

Two-Stage Multitask-Lerner for PCG Murmur Location Detection

Maurice Rohr, Benedikt Müller, Sebastian Dill, Gökhan Güney, Christoph Hoog Antink

KIS*MED – AI Systems in Medicine,
Technische Universität Darmstadt, Darmstadt, Germany

Abstract

Pre-screening for heart conditions is particularly challenging in low- and middle-income countries due to the lack of expensive equipment and a shortage in medically trained professionals. As heart sounds can be captured easily by smartphones or similar devices, their automated analysis may provide a cost-efficient alleviation of this problem. One potential symptom for cardiac diseases that can be detected through heart sound analysis are so-called heart murmurs.

In this study, we present an approach for detecting heart murmurs that utilizes a Pooling-based Artificial Neural Network (PANN) structure for extracting features from audio waveforms of arbitrary lengths. It can classify single recordings based on recording location and the extracted features in an end-to-end manner. The approach is inspired by the multiple instance learning framework.

We performed a 10-fold stratified cross-validation and report the calculated evaluation measures as average (standard deviation): Murmur weighted accuracy 0.715 (0.077), Outcome-Metric 13640 (2401). The official murmur weighted accuracy and outcome validation score were 0.720, 9135 respectively.

1. Introduction

Cardiovascular diseases are the cause of approximately one third of all deaths globally [1] and a major focus of risk factor analysis and screening. Heart murmurs are indicators of heart diseases and have a high prevalence, yet recognizing them requires strong cardiac auscultation skills, which for many physicians are sub-optimal [2]. Therefore, technical and automated solutions for heart murmur detection are required.

Heart murmurs are essentially audible vibrations caused by perturbations of the blood flow such as strong pressure gradients or velocity changes. Mostly, they arise when heart valves are not opening or closing correctly. Our aim which is also the goal of the George B. Moody Challenge 2022[3], thus was to predict heart murmur from Phonocardiograms (PCG).

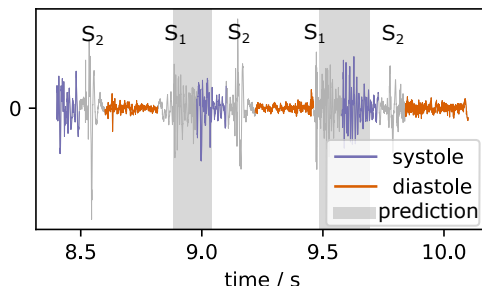


Figure 1. Murmur segmentation generated by MILU-Net on recording 9979_AV [5]. According to the medical labeling there should be systolic murmur. The excitation detected in the model output is marked as grey

PCGs are recordings of all sounds of the heart during a cardiac cycle. This includes sounds such as the first (S1) and second (S2) heart sound, but also murmurs. Commonly a segmentation based on S1/S2 is performed before computing features using e.g. the (mel-)spectrogram for PCG analysis. Khan et al. [4] show that deep learning methods that do not rely on PCG segmentation achieve high accuracy (around 90%) for murmur detection and consistently achieve better results than conventional machine learning algorithms.

By formulating the detection of heart murmurs as a two stage multiple instance learning (MIL) problem with “weakly labeled” data, as known from sound event detection [6], we can show that a machine learning model can find the relevant segments for PCG analysis on its own. Weakly labeled in this context means that for each sound recording, only a single tag is provided without knowing the exact onset and offset times of the relevant sound bites. More specifically in MIL [6], a weakly labeled dataset $D = \{B_n, y_n\}^N$ consists of a set of bags $B_n = \{\mathbf{x}_n^1, \dots, \mathbf{x}_n^{T_n}\}$, where each bag is a collection of instances. T_n is the number of instances in the n -th bag and each instance has a length (duration) d . Given a particular sound class k , a bag is considered *positive* ($y_{k,n} = 1$) if it contains at least one positive instance and *negative* ($y_{k,n} = 0$) if it contains no positive instance at all. Thus,

each PCG can be thought of as a bag of instances where the presence of murmur defines a positive instance. The problem formulation heavily implies that, given a prediction for each instance of a bag, the condensed label is given by the maximum of all instances of the bag. Using the same abstraction on the patient level, the same is true for different recording locations. For both hierarchically ordered sub-problems “weak labels” are available from the CirCor dataset [5].

Our two-stage approach makes it possible to guide the physician to the location where the murmur is most audible, while also giving hints to where in time murmur can be heard in each recording.

