable inhalation event, "E"
439    represents an identifiable exhalation event, and the black areas represent pause phases.

440

441        The developed database contains imbalanced numbers of inhalation and exhalation labels

442    (34,095 and 18,349, respectively) because not every exhalation was heard and labeled. In addition,

443    the proposed models may possess the capability of learning the rhythmic rise and fall of breathing

444    signals but not the capability of learning acoustic or texture features that can distinguish an

445    inhalation from an exhalation. This may thus explain the models' poor performance in exhalation

446    detection. However, these models are suitable for respiratory rate estimation and apnea detection as

447    long as appropriate inhalation detection is achieved. Furthermore, for all labels, the summation of the

448    event duration was smaller than that of the background signal duration (these factors had a ratio of

449    approximately 1:2.5 to 1:7). The aforementioned phenomenon can be regarded as foreground–

450    background class imbalance [51] and will be addressed in future studies.

451    Most of the sounds in the established database were not recorded during the patients performed

452    deep breathing; thus, the signal quality was not maximized. However, training models with such

453    nonoptimal data increase their adaptability to real-world scenarios. Moreover, the SNR may be

454    reduced by noise, such as human voices; music; sounds from bedside monitors, televisions, air

455    conditioners, fans, and radios; sounds generated by mechanical ventilators; electrical noise generated

456    by touching or moving the parts of acoustic sensors; and friction sounds generated by the rubbing of

457    two surfaces together (e.g., rubbing clothes with the skin). A poor SNR of audio signals can lead to

458    difficulties in labeling and prediction tasks. The features of some noise types are considerably similar

459    to those of adventitious sounds. The poor performance of the proposed models in CAS detection can

460    be partly attributed to the noisy environment in which the lung sounds were recorded. In particular,

461    the sounds generated by ventilators caused numerous FP events in the CAS detection tasks. Thus,

462    additional effort is required to develop a superior preprocessing algorithm that can filter out

463    influential noise or to identify a strategy to ensure that models focus on learning the correct CAS

464    features. Furthermore, the integration of active noise-canceling technology [52] or noise suppression

465    technology [53] into respiratory sound monitors can help reduce the noise from auscultatory signals.

466    The sound recordings in the HF_Lung_V1 database were labeled by only one labeler; thus,

467    some noisy labels, including class and attribute noise, may exist in the database [54]. These noisy

468    labels are attributable to (1) the different hearing abilities of the labeler, which can cause differences

469    in the labeled duration; (2) the absence of clear criteria for differentiating between target and

470 confusing events; (3) individual human errors; (4) tendency to not label events located close to the

471 beginning and end of a recording; and (5) confusion caused by unusual breath patterns and poor

472 SNRs. However, deep learning models exhibit high robustness to noisy labels [55]. Accordingly, we

473 are currently working toward establishing better ground-truth labels.

474       Breathing generates CASs and DASs under abnormal respiratory conditions. This means that

475 the breathing sound, CAS, and DAS might overlap with one another during the same period. This

476 sound overlapping, along with the data imbalance, makes the CAS and DAS detection models learn

477 to read the rise and fall of the breathing energy and falsely identify an inhalation or exhalation as

478 CAS or DAS, respectively. This FP detection was observed in our benchmark results. In the future,

479 strategies must be adopted to address the problem of sound overlap.

480

481 **6   Conclusions**

482       We established a large open-access lung sound database, namely HF_Lung_V1

483 (https://gitlab.com/techsupportHF/HF_Lung_V1), that contains 9,765 audio files of lung sounds

484 (each with a duration of 15 s), 34,095 inhalation labels, 18,349 exhalation labels, 13,883 CAS labels

485 (comprising 8,457 wheeze labels, 686 stridor labels, and 4,740 rhonchus labels), and 15,606 DAS

486 labels (all of which are crackles).

487       We also investigated the performance of eight RNN-based models in terms of inhalation,

488 exhalation, CAS detection, and DAS detection in the HF_Lung_V1 database. We determined that the

30

489   bidirectional models outperformed the unidirectional models in lung sound analysis. Furthermore,

490   the addition of a CNN to these models further improved their performance.

491       Future studies can develop more accurate respiratory sound analysis models. First, highly

492   accurate ground-truth labels should be established. Second, researchers should investigate the

493   performance of RNN-based models containing state-of-the-art convolutional layers. Third, regional

494   CNN variants can be adopted in lung sound analysis if the labels are expanded to two-dimensional

495   bounding boxes [27]. Fourth, wavelet-based approaches, empirical mode decomposition, and other

496   methods that can extract different features should be investigated [4, 56]. Finally, respiratory sound

497   monitors should be equipped with the capability of tracheal breath sound analysis [52].

498

## References

507    1.    Bohadana A, Izbicki G, Kraman SS. Fundamentals of lung auscultation. New England Journal
508    of Medicine. 2014;370(8):744-51.

509    2.    Goettel N, Herrmann MJ. Breath Sounds: From Basic Science to Clinical Practice. Anesthesia
510    & Analgesia. 2019;128(3):e42.

511    3.    Sarkar M, Madabhavi I, Niranjan N, Dogra M. Auscultation of the respiratory system. Annals
512    of thoracic medicine. 2015;10(3):158.

