

MODEL-DRIVEN HEART RATE ESTIMATION AND HEART MURMUR DETECTION BASED ON PHONOCARDIOGRAM

Jingping Nie, Ran Liu, Behrooz Mahasseni, Erdrin Azemi, and Vikramjit Mitra

Apple Inc

ABSTRACT

Acoustic signals are crucial for health monitoring, particularly heart sounds which provide essential data like heart rate and detect cardiac anomalies such as murmurs. This study utilizes a publicly available phonocardiogram (PCG) dataset to estimate heart rate using model-driven methods and extends the best-performing model to a multi-task learning (MTL) framework for simultaneous heart rate estimation and murmur detection. Heart rate estimates are derived using a sliding window technique on heart sound snippets, analyzed with a combination of acoustic features (Mel spectrogram, cepstral coefficients, power spectral density, root mean square energy). Our findings indicate that a 2D convolutional neural network (**2dCNN**) is most effective for heart rate estimation, achieving a mean absolute error (*MAE*) of 1.312 bpm. We systematically investigate the impact of different feature combinations and find that utilizing all four features yields the best results. The MTL model (**2dCNN-MTL**) achieves accuracy over 95% in murmur detection, surpassing existing models, while maintaining an *MAE* of 1.636 bpm in heart rate estimation, satisfying the requirements stated by Association for the Advancement of Medical Instrumentation (AAMI).

Index Terms— Heart Rate Estimation, Phonocardiogram (PCG), Health Care, Heart Murmur Detection, Machine Learning

1. INTRODUCTION

The human heart, a symphony of rhythmic beats, encapsulates a wealth of physiological information. Heart sound signals, phonocardiogram (PCG) provide fundamental insights into heart rate (HR), cardiovascular diseases, stress assessment, and overall well-being [1]. In an era where artificial intelligence of things (AIoT) has become an integral part of our daily life, the ability to non-invasively and accurately monitor HR and cardiovascular disease from diverse contexts empowers individuals to make informed decisions about their health [2]. Emerging commercial-off-the-shelf customer-facing digital stethoscopes and research projects aim to empower users to self-monitor their cardiovascular activities from PCG [3].



Fig. 1. The overall goal of this project.

PCG produces two primary heart sounds, S1 and S2, caused by the closure of the atrioventricular and semilunar valves, respectively. Additional sounds like S3, S4, and heart murmurs can indicate cardiovascular diseases [4]. The complex real-world environments, coupled with ambient noise and noise introduced by body movement pose significant challenges to accurate HR estimation and heart murmur detection from PCGs. Traditional signal processing methods struggle with these noises due to strong distributional shifts and large variations across different scenarios [5].

Deep learning, with its ability to learn intricate patterns, is emerging as a powerful tool for HR estimation, heart sound segmentation, and heart murmur detection from noisy PCGs [6]. Some studies use deep recurrent neural networks (RNNs) to detect S1 and S2 sounds for PCG segmentation and subsequent heart rate estimation [7, 8]. Recently, deep convolutional neural networks (CNNs) have shown superior performance in detecting abnormal heart sounds using PCGs [9, 10]. Despite these advancements, the systematic exploration of CNN-based methods for heart rate estimation remains limited [11]. This absence presents a notable challenge in the development of multi-task architectures for learning heart sounds.

Considering the aforementioned challenges and opportunities, this paper aims to investigate model-driven approaches to estimate HR and detect heart murmur from short segments of heart sounds, utilizing the *CirCor DigiScope Phonocardiogram dataset*. This work first designs and compares different CNN-based models to enable reliable HR estimation and then proposes a machine learning model (**2dCNN-MTL**) to enable reliable HR estimation and heart murmur detection. Furthermore, the implications of this study reverberate across telemedicine, remote patient monitoring, and resource-constrained environments for better disease management, early diagnosis, and treatment.