2. Methods

We split the methods part in two problems: Primarily we want to construct a model (*MILU-Net*) which, based on weak labels, produces a fine resolution murmur detection which manages to explain the final decision about the presence of murmur for a single recording. Secondly, based on that structure, we design a model (*PANN*) that achieves good murmur prediction accuracy for a set of multiple recordings of a particular patient at the cost of losing interpretability with respect to a single recording. Both approaches rely on the same *pre-processing*.

Pre-processing. The recordings are pre-processed independently of location by removing segments with low signal quality based on signal-to-noise ration and saturation [7] and applying bandpass filtering (10 to 800 Hz). Only during training, all signals are cut or zero padded to a fixed length of 8.2 s to increase efficiency.

MILU-Net. The MILU-Net model consists of a simple U-Net [8] structure for feature generation and a part which facilitates MIL. The U-Net guarantees that its output is the same dimension as the input and thus provides a “strong” label for each sample of the input. It consists of 5 down- and up-sampling convolutional layers with batch norm and ReLU activations and a final output convolutional layer that summarizes the features into a 1D signal. Most importantly, the output of the U-Net is used in a softmax-pooling layer [6] with sigmoid activation to obtain a scalar output murmur probability. By using the softmax-pooling layer, we loosen the MIL assumption, which implies max-pooling. In return, we achieve a more robust training that depends less on the initialization, because the output depends on all instances instead of single particular instances that are chosen randomly due to parameter initialization. The model is overfitted to the available data as the goal is not only to predict the correct label, but also find out if it can learn the unique attributes of murmur. The excitation of the U-Net is then directly related to the relevance and “murmurness” of the respective signal part. An example prediction is shown in Fig. 1.

PANN. The PANN model in Fig. 2 follows the MIL idea loosely by widening the single output signal to an array of feature signals. It employs a convolutional encoder consisting of 6 blocks of 1d-convolutions (kernel size=5, padding=same), batch norm, 1d-convolution, dropout and max pooling (stride=2) to encode the signals in a feature-rich presentation. These features, which can be thought of as time signals, are then processed by adaptive pooling layers which produce a fixed-size output of 30 features in total (max-pooling=15, average-pooling=10, min-pooling=5). The intuition of the convolution block is that it gets activated by murmurs in the respective parts in each segment. The approach then takes advantage of the periodicity of murmurs by employing average and max pooling layers that collect and summarize the information about murmur appearance from all segments of the signal, rendering the output feature dimensions independent of the input length. The output features of the pooling layer are then combined in a fully connected layer (30x64+5 input, 100 hidden and 20 output neurons) with the one-hot encoded recording locations. These outputs are then evaluated in a linear decision layer with softmax activation function. By processing each recording location separately and combining the detection results based on simple decision rules, we enable the user to verify the suspected murmur position.

As depicted in Fig. 3, in a second stage the multi-label model is fed with the features and the encoded output of the PANN model which includes (1) information from statistical features from classical audio processing, (2) the one-hot encoded recording location, and (3) demographic features. The multi-label model is a simple feed forward neural network with 5 feed forward layers with (123, 492, 246, 20) neurons, batch norm after each hidden layer and leaky ReLU activations and one dropout layer at the end. The output layers consist of 3 neurons for murmur prediction and 2 for outcome prediction with softmax activation each.

Augmentation. Due to the small training dataset we employ data augmentation. During training and after pre-processing, one or multiple augmentations are performed at random: *scaling*, *gaussian noise*, *drop*, *cutout*, *shift*, *resampling*, *random resampling*, *sine wave*, *bandpass filtering*. Scaling randomly rescales the signal. Gaussian noise adds gaussian noise to the signal. Drop randomly sets signal values to zero. Cutout randomly sets signal intervals to zero. Shift randomly shifts the signal in time (creating zeros at either end). Random resampling creates smooth time offsets simulating a changing heart rate. Resampling linearly resamples the signal to another sampling frequency, simulating another heart rate. Sine wave adds a random sine wave to the signal. Bandpass filtering randomly applies a bandpass filter between 0.2 and 45 Hz.