513    4.    Pramono RXA, Bowyer S, Rodriguez-Villegas E. Automatic adventitious respiratory sound
514    analysis: A systematic review. PloS one. 2017;12(5):e0177926.

515    5.    Wang B, Liu Y, Wang Y, Yin W, Liu T, Liu D, et al. Characteristics of Pulmonary auscultation
516    in patients with 2019 novel coronavirus in china. 2020.

517    6.    Raj V, Renjini A, Swapna M, Sreejyothi S, Sankararaman S. Nonlinear time series and principal
518    component analyses: Potential diagnostic tools for COVID-19 auscultation. Chaos, Solitons &
519    Fractals. 2020;140:110246.

520    7.    Sovijärvi A, Vanderschoot J, Earis J. Standardization of computerized respiratory sound
521    analysis. Crit Care Med. 1997;156:974-87.

522    8.    Berry MP, Martí J-D, Ntoumenopoulos G. Inter-rater agreement of auscultation, palpable
523    fremitus, and ventilator waveform sawtooth patterns between clinicians. Respiratory care.
524    2016;61(10):1374-83.

525    9.    Grunnreis FO. Intra-and interobserver variation in lung sound classification. Effect of training:
526    UiT Norges arktiske universitet; 2016.

527    10.  Gurung A, Scrafford CG, Tielsch JM, Levine OS, Checkley W. Computerized lung sound
528    analysis as diagnostic aid for the detection of abnormal lung sounds: a systematic review and
529    meta-analysis. Respiratory medicine. 2011;105(9):1396-403.

530    11.  Huq S, Moussavi Z. Acoustic breath-phase detection using tracheal breath sounds. Medical &
531    biological engineering & computing. 2012;50(3):297-308.

532    12.  Mesaros A, Heittola T, Virtanen T. Metrics for polyphonic sound event detection. Applied
533    Sciences. 2016;6(6):162.

534    13.  Pasterkamp H, Kraman SS, Wodicka GR. Respiratory sounds: advances beyond the
535    stethoscope. American journal of respiratory and critical care medicine. 1997;156(3):974-87.

536    14.  Chambres G, Hanna P, Desainte-Catherine M, editors. Automatic detection of patient with
537    respiratory diseases using lung sound analysis. 2018 International Conference on Content-Based
538    Multimedia Indexing (CBMI); 2018: IEEE.

539    15.  Demir F, Sengur A, Bajaj V. Convolutional neural networks based efficient approach for
540    classification of lung diseases. Health Inf Sci Syst. 2020;8(1):4.

33

541    16.  Hosseini M, Ren H, Rashid H-A, Mazumder AN, Prakash B, Mohsenin T. Neural Networks for
542    Pulmonary Disease Diagnosis using Auditory and Demographic Information. arXiv preprint
543    arXiv:201113194. 2020.

544    17.  Perna D, Tagarelli A. Deep Auscultation: Predicting Respiratory Anomalies and Diseases via
545    Recurrent Neural Networks. 2019 IEEE 32nd International Symposium on Computer-Based
546    Medical Systems (CBMS)2019. p. 50-5.

547    18.  Pham L, McLoughlin I, Phan H, Tran M, Nguyen T, Palaniappan R. Robust Deep Learning
548    Framework For Predicting Respiratory Anomalies and Diseases. arXiv preprint arXiv:200203894.
549    2020.

550    19.  Acharya J, Basu A. Deep Neural Network for Respiratory Sound Classification in Wearable
551    Devices Enabled by Patient Specific Model Tuning. IEEE transactions on biomedical circuits and
552    systems. 2020;14(3):535-44.

553    20.  Aykanat M, Kılıç Ö, Kurt B, Saryal S. Classification of lung sounds using convolutional neural
554    networks. EURASIP Journal on Image and Video Processing. 2017;2017(1).

555    21.  Bardou D, Zhang K, Ahmad SM. Lung sounds classification using convolutional neural
556    networks. Artif Intell Med. 2018;88:58-69.

557    22.  Chen H, Yuan X, Pei Z, Li M, Li J. Triple-Classification of Respiratory Sounds Using
558    Optimized S-Transform and Deep Residual Networks. IEEE Access. 2019;7:32845-52.

559    23.  Grzywalski T, Piecuch M, Szajek M, Breborowicz A, Hafke-Dys H, Kocinski J, et al. Practical
560    implementation of artificial intelligence algorithms in pulmonary auscultation examination. Eur J
561    Pediatr. 2019;178(6):883-90.

562    24.  Kochetov K, Putin E, Balashov M, Filchenkov A, Shalyto A. Noise Masking Recurrent Neural
563    Network for Respiratory Sound Classification. Artificial Neural Networks and Machine Learning –
564    ICANN 2018. Lecture Notes in Computer Science2018. p. 208-17.

565    25.  Li L, Xu W, Hong Q, Tong F, Wu J, editors. Classification between normal and adventitious
566    lung sounds using deep neural network. 2016 10th International Symposium on Chinese Spoken
567    Language Processing (ISCSLP); 2016: IEEE.

568    26.  Hsiao C-H, Lin T-W, Lin C-W, Hsu F-S, Lin FY-S, Chen C-W, et al., editors. Breathing Sound
569    Segmentation and Detection Using Transfer Learning Techniques on an Attention-Based
570    Encoder-Decoder Architecture. 2020 42nd Annual International Conference of the IEEE Engineering
571    in Medicine & Biology Society (EMBC); 2020: IEEE.

