

# Classification of Phonocardiogram Recordings using Vision Transformer Architecture

Joonyeob Kim, Gibeom Park, Bongwon Suh

Department of Intelligence and Information, Seoul National University, Seoul, South Korea

## Abstract

*We participated in the George B. Moody PhysioNet Challenge 2022 to make a model which detects the presence or absence of murmurs from multiple heart sound recordings from multiple auscultation locations, as well as detecting the clinical outcomes from phonocardiograms (PCGs) well. Our team, HCCL, developed a model with a visual approach for deriving a high-performance model. The model converts heart sound signals into spectrograms without requiring resampling or signal filtering. The result shows a weighted accuracy score of 0.671 (ranked 29th out of 62 teams) for the murmur detection classification on the hidden validation data. For the clinical outcome identification task on the hidden validation data, it shows a Challenge cost score of 9903 (ranked 23th out of 62 teams)*

## 1. Introduction

Heart conditions and heart disease usually can be diagnosed through echocardiography or electrocardiogram. However, in underdeveloped countries, cardiologists and equipment are scarce, so they can only use PCGs (phonocardiograms) which can be easily measured with a stethoscope. Nevertheless, interpreting sounds requires experts to interpret the results, so developing automated approaches for detecting abnormal heart function from multi-location PCGs recordings of heart sounds at The George B. Moody PhysioNet Challenge 2022 will help diagnose and treat heart conditions [1–3].

In recent deep learning research, applying Transformer architecture [4] has shown excellent results for natural language processing tasks such as BERT [5] as well as for computer vision. In the field of computer vision, image classification models based on Transformer, have achieved better performances than traditional CNN-mixed architectures. In addition, self-attention which is an application of the Transformer model could generate visualization that could help understand the results. ViT (Vision Transformer) [6] is introduced as an image classification model based on Transformer. The model splits an image into mul-

iple patches and uses embedding of individual patches and exhibited excellent performances.

Although deep learning models have been used for medical data analysis in many studies, there has been no case of modeling PCG recordings of heart sounds with a visual transformer so far. Inspired by this idea, we aim to explore the application of ViT architecture to classify PCG recordings with heart murmur patterns. Our final method is composed of time masking, signal segmentation without overlap, and converting a PCG into multiple spectrogram patches. In addition, the self-attention part of the model generates the attention map which could help understand the model’s results. This approach allows us to distinguish heart murmurs recorded in PCGs with higher accuracy as well as visualize the attention of the model on a spectrogram.

## 2. Methods

Our main research goal is to derive a high performance classifier based on the ViT model, which is no need for complex pre-processing of electrocardiogram [7] signals as well as finding visual features. First, our method starts by converting heart sound to spectrogram which visualize the representation of the spectrum of frequencies of a signal varies with time (Fig 1.(a)). Using the spectrogram, we perform additional training using the pre-trained model which was trained on ImageNet-21K datasets, intending to allow the model to learn various features quickly and even with a small amount of data [6]. Second, to address the challenge of small-sized training dataset, we used data augmentation method suitable for the model to avoid overfitting. Third, we pre-process demographic information such as gender and age, to put it into the classifier.

### 2.1. Dataset

The dataset of the PhysioNet/CinC Challenge 2022 contains one or more heart sound recordings for 1568 patients and routine demographic information about the patients [2, 3]. We do not use any extra data other than the training dataset provided by the George B. Moody PhysioNet

Challenge 2022.

## 2.2. Pre-processing

We perform the following three pre-processing. During the process, we did not exclude or relabel the training data because our team has no experts on the given data. Also, no signal filtering or resampling was performed.

### Signal segmentation without overlap

To make spectrograms of identical size, we segment PCG recordings. According to the paper of Raza, A. [8], in the case of PCGs, the best performance was achieved when 12.5s was used as the input of the deep learning model. Based on the observation, we also segment the original sound signal into 12.5s. Next, we convert each segment to spectrogram with the following conversion parameters - window size of 0.11s and the overlap ratio of 0.5. The process produces spectrograms of a  $224 \times 224$  matrix. When the size of the segment is smaller than 12.5s, the image is padded by black pixels, giving the signal the same effect as zero padding. No resampling or filtering was performed to minimize the loss of original data. The ViT model could be used without segmentation accepting full-length signals. However, we got better results when we used segmented signals. It might be due to the data augmentation effect.

Since the pre-trained ViT model accepts  $224 \times 224$ , we create spectrogram images into  $224 \times 224$  pixels to avoid the loss of information.

### Data Augmentation

While model training used only the data provided in the competition, the size of data was insufficient and overfitting occurred. Thus, we tried to solve this problem with data augmentation. Commonly used image augmentation methods are scaling, cropping, flipping, rotation, contrast, and saturation. However, these augmentation methods could compromise the important information in the spectrogram. Therefore, we applied the SpecAugment which is an augmentation method that works on the spectrogram of input signal [9]. SpecAugment includes various augmentation methods such as frequency masking, time masking, fade in and fade out. In the hidden validation data, only using time masking showed a higher score than any augmentation.