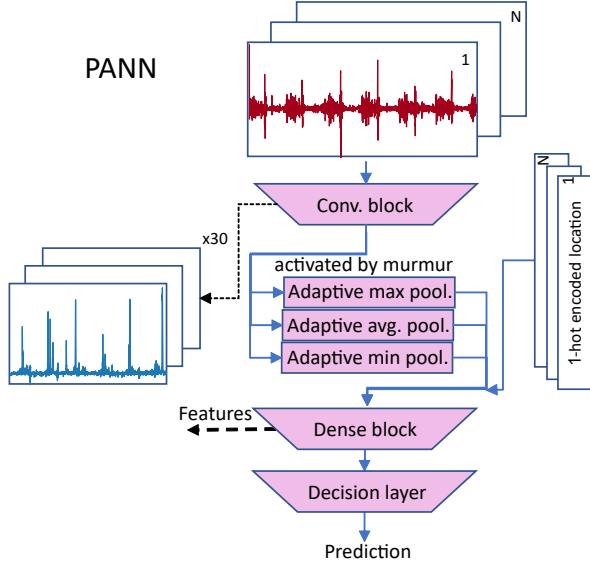


Figure 2. General Structure of PANN Model for murmur detection

Training. For training we employ cross entropy loss for both model stages. First the PANN model is trained using a learning rate of 0.001 which is decreased by factor of 0.3 after each 30 epochs. The best model is picked based on the mean of training and validation accuracy.

3. Results

We test the PANN model by computing the accuracy for each auscultation location separately. The labels are created by considering all recordings of a subject with absent/unknown murmur as “Absent”/“UNK”. If murmur is present, the locations marked as hearable are labeled “Present”, the rest “UNK”. The accuracy of the intermediate decisions of PANN model (Tab. 1) for a confidence threshold of 0.8 show that murmur is detected with an accuracy of $> 90\%$ for all recording locations separately. Given the true label for the subject is “Present”, the accuracy in predicting the correct murmur locations is reduced.

We performed a 5-fold stratified cross-validation for each of the signal augmentations we used during training, by applying each singular augmentation with a probability of 0.25. The effect on the final murmur predictions of each augmentation as well as the validation measures for no augmentation (“none”) is shown in Table 2.

We performed a 10-fold stratified cross-validation of the final model (Fig. 3) resulting in the following scores listed as mean (standard deviation): **Murmur** AUROC 0.831 (0.038), AUPRC 0.657 (0.059), F-measure 0.572 (0.076), Accuracy 0.722 (0.088), Weighted Accuracy 0.715 (0.077); **Outcome** AUROC 0.628 (0.047), AUPRC 0.634 (0.047), F-measure 0.580 (0.056), Accuracy 0.590

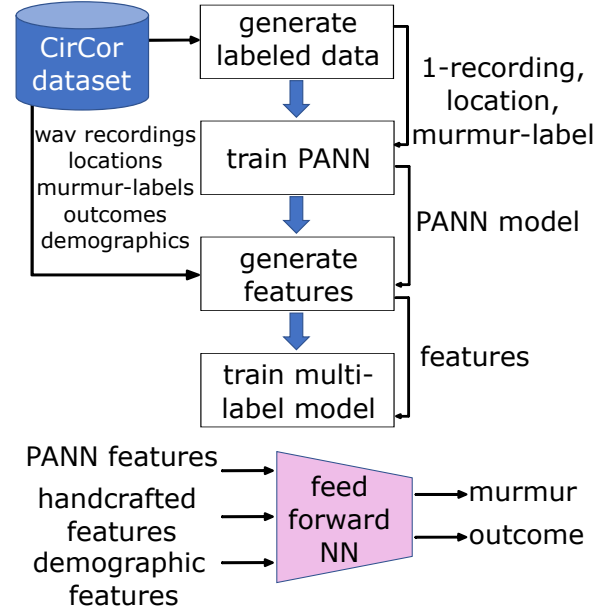


Figure 3. Training of multi label model

Table 1. PANN prediction for presence of murmur per auscultation location [5] given the experts label says “Present” and the total of high signal quality recordings (n=2803).

accuracy	AV	PV	TV	MV	total
present	0.733	0.754	0.800	0.756	0.761
total	0.909	0.911	0.911	0.921	0.913

(0.057), Cost 13640 (2401).

4. Discussion

The task was to learn a high resolution murmur detection given only a weak label for the whole recording. As can be seen in Fig. 1 this works in principle. The main problem is that generalizability is difficult to achieve and often times the murmur locations are not found, which is in part due to the fact, that one case of murmur is enough for the network to be correct about a complete recording. While the accuracy for murmur prediction on the recording level of PANN is quite good ($> 90\%$), for the subject level a rule-based system that followed the MIL idea (regarding murmur as “Present” if there was murmur in at least one location, “Absent” if murmur is absent in all locations and “UNK” otherwise) turned out to score badly (weighted accuracy of 0.56).