572    27.  Jácome C, Ravn J, Holsbø E, Aviles-Solis JC, Melbye H, Ailo Bongo L. Convolutional neural
573    network for breathing phase detection in lung sounds. Sensors. 2019;19(8):1798.

574    28.  Liu Y, Lin Y, Gao S, Zhang H, Wang Z, Gao Y, et al., editors. Respiratory sounds feature
575    learning with deep convolutional neural networks. 2017 IEEE 15th Intl Conf on Dependable,
576    Autonomic and Secure Computing, 15th Intl Conf on Pervasive Intelligence and Computing, 3rd Intl

Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech); 2017: IEEE.

29. Messner E, Fediuk M, Swatek P, Scheidl S, Smolle-Juttner F-M, Olschewski H, et al., editors. Crackle and breathing phase detection in lung sounds with deep bidirectional gated recurrent neural networks. 2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC); 2018: IEEE.

30. Rocha B, Filos D, Mendes L, Vogiatzis I, Perantoni E, Kaimakamis E, et al., editors. A respiratory sound database for the development of automated classification. International Conference on Biomedical and Health Informatics; 2017: Springer.

31. Hestness J, Narang S, Ardalani N, Diamos G, Jun H, Kianinejad H, et al. Deep learning scaling is predictable, empirically. arXiv preprint arXiv:171200409. 2017.

32. Sun C, Shrivastava A, Singh S, Gupta A, editors. Revisiting unreasonable effectiveness of data in deep learning era. Proceedings of the IEEE international conference on computer vision; 2017.

33. Elman JL. Finding structure in time. Cognitive science. 1990;14(2):179-211.

34. Hochreiter S, Schmidhuber J. Long short-term memory. Neural computation. 1997;9(8):1735-80.

35. Cho K, Van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv preprint arXiv:14061078. 2014.

36. Graves A, Schmidhuber J. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. Neural networks. 2005;18(5-6):602-10.

37. Schuster M, Paliwal KK. Bidirectional recurrent neural networks. IEEE transactions on Signal Processing. 1997;45(11):2673-81.

38. Khandelwal S, Lecouteux B, Besacier L. Comparing GRU and LSTM for automatic speech recognition. 2016.

39. Li L, Wu Z, Xu M, Meng HM, Cai L, editors. Combining CNN and BLSTM to Extract Textual and Acoustic Features for Recognizing Stances in Mandarin Ideological Debate Competition. INTERSPEECH; 2016.

40. Parascandolo G, Huttunen H, Virtanen T, editors. Recurrent neural networks for polyphonic sound event detection in real life recordings. 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP); 2016: IEEE.

41. Deng M, Meng T, Cao J, Wang S, Zhang J, Fan H. Heart sound classification based on improved MFCC features and convolutional recurrent neural networks. Neural Networks. 2020.
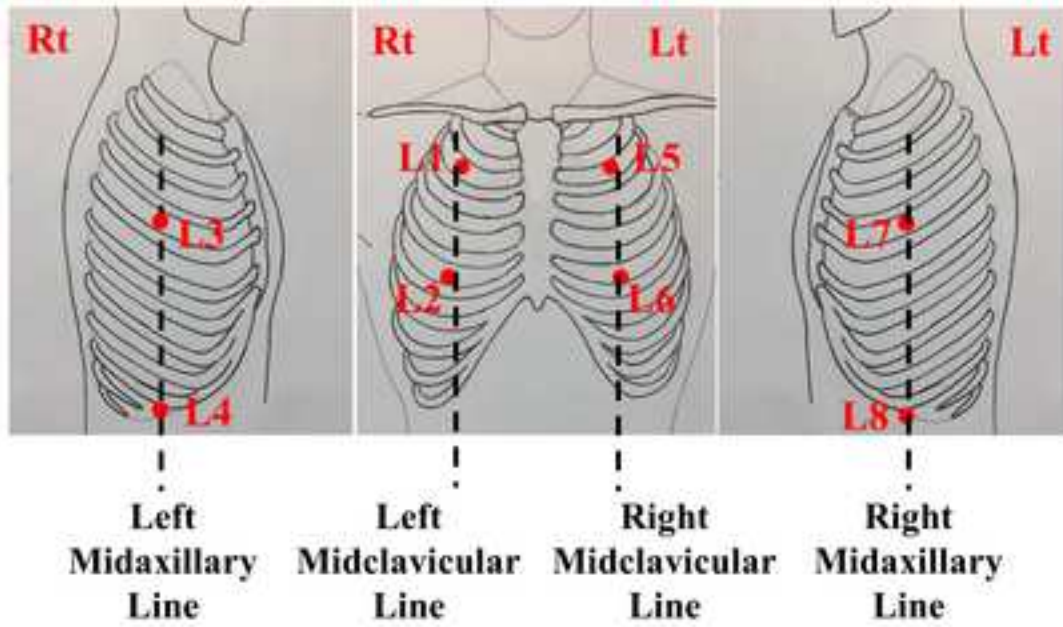
42. Zhao H, Zarar S, Tashev I, Lee C-H, editors. Convolutional-recurrent neural networks for speech enhancement. 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP); 2018: IEEE.