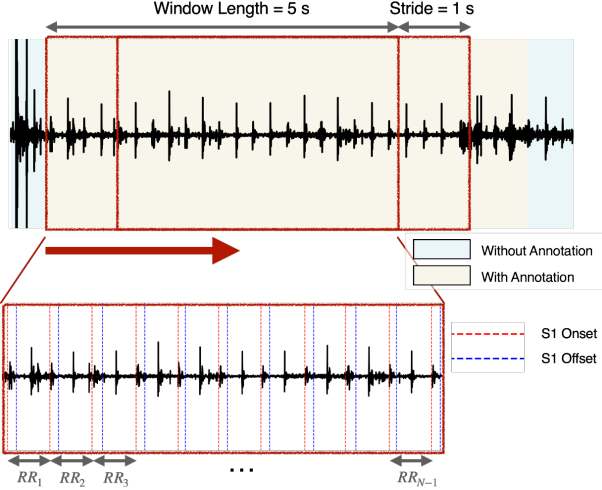


Fig. 2. The training data preparation process.

2. DATA

Datasets. The PCG dataset used in our study is the *CirCor DigiScope Phonocardiogram dataset*, which includes 3,163 heart sound recordings from 942 subjects, each spanning 5.1 to 64.5 seconds and totaling about 20 hours, with half annotated, collected across four main auscultation sites in hospitals [12]. The heart sound recordings are low-pass filtered (with a cutoff frequency of 2,000 Hz). The cardiac murmur in the dataset has been annotated in detail by an expert annotator. The segmentation annotations (onsets and offsets) regarding the location of fundamental heart sounds (S1 and S2), were obtained through a semi-supervised approach, leveraging a voting mechanism that involved three machine-learning approaches. The characteristics of PCG audio files and annotations of this dataset pose several challenges for realizing a robust HR estimation and murmur detection model. First of all, the PCG audio files are collected in the wild, where there are different low-frequency noises (e.g., environmental background noises). Second, there are annotation bias and errors in the segmentation annotations. In addition, only part of each heart sound recording is annotated.

Data Preparation. To prepare a dataset with a decent amount of labeled data, a sliding window with $window_length = 5s$ and $stride = 1s$ is applied to the raw PCG audio files with the annotated period longer than 5s, as shown in Figure 2. As such, 23,381 heart sound recording snippets are generated. The average HR (\overline{HR}) in beat per minute (*bpm*) of each audio snippet is calculated by $\overline{HR} = \frac{1}{N-1} \sum_{n=1}^{N-1} \frac{60}{RR_n}$, where RR_n is the interbeat interval between the adjacent onsets of S1 waves and N is the number of S1 waves that appear in the audio snippet. The appearance of the heart murmur is assigned to each sound snippet ($Murmur \in \{Absent, Present, Unknown\}$). The dataset is split into a training set (80%), a validation set (10%), and a test set (10%). The audio snippets in each set are from different subjects. Note that the murmur detection only applies to

the audio snippets with the murmur labels as *Absent* or *Present*.

3. MODEL DESIGN AND MICROBENCHMARKS

In this section, we discuss the acoustic features, model design reasonings, and microbenchmarks. We start with exploring, designing, and micro benchmarking the model for the more challenging task, HR estimation, and then moving towards multi-task learning (MTL) in Section 4.

3.1. Acoustic Features

In this study, we have employed a selection of prominent acoustic features to characterize audio signals. These features contain Mel spectrogram (Mel), Mel-frequency cepstral coefficients (MFCC), power spectral density (PSD), and the root mean square energy (RMS) of the audio signal, which are well-established and widely employed in the field of audio signal processing [13, 14]. PSD and RMS contain the temporal information of each sound snippet, while Mel and MFCC provide insights for both temporal and spatial information.

To generate the acoustic features in this study, the audio is resampled from 22,050 Hz to 16,000 Hz. For Mel and MFCC, the number of Mel bands and MFCCs is set to 40, the highest frequency is set to 2,000 Hz, the window size for the short-time Fourier transform (STFT) is set to be 1,024, and the hop length is set to be 160.