Random resampling augmentation appears to reduce the performance of the model significantly in one of the folds while providing no improvement in the others. It might distort the signal sufficiently to make normal signals look similar to murmurs and thus random resampling cannot be recommended in this case. Rescaling the signals or adding

Table 2. 5-fold cross-validation for different augmentations ordered descendingly by improvement on murmur prediction performance on training.

augmentation	AUROC	AUPRC	F-measure
g. noise	0.839 (0.047)	0.664 (0.043)	0.546 (0.110)
scaling	0.840 (0.020)	0.659 (0.050)	0.537 (0.068)
shift	0.835 (0.029)	0.647 (0.035)	0.536 (0.047)
resample	0.832 (0.046)	0.652 (0.036)	0.532 (0.057)
bandpass	0.803 (0.062)	0.632 (0.068)	0.545 (0.043)
cutout	0.791 (0.095)	0.620 (0.082)	0.557 (0.090)
none	0.797 (0.077)	0.616 (0.064)	0.532 (0.080)
sine wave	0.808 (0.062)	0.617 (0.062)	0.514 (0.047)
drop	0.814 (0.031)	0.609 (0.037)	0.500 (0.041)
r. resample	0.782 (0.112)	0.598 (0.127)	0.505 (0.121)

Gaussian noise during training are consistently improving generalization of the model during training. All other augmentation techniques did not result in a significant difference, although more variance in the data certainly helps training more generalized models. Thus, we decided to use them either way but to a lesser extent. Surprisingly, while rescaling reverts normalization of the signals which is standard procedure during training of most machine learning models, it improves generalization. We assume that signals with outliers in amplitudes might be positively affected. Adapting the probability of an augmentation being applied based on ranking table 2 leads to an improvement in all relevant metrics.

The results of the 10-fold stratified cross-validation seem surprisingly low compared to the official validation results. The training of the model so far has shown to be highly sensitive to initialization and the amount of training data. The sensitivity to initialization is an immanent problem of pooling based neural networks [6]. Training based on cost functions for both classification tasks simultaneously reduced the classification accuracy for murmur only by 3.6 %.

5. Conclusion

We present a model based on the MIL network to create fine-granular murmur detection in PCG recordings. While we can show that a basic model generates good results in

finding specific murmur locations only from training on weak labels, this learning does not yet translate to good prediction accuracy on unseen data. An adapted model based on the same pooling ideas but with decreased level of interpretability achieves competitive results in both murmur detection and clinical outcome prediction.

References

- [1] Roth GA, Mensah GA, Fuster V. The global burden of cardiovascular diseases and risks: a compass for global action. *Journal of the American College of Cardiology* 2020; 76(25):2980–2981.
- [2] Vukanovic-Criley JM, Criley S, Warde CM, Boker JR, Guevara-Matheus L, Churchill WH, Nelson WP, Criley JM. Competency in cardiac examination skills in medical students, trainees, physicians, and faculty: a multicenter study. *Archives of internal medicine* 2006;166(6):610–616.
- [3] Reyna MA, Kiarashi Y, Elola A, Oliveira J, Renna F, Gu A, Perez-Alday EA, Sadr N, Sharma A, Mattos S, Clifford GD. Heart Murmur Detection from Phonocardiogram Recordings: The George B. Moody PhysioNet Challenge 2022. *medRxiv* 2022;.
- [4] Khan KN, Khan FA, Abid A, Olmez T, Dokur Z, Khandakar A, Chowdhury ME, Khan MS. Deep learning based classification of unsegmented phonocardiogram spectrograms leveraging transfer learning. *Physiological measurement* 2021; 42(9):095003.
- [5] Oliveira J, Renna F, Costa PD, Nogueira M, Oliveira C, Ferreira C, Jorge A, Mattos S, Hatem T, Tavares T, Elola A, Rad AB, Sameni R, Clifford GD, Coimbra MT. The CirCor DigiScope Dataset: From Murmur Detection to Murmur Classification. *IEEE Journal of Biomedical and Health Informatics* 2022;26(6):2524–2535.
- [6] McFee B, Salamon J, Bello JP. Adaptive pooling operators for weakly labeled sound event detection. *IEEEACM Transactions on Audio Speech and Language Processing* 2018; 26(11):2180–2193.
- [7] Plesinger F, Viscor I, Halamek J, Jurco J, Jurak P. Heart sounds analysis using probability assessment. *Physiological measurement* 2017;38(8):1685.
- [8] Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015; 234–241.

Address for correspondence:

Maurice Rohr
 KIS*MED, TU Darmstadt
 Merckstr. 25, 64283 Darmstadt, Germany
 rohr@kismed.tu-darmstadt.de