613    43.    Pasterkamp H, Brand PL, Everard M, Garcia-Marcos L, Melbye H, Priftis KN. Towards the
614    standardisation of lung sound nomenclature. European Respiratory Journal. 2016;47(3):724-32.
615    44.    Hsu F-S, Huang C-J, Kuo C-Y, Huang S-R, Cheng Y-R, Wang J-H, et al. Development of a
616    respiratory sound labeling software for training a deep learning-based respiratory sound analysis
617    model. arXiv preprint arXiv:210101352. 2021(, ).
618    45.    Chamberlain D, Kodgule R, Ganelin D, Miglani V, Fletcher RR, editors. Application of
619    semi-supervised deep learning to lung sound analysis. 2016 38th Annual International Conference of
620    the IEEE Engineering in Medicine and Biology Society (EMBC); 2016: IEEE.
621    46.    Chung J, Gulcehre C, Cho K, Bengio Y. Empirical evaluation of gated recurrent neural
622    networks on sequence modeling. arXiv preprint arXiv:14123555. 2014.
623    47.    Shewalkar AN. Comparison of rnn, lstm and gru on speech recognition data. 2018.
624    48.    Thille AW, Rodriguez P, Cabello B, Lellouche F, Brochard L. Patient-ventilator asynchrony
625    during assisted mechanical ventilation. Intensive care medicine. 2006;32(10):1515-22.
626    49.    Blanch L, Bernabé F, Lucangelo U. Measurement of air trapping, intrinsic positive
627    end-expiratory pressure, and dynamic hyperinflation in mechanically ventilated patients. Respiratory
628    care. 2005;50(1):110-24.
629    50.    Miller WT, Chatzkel J, Hewitt MG. Expiratory air trapping on thoracic computed tomography.
630    A diagnostic subclassification. Annals of the American Thoracic Society. 2014;11(6):874-81.
631    51.    Oksuz K, Cam BC, Kalkan S, Akbas E. Imbalance problems in object detection: A review.
632    IEEE Transactions on Pattern Analysis and Machine Intelligence. 2020.
633    52.    Wu Y, Liu J, He B, Zhang X, Yu L. Adaptive Filtering Improved Apnea Detection Performance