3.2. Model Training and Evaluation

As HR is usually presented in an integer format and usually ranges from 40 to 180 *bpm* [15]. As such, we treated the HR estimation as a 141-class classification problem and the murmur detection as a binary classification task. The weighted cross entropy (*CE*) loss and binary cross-entropy (*BCE*) are used for HR (HR) estimation and heart murmur (MM) detection tasks, respectively:

$$CE_{HR} = w_{HR} \sum_{a=1}^A \left(- \sum_{c=1}^C \log \frac{\exp(x_{a,c})}{\sum_{i=1}^C \exp(x_{n,i})} y_{a,c} \right) \quad (1)$$

$$BCE_{MM} = w_{MM} \sum_{b=1}^B [-y_b \log x_b + (1 - y_b) \log(1 - x_b)] \quad (2)$$

$$\mathcal{L} = CE_{HR} + BCE_{MM}, \quad (3)$$

where x is the input, y is the target, A is the number of audio snippets, B is the number of audio snippets containing heart murmur labels ($y_b = 0$ (*Absent*) or $y_b = 1$ (*Present*)), C is the number of HR estimation classes, \mathcal{L} is the training objective for MTL, and w_{HR} and w_{MM} are the weights as hyperparameters. The models were trained with a mini-batch size of 16 for 100 epochs. The initial learning rate was 0.001 for all models with Adam as the optimizer.

In addition, mean absolute error (MAE_{HR}) was used to evaluate the model performance for the HR estimation,

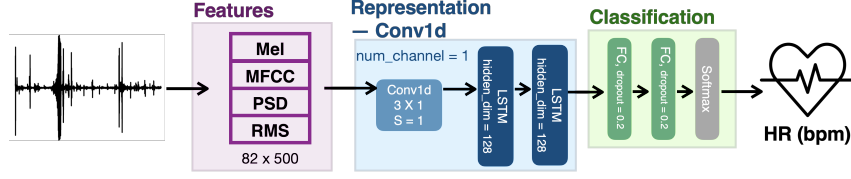


Fig. 3. HR estimation based on time convolutional (1D) neural network and LSTM (**TCNN-1stm**).

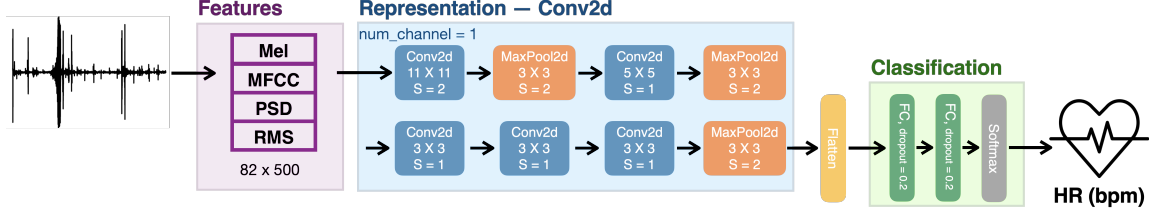


Fig. 4. HR estimation based on 2D convolutional neural network (**2dCNN**).

while accuracy (ACC_{MM}), precision ($Precision_{MM}$), and recall ($Recall_{MM}$) were used to check the performance of murmur detection tasks. MAE_{HR} and ACC_{MM} are defined as:

$$MAE_{HR} = \frac{1}{M} \sum_{i=1}^M |HR_{predicted,i} - HR_{target,i}| \quad (4)$$

$$ACC_{MM} = \frac{1}{M} \sum_{i=1}^M \mathbf{1}\{MM_{predicted,i} = MM_{target,i}\}, \quad (5)$$

where M is the total number of samples in the dataset.

3.3. HR Estimation Model Architecture

The learning objective was CE_{HR} . MAE_{HR} was used to evaluate HR estimation models. The acoustic features (Mel, MFCC, PSD, and RMS) were vertically concatenated to expand the spatial information (*Features* module). As the acoustics features contain temporal and spatial information, inspired by [9] and [16], we first investigated the temporal information by constructing a time-convolutional long short-term memory network (**TCNN-LSTM**) model as shown in Figure 3. Specifically, there is one 1D temporal convolutional layer followed by two layers of LSTM for representation extraction (*Representation* module). In addition, there are two fully connected layers with dropout to reduce overfitting and a softmax layer in the *Classification* module for HR estimation.