634    Using Tracheal Sounds in Noisy Environment: A Simulation Study. BioMed Research International.
635    2020;2020.
636    53.    Emmanouilidou D, McCollum ED, Park DE, Elhilali M. Computerized lung sound screening
637    for pediatric auscultation in noisy field environments. IEEE Transactions on Biomedical
638    Engineering. 2017;65(7):1564-74.
639    54.    Zhu X, Wu X. Class noise vs. attribute noise: A quantitative study. Artificial intelligence
640    review. 2004;22(3):177-210.
641    55.    Rolnick D, Veit A, Belongie S, Shavit N. Deep learning is robust to massive label noise. arXiv
642    preprint arXiv:170510694. 2017.
643    56.    Pramono RXA, Imtiaz SA, Rodriguez-Villegas E. Evaluation of features for classification of
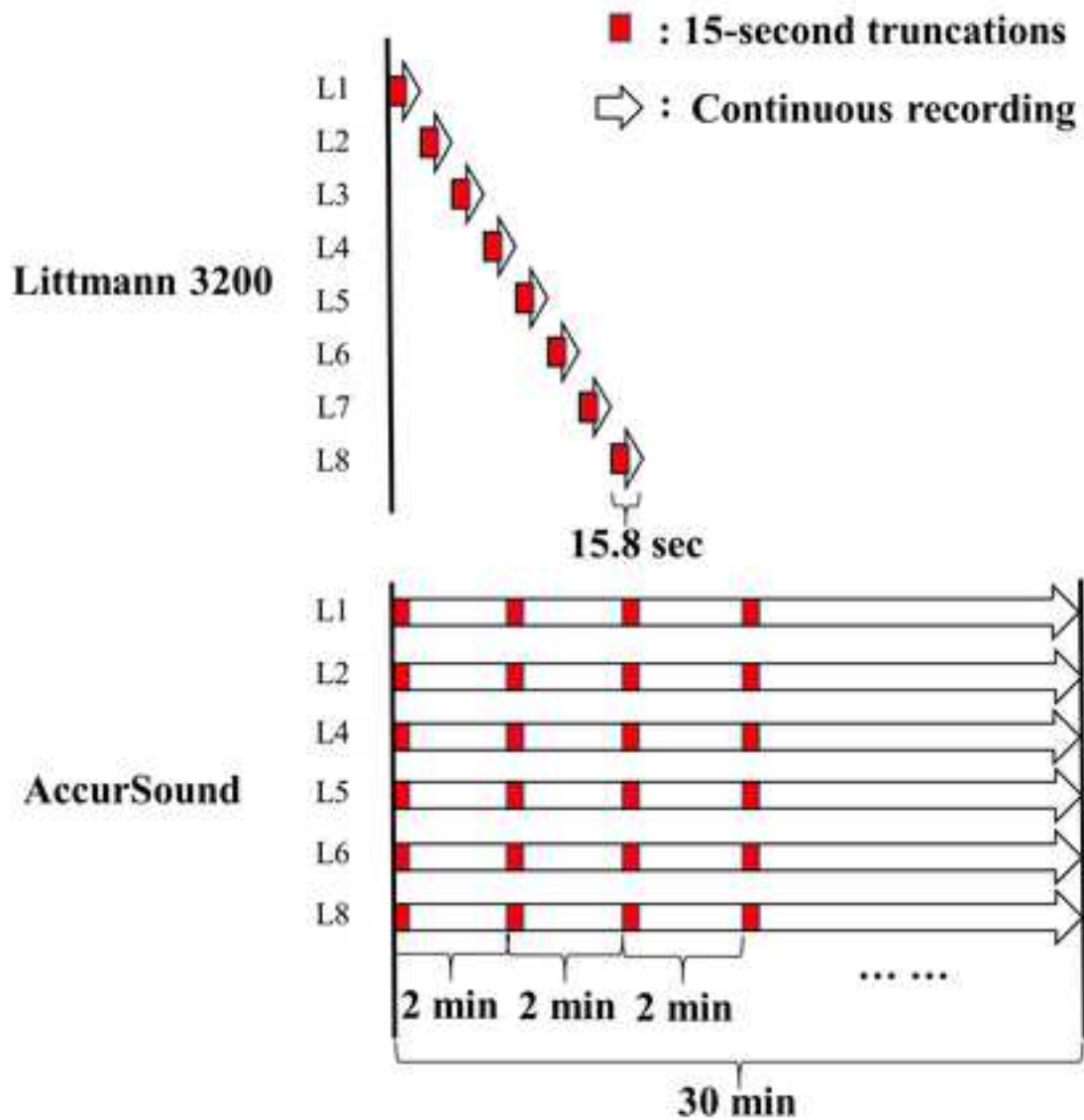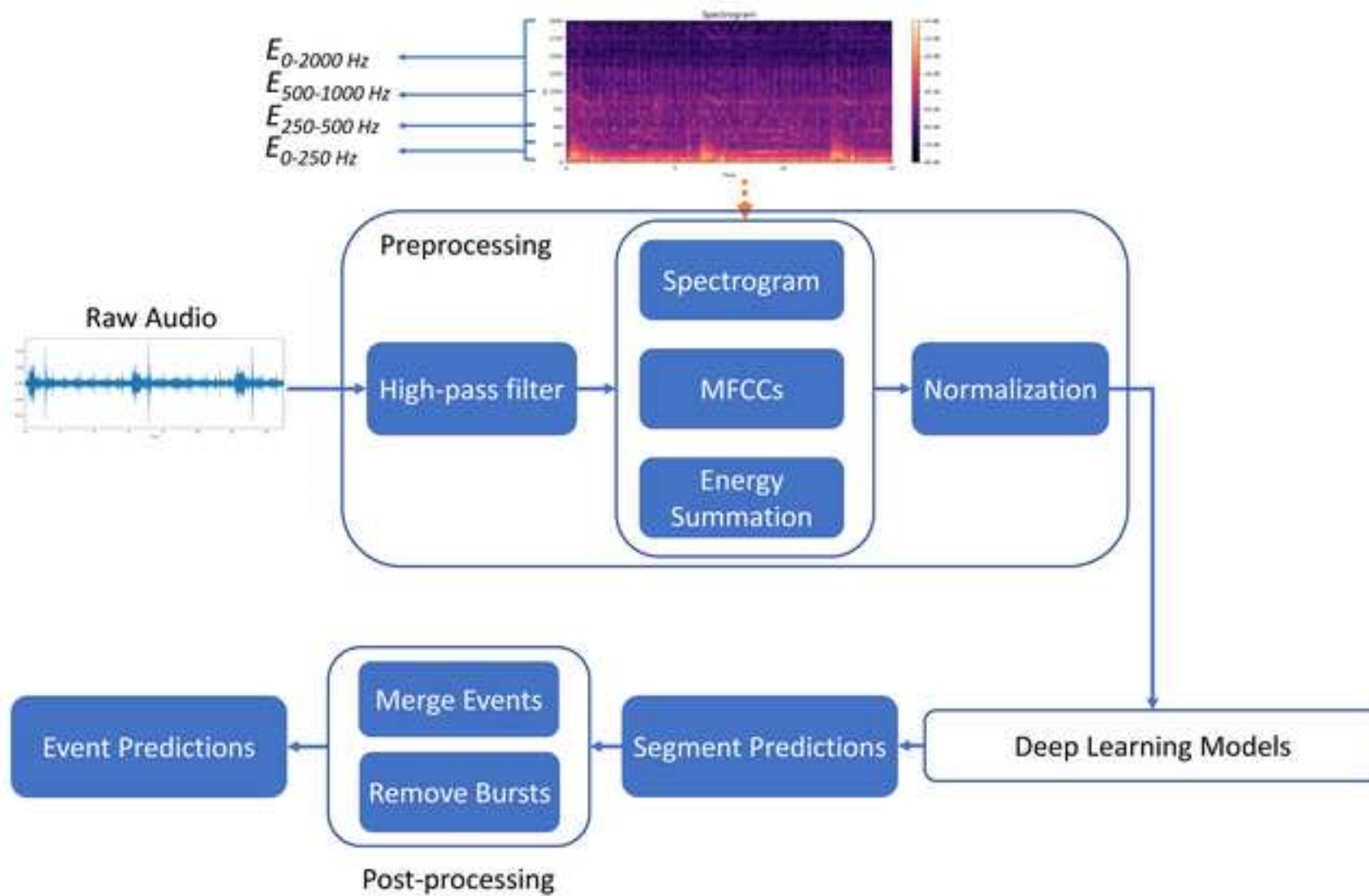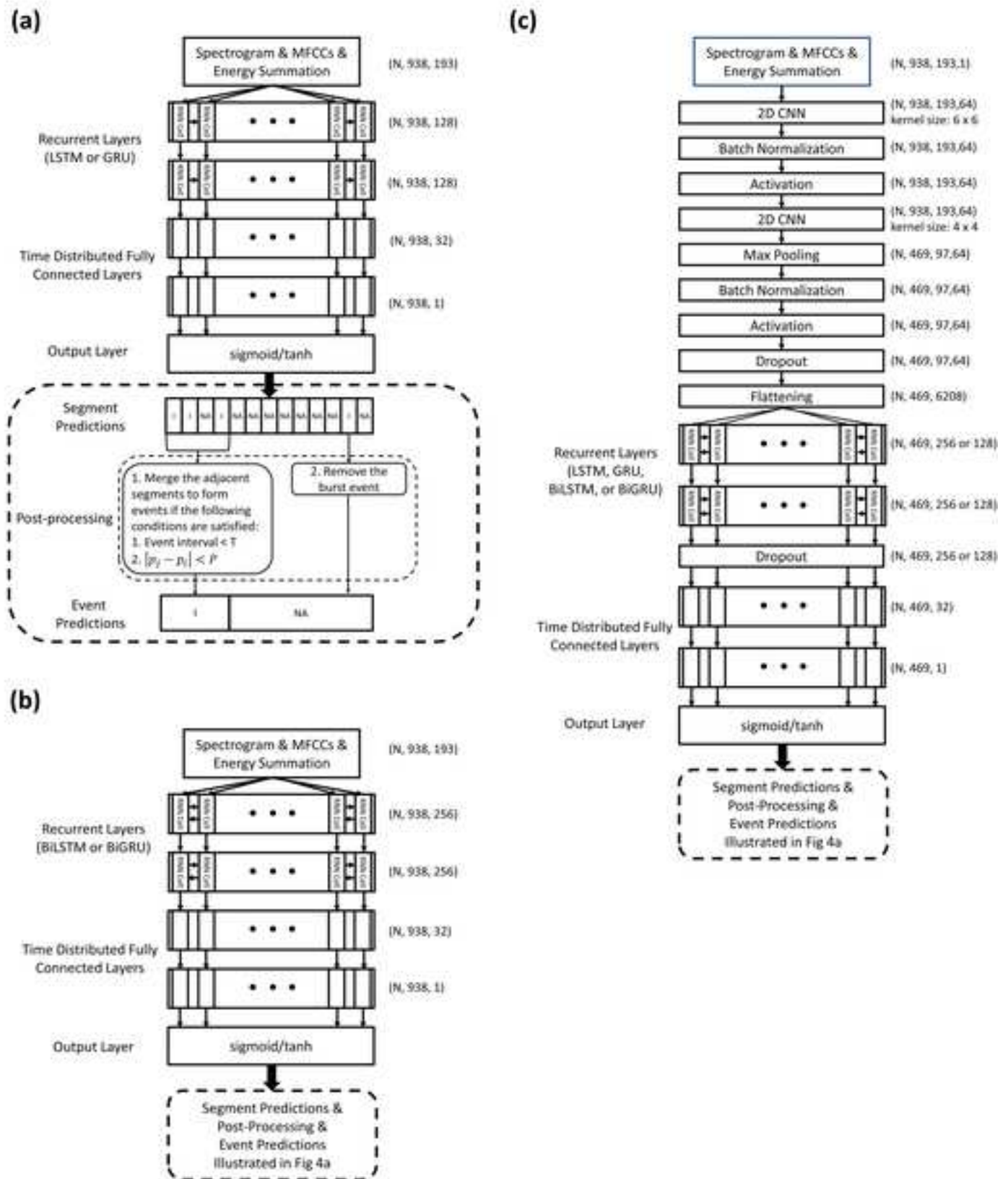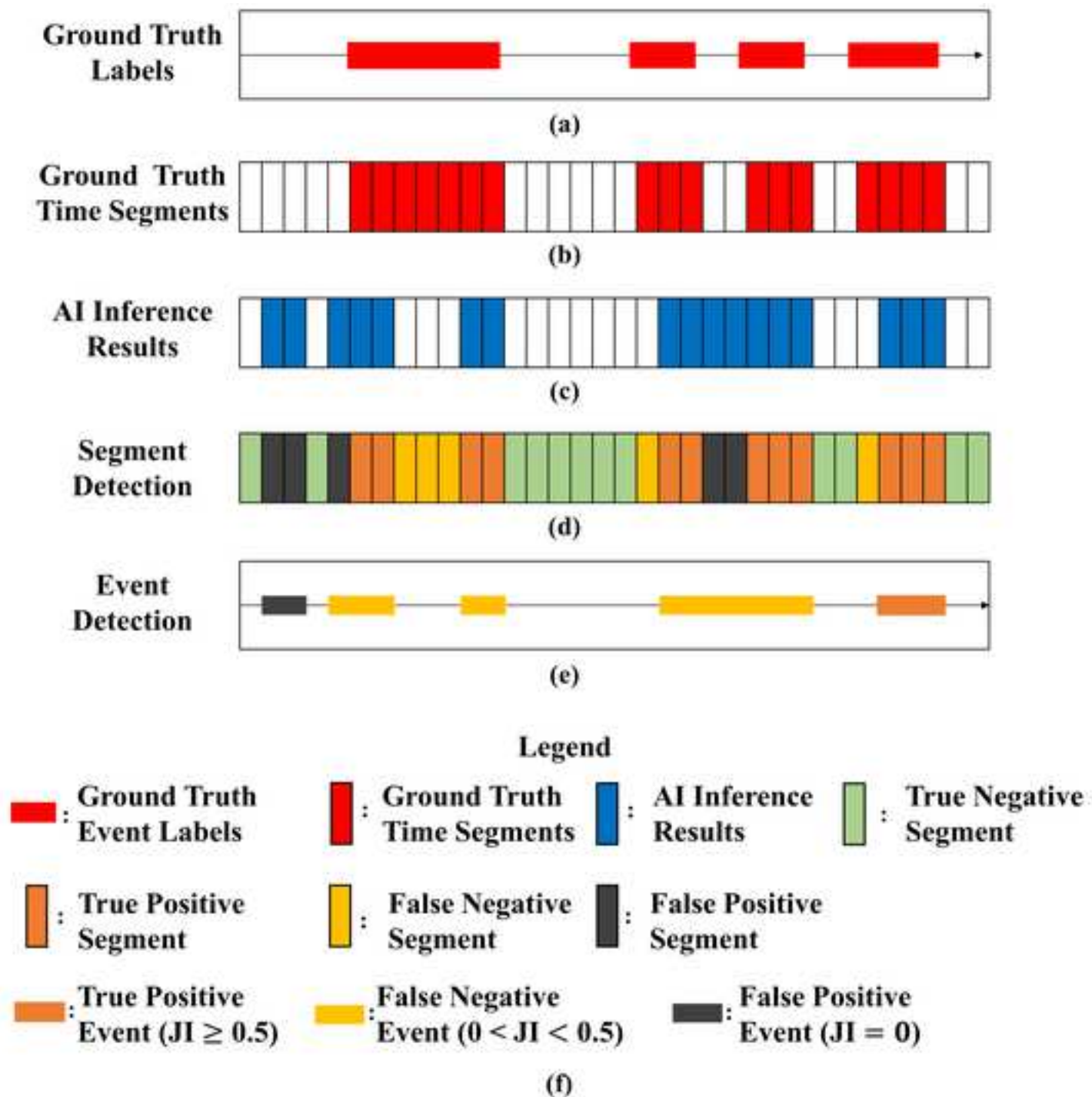644    wheezes and normal respiratory sounds. PloS one. 2019;14(3):e0213659.

645

646

Fig 1

Tablet PC

HF-Type-1

Acoustic Sensors

Fig 2

Click here to access/download;Figure;Fig2.tif ±


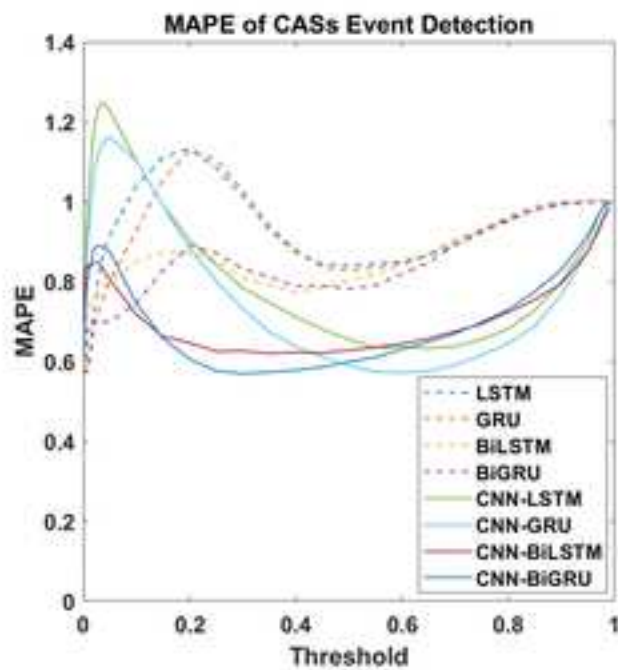
(a)

(b)