We then explored the information from both time- and frequency domains by designing an AlexNet-based 2D-convolutional (**2dCNN**) model inspired by [17] (see Figure 4). In particular, the *Representation* module of this **2dCNN** model has five convolutional layers, each followed by max-pooling layers, where the convolutional layers use different filter sizes and strides and the activation function used is the rectified linear unit (ReLU). Then, the multi-dimensional output from the convolutional and pooling layers is transformed into a one-dimensional vector through a flattening operation before sending it to the final *Classification* module. As illustrated in Figure 6, the **2dCNN** model exhibits a smaller mean absolute error ($MAE_{2dCNN} = 1.56$) compared to the **TCNN-LSTM** model ($MAE_{TCNN-LSTM} = 1.63$), demonstrat-

ing a statistically significant improvement with the 2D CNN architecture.

In addition, to further probe if extending the temporal features would increase the model performance, a fusion 2D-convolutional model (**2dCNN-Fusion**) was designed as illustrated in Figure 5. In this **2dCNN-Fusion** model, the vertical concatenations of Mel and PSD as well as MFCC and RMS are separately fed into two *Representation* modules, each followed by flattening operations. The resultant flattened representations are then directed to the *Classification* module. The **2dCNN-Fusion** model demonstrates a slightly improved performance ($MAE_{2dCNN-Fusion} = 1.41$) compared to the baseline **2dCNN** model ($MAE_{2dCNN} = 1.56$). However, it is important to highlight that the **2dCNN** model’s performance outshines all other models, achieving the best results when a step-wise learning rate (LR) scheduling strategy is applied. This LR scheduler is activated once the validation set’s MAE drops below 2, with a step size set at 2 and a decay rate of 0.1. With this adaptive learning rate strategy in place, the **2dCNN** model yields superior results, with an even lower MAE of 1.312. Considering the performance and model size, **2dCNN** is selected as the model and is leveraged to enable MTL (HR estimation and murmur detection).

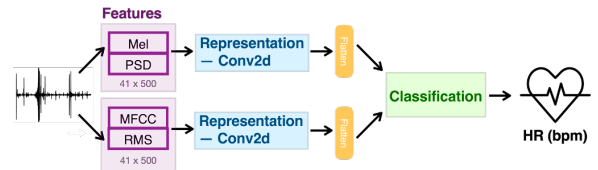


Fig. 5. HR estimation based on a fusion 2D-convolutional neural network (**2dCNN-Fusion**).

3.4. Selection of Acoustic Features

We extensively investigated the contribution of various acoustic features to the performance of the **2dCNN** model. In particular, as illustrated in Figure 7, we vertically concatenated different combinations of Mel, MFCC, PSD, and RMS

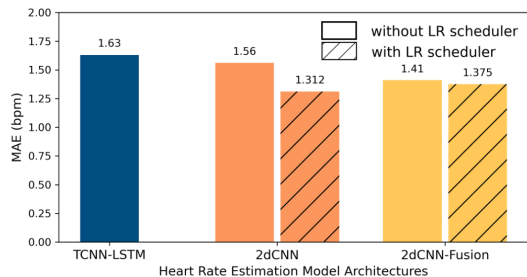


Fig. 6. MAE of different HR estimation models.

to form the *Features* module for the **2dCNN** model. When considering the use of solely Mel features within the *Features* module—an approach frequently adopted in acoustics-based biosignals estimators [18]—the model achieved an accuracy resulting in an $MAE_{2dCNN} = 2.413$ on the testing dataset. Although some of the other feature combinations achieve slightly better *MAE*, with the LR scheduler, using all four features leads to the best model performance.

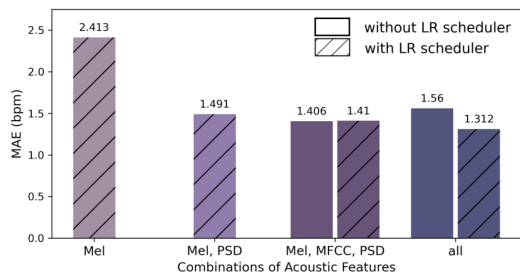


Fig. 7. MAE of **2dCNN** model with different acoustic feature combinations.