Fig 3

Fig 4

Fig 5
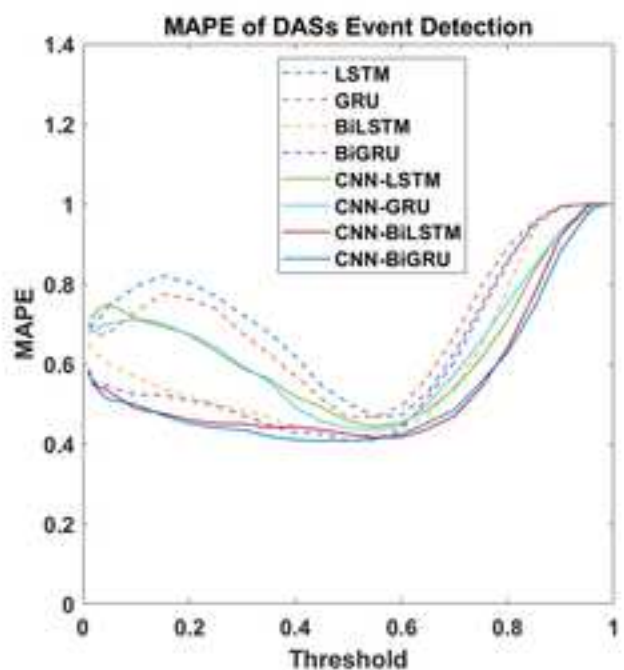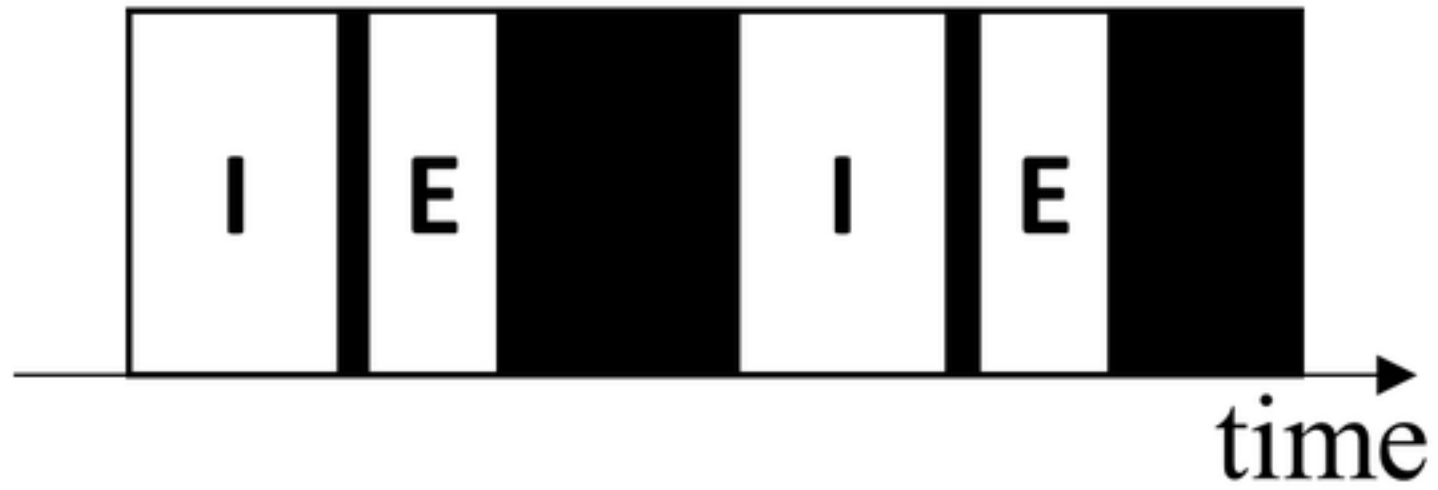
Fig 6
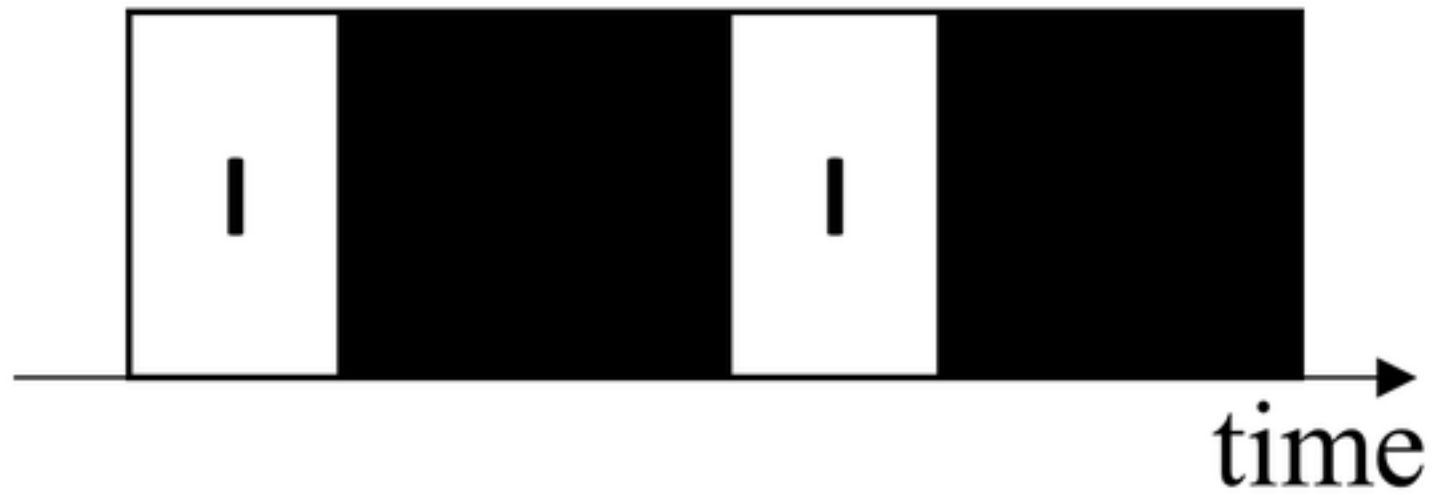
(a)

(b)

(c)

(d)

Fig 7

(a)

(b)

(c)

(d)

Fig 8

(a)

(b)

Click here to access/download
**Supporting Information**
Supporting information.docx