4. MTL MODEL ARCHITECTURE AND TRAINING

To facilitate HR estimation and heart murmur detection within the same model, we propose the **2dCNN-MTL** model shown in Figure 8. This model incorporates an additional *Classification* module after the flatten layer of the **2dCNN** model, specifically designed for murmur detection. The training objective for **2dCNN-MTL** model is described by Equation 3. We studied different weights (w_{HR} and w_{MM}) and the effect of adding the LR scheduler as listed in Table 1.

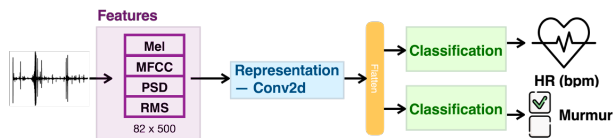


Fig. 8. HR estimation and murmur detection based on 2D convolutional neural network (**2dCNN-MTL**).

Since the **2dCNN-MTL** model is tasked with simultaneously optimizing HR estimation and murmur detection performance, a fundamental trade-off arises when selecting the best model. As shown in Table 1, when setting both w_{HR} and w_{MM} to 1 in the MTL loss (see Equation 3), the best

2dCNN-MTL models for HR estimation perform comparably to the **2dCNN** model, with ACC_{MM} reaching 92.39% and 95.19% with and without the step-wise LR scheduler, respectively. Adding the LR scheduler significantly increases the $Precision_{MM}$ and $Recall_{MM}$. The best **2dCNN-MTL** models for murmur detection achieve an accuracy of 97.49% but slightly decrease HR estimation performance. Furthermore, by increasing w_{MM} , the weights for BCE_{MM} in the MTL loss from 1 to 2, we observed an improvement in the model performance. Specifically, the accuracy (ACC_{MM}) of the best **2dCNN-MTL** models increased by 0.74% in heart rate estimation and by 1.7% in murmur detection. Additionally, there were increases in both the precision ($Precision_{MM}$) and recall ($Recall_{MM}$) metrics. However, this increase in w_{MM} results in a marginal reduction in MAE_{HR} .

Association for the Advancement of Medical Instrumentation (AAMI) states that “HR monitors should be able to compute the HR to within 10% of the reference HR, or within five beats per minute (bpm), whichever is larger” [19]. Our **2dCNN-MTL** model satisfies the AAMI requirements to proceed with HR estimation on the 5-second PCG snippets and surpasses the performance of the HR estimator for PCG recordings [5]. Furthermore, across all MTL weight configurations, the murmur detection accuracy of the **2dCNN-MTL** model is comparable to other murmur detection models on the same dataset [11, 20]. It is worth noting that our murmur detection accuracy is computed using audio snippets with labels indicating the presence (1) or absence (0) of murmurs. In contrast, the accuracy figures reported by other projects primarily rely on the detection outcomes derived from complete audio recordings within the hidden test set.

5. DISCUSSION

HR Estimation. Figure 9 shows the predicted HR alongside the target HR, using the **2dCNN** model that yielded the best performance on the heart sound snippets within the test set ($MAE = 1.312$). In general, the predicted results exhibit a strong correlation with the target HR, while the **2dCNN** model exhibits lower accuracy in estimating the HR for sound snippets with lower target HR. Notably, the outliers marked by the red circle in Figure 9 pertain to sound snippets from a participant with arrhythmia (an irregular heartbeat). These snippets have interbeat intervals (RR) ranging from 0.68 *ms* to 1.09 *ms*. Note that the annotation for the *CirCor* dataset lacks information regarding whether the participants have concurrent heart diseases alongside heart murmurs.

In addition, we deployed a baseline signal processing-based approach to HR detection using PSD. This involved processing the PSD with a Butterworth low-pass filter at a cutoff frequency of 6 *Hz* to isolate the HR frequency band. Peak detection was performed on the filtered PSD signal to identify significant peaks corresponding to S1 and S2 waves. Peak detection was then applied to the filtered PSD to iden-

Table 1. The performance of **2dCNN-MTL** with different weights in learning objective and w/ and w/o LR scheduler.

	$w_{HR} = 1, w_{MM} = 1$, No Scheduler		$w_{HR} = 1, w_{MM} = 1$, LR scheduler		$w_{HR} = 1, w_{MM} = 2$, LR scheduler	
	Best Model _{HR}	Best Model _{MM}	Best Model _{HR}	Best Model _{MM}	Best Model _{HR}	Best Model _{MM}
MAE_{HR}	1.542	1.636	1.338	1.368	1.400	1.455
ACC_{MM}	92.39%	97.49%	95.19%	95.63%	95.93%	97.33%
$Precision_{MM}$	55.11%	85.16%	68.34%	71.6%	68.20%	84.28%
$Recall_{MM}$	81.46%	86.82%	86.34%	84.88%	86.82%	86.82%

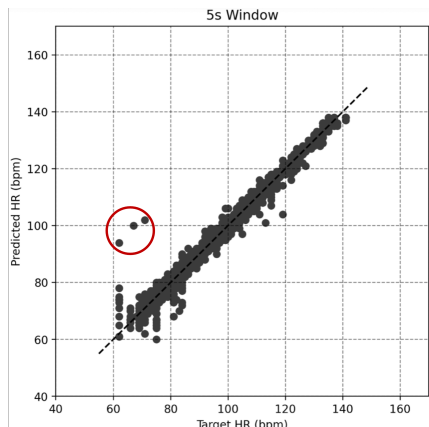


Fig. 9. The predicted HR versus the target HR for the test set with **2dCNN** model.

tify significant S1 and S2 peaks, ensuring a minimum distance of 8 samples between peaks to exclude spurious detections. We calculated inter-peak intervals for S1 and S2, converting these to beats per minute (*bpm*) to estimate the HR. Using this method, the mean absolute error is 8.2757 *bpm*, and this method can hardly estimate the sound snippets with low and high HR correctly.

Data Skewness. To further investigate the performance of the **2dCNN** model across various HR ranges, we exam the distribution of sound snippets within these ranges. The histogram depicted in Figure 10 illustrates the distribution of sound snippets across different HR ranges in the test set, revealing a degree of data skewness. In particular, the counts of sound snippets are not evenly distributed among the various HR ranges, with a notable concentration of sound snippets falling within the 90 to 100 *bpm* HR range. Furthermore, the MAE of the predicted HR for sound snippets within each HR range is depicted as green scatter points in Figure 10. The **2dCNN** model demonstrates effective performance with sound snippets having a target HR exceeding 70 *bpm*, yielding an average MAE of less than 2 *bpm*. However, its performance is less reliable when the target HR falls below 70 *bpm*. The uneven data distribution among HR ranges is a potential contributing factor to the model’s diminished performance at lower or higher target HR. Interestingly, despite having fewer snippets in the 120–140 *bpm* HR range compared to the 60–80 *bpm* range, the model exhibits more errors for the latter. This phenomenon may be linked to the scarcity of S1 and S2 segments within the 5-second window for lower HR.

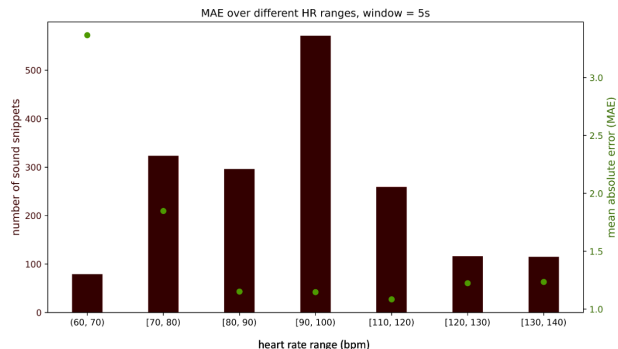


Fig. 10. The **2dCNN** model performance on different HR ranges and the number of sound snippets in each HR range in the test set.

TCNN-based MLT. We also deployed a **TCNN-LSTM-MLT** model follow the similar design of **2dCNN-MLT** model shown in Figure 8, which the *Representation 2D* module replaced by the *Representation 1D* module depicted in Figure 3. With $w_{HR} = 1$ and $w_{MM} = 1$, for the best HR model, **TCNN-LSTM-MLT** achieves MAE_{HR} of 1.9138 *bpm* and 2.0729 *bpm*, and ACC_{MM} of 90.14% and 89.7% with and without the LR scheduler. **2dCNN-MLT** outperforms **TCNN-LSTM-MLT** in the multi-task learning setup.

6. LIMITATION AND FUTURE WORK

This section summarizes current limitations and areas for further exploration. The *CirCor* dataset lacks annotations for environmental noise and respiration rate and contains low-pass filtered PCG audio files, so our method does not include explicit source separation or noise suppression steps. It would be beneficial for future studies to investigate and incorporate heart sound source separation method to remove low-frequency noises without losing acoustic features for heart murmurs. Additionally, considering the duration of PCG files in the *CirCor* dataset, as detailed in Section ??, we set the window and stride lengths to 5 *s* and 1 *s*, respectively, to generate an adequate number of heart sound snippets for training. However, reducing the window size to 3s with the same **2dCNN-MTL** model increases the MAE for HR estimation to 3.295 *bpm*. We also plan to implement a custom loss function with a penalty term weighted by the difference between predicted and true heart rates to ensure larger errors are penalized more heavily. In current model settings, treating heart rate estimation as a regression problem has

underperformed compared to treating it as a classification problem. We aim to explore more regression models and perform hyperparameter tuning to investigate the feasibility of using regression for heart rate estimation. Furthermore, it is worth noting that the PCG recordings in the *CirCor* dataset are resting heart sounds. The exploration of non-steady-state PCG data, such as post-exercise heart sounds, could significantly enhance model adaptability across various everyday scenarios and enable more applications.

7. CONCLUSION

This study presents a significant contribution to the field of health monitoring and cardiac assessment through its novel model-driven approach to heart rate estimation and heart murmur detection based on phonocardiogram (PCG) analysis. Utilizing a publicly available PCG dataset, the research demonstrated the efficacy of the proposed 2D convolutional neural network (**2dCNN**) for heart rate estimation. The model, with a mean absolute error (*MAE*) of 1.312 *bpm*, effectively integrates diverse acoustic features: *Me1*, MFCC, PSD, and RMS. This work extended to a multi-task learning (MTL) framework, encapsulated in the **2dCNN-MTL** model, which concurrently achieved heart rate estimation and murmur detection. The **2dCNN-MTL** model's accuracy exceeds 95%, surpassing existing models in both accuracy and efficiency, with a maintained *MAE* below 1.636 *bpm* in heart rate estimation. We envision the integration of these techniques to revolutionize remote patient monitoring and self-care.

8. REFERENCES

- [1] A Shyam, Vignesh Ravichandran, SP Preejith, Jayaraj Joseph, and Mohanasankar Sivaprakasam, "Ppgnet: Deep network for device independent heart rate estimation from photoplethysmogram," in *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2019, pp. 1899–1902.
- [2] Goutam Kumar Sahoo, Keerthana Kanike, Santos Kumar Das, and Poonam Singh, "Machine learning-based heart disease prediction: A study for home personalized care," in *2022 IEEE 32nd International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2022, pp. 01–06.
- [3] Kaiyuan Hou, Stephen Xia, Emily Bejerano, Junyi Wu, and Xiaofan Jiang, "ARSteth: Enabling home self-screening with ar-assisted intelligent stethoscopes," in *Proceedings of the 22nd International Conference on Information Processing in Sensor Networks*, 2023, pp. 205–218.
- [4] Tomoya Koike, Kun Qian, Qiuqiang Kong, Mark D Plumbley, Björn W Schuller, and Yoshiharu Yamamoto, "Audio for audio is better? an investigation on transfer learning models for heart sound classification," in *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 2020, pp. 74–77.
- [5] David B Springer, Thomas Brennan, Jens Hitzeroth, Bongani M Mayosi, Lionel Tarassenko, and Gari D Clifford, "Robust heart rate estimation from noisy phonocardiograms," in *Computing in Cardiology 2014*. IEEE, 2014, pp. 613–616.
- [6] John S Chorba, Avi M Shapiro, Le Le, John Maidens, John Prince, Steve Pham, Mia M Kanzawa, Daniel N Barbosa, Caroline Currie, Catherine Brooks, et al., "Deep learning algorithm for automated cardiac murmur detection via a digital stethoscope platform," *Journal of the American Heart Association*, vol. 10, no. 9, pp. e019905, 2021.
- [7] Sofia Monteiro, Ana Fred, and Hugo Plácido da Silva, "Detection of heart sound murmurs and clinical outcome with bidirectional long short-term memory networks," in *2022 Computing in Cardiology (CinC)*. IEEE, 2022, vol. 498, pp. 1–4.
- [8] Tharindu Fernando, Houman Ghaemmaghami, Simon Denman, Sridha Sridharan, Nayyar Hussain, and Clinton Fookes, "Heart sound segmentation using bidirectional lstms with attention," *IEEE journal of biomedical and health informatics*, vol. 24, no. 6, pp. 1601–1609, 2019.
- [9] Ahmed Imtiaz Humayun, Shabnam Ghaffarzadegan, Zhe Feng, and Taufiq Hasan, "Learning front-end filter-bank parameters using convolutional neural networks for abnormal heart sound detection," in *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2018, pp. 1408–1411.
- [10] Mohanad Alkhdari, Murad Almadani, Samit Kumar Ghosh, and Ahsan H Khandoker, "Fhsu-net: Deep learning-based model for the extraction of fetal heart sounds in abdominal phonocardiography," in *2023 IEEE 33rd International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2023, pp. 1–6.
- [11] Matthew A Reyna, Yashar Kiarashi, Andoni Elola, Jorge Oliveira, Francesco Renna, Annie Gu, Erick A Perez Alday, Nadi Sadr, Ashish Sharma, Sandra Mattos, et al., "Heart murmur detection from phonocardiogram recordings: The george b. moody physionet challenge 2022," in *2022 Computing in Cardiology (CinC)*. IEEE, 2022, vol. 498, pp. 1–4.
- [12] Jorge Oliveira, Francesco Renna, Paulo Costa, Marcelo Nogueira, Ana Cristina Oliveira, Andoni Elola, Carlos Ferreira, Alipio Jorge, Ali Bahrami Rad, Matthew Reyna, et al., "The circor digiscope phonocardiogram dataset," *version 1.0. 0*, 2022.
- [13] Vikramjit Mitra, Jingping Nie, and Erdrin Azemi, "Investigating salient representations and label variance in dimensional speech emotion analysis," *arXiv preprint arXiv:2312.16180*, 2023.
- [14] Stephen Xia, Jingping Nie, and Xiaofan Jiang, "Csafe: An intelligent audio wearable platform for improving construction worker safety in urban environments," in *Proceedings of the 20th International Conference on Information Processing in Sensor Networks*, 2021, pp. 207–221.
- [15] Thomas Pursche, Jarek Krajewski, and Reinhard Moeller, "Video-based heart rate measurement from human faces," in *2012 IEEE international conference on consumer electronics (ICCE)*. IEEE, 2012, pp. 544–545.
- [16] Yuanhang Su and C-C Jay Kuo, "On extended long short-term memory and dependent bidirectional recurrent neural network," *Neurocomputing*, vol. 356, pp. 151–161, 2019.
- [17] Tarik Alafif, Mehrez Boulares, Ahmed Barnawi, Talal Alafif, Hassan Althobaiti, and Ali Alferaidi, "Normal and abnormal heart rates recognition using transfer learning," in *2020 12th International Conference on Knowledge and Systems Engineering (KSE)*. IEEE, 2020, pp. 275–280.
- [18] Agni Kumar, Vikramjit Mitra, Carolyn Oliver, Adeeti Ullal, Matt Bidulph, and Irida Mance, "Estimating respiratory rate from breath audio obtained through wearable microphones," in *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 2021, pp. 7310–7315.
- [19] ANSI/AAMI/IEC, "Ansi/aami ec13-2002: Cardiac monitors, heart rate meters, and alarms," 2002.
- [20] Hui Lu, Julia Beatriz Yip, Tobias Steigleder, Stefan Griebhammer, Maria Heckel, Naga Venkata Sai Jitin Jami, Bjoern Eskofier, Christoph Ostgathe, and Alexander Koelpin, "A lightweight robust approach for automatic heart murmurs and clinical outcomes classification from phonocardiogram recordings," in *2022 Computing in Cardiology (CinC)*. IEEE, 2022, vol. 498, pp. 1–4.