### Demographic information

Our team used age, sex, height, weights, pregnancy information of the challenge dataset. In case of categorical data, it was converted using a label encoder, and each encoder was saved and used for the conversion of test data. For missing values, they were classified and converted into a new class. In the case of numeric data, it was transformed

using the MinMax scaler, and each scaler was saved and used for the transformation of the test data. For a missing value, the mode was used for filling up the missing value. These values were also used for the test dataset.

## 2.3. Model

We used the ViT model of which patch size is  $16 \times 16$  pixel and image size is  $224 \times 224$  pixel and pre-trained on ImageNet-21K (14 million images, 21843 classes) as the baseline model [6]. We modified the classifier portion to add demographic information to the model. We used most of the same hyperparameters for the murmur task and clinical outcome task such as an initial learning rate of 0.0001 with AdamW and LambdaLR scheduler. In case of murmur classification, batch size was 64, and saving steps and evaluation steps were 100. However, we used batch size of 32, and saving steps and evaluation steps were 50 for outcome classification. The larger the batch size of ViT, the better the performance, so the maximum possible value in the test environment was used, and in the case of step, the optimal value was found experimentally.

### Vision Transformer

ViT splits the image into small patches and performs image classification. The ViT model shows higher performance with less data during fine-tuning than CNN-based models. Furthermore, it has the advantage of self-attention that can generate an attention map that could help understand the model's results. In the attention mask (Fig 1.(b)), the yellow parts which are bright than the other parts are the attention part. It mainly paid attention to low frequencies, but occasionally attended to high frequencies, and this information might be helpful for classification.

We used pre-trained with ImageNet-21K model because it is impossible to predict good performance without pre-training with a large amount of data similar to Transformer.

### Model architecture

Our model is as shown in Fig 1.(b). Pre-processed heart sound signals are converted into images, and it is input to ViT Feature extractor. As the output of ViT feature extractor, the encoded feature is output in batch units, and it is input to the Transformer encoder of ViT for Image Classification. The input features go through a total of 12 ViT layers, and the pre-processed demographic information features are concatenated into the output. The concatenated feature is input to the last layer, the classifier. Finally, the class is learned and inference is performed.

### Weighted categorical cross-entropy loss

Cross-entropy is used to measure the difference between different probability distributions and is used as a loss function for classification in deep learning and machine

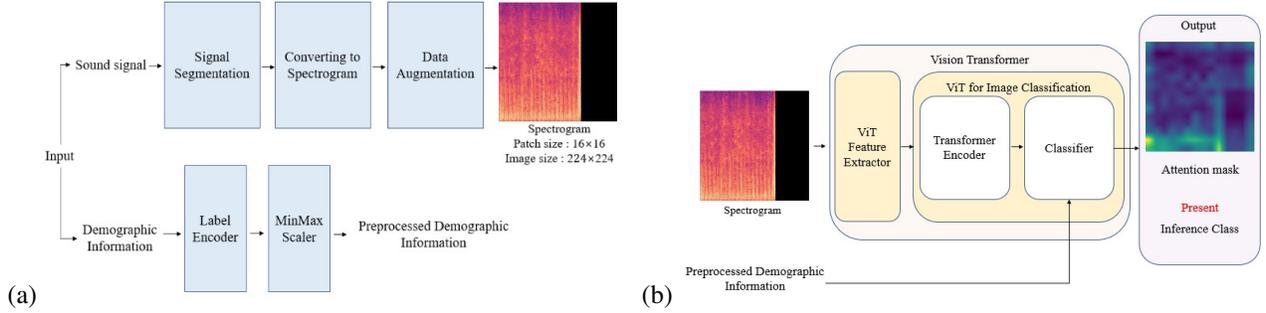


Figure 1. (a) Pre-processing overview. (b) Model overview.

learning. Since the class of the dataset of the challenge is imbalanced, a weighted categorical cross-entropy loss is applied. The weighted categorical cross-entropy loss is shown as (1):

$$\mathcal{L} = -\frac{1}{M} \sum_{k=1}^K \sum_{m=1}^M y_m^k \times w_k \times \log(h_\theta(x_m, k)) \quad (1)$$

where  $M$  is number of training examples,  $K$  is number of classes (3 at murmur, 2 at outcome),  $y_m^k$  is target label for class  $k$  at each training example  $m$ ,  $w_k$  is weight (1, 5, 3 for Absent, Present, Unknown class at murmur, 1.2, 1 for Ab-normal, Normal at outcome),  $h_\theta$  is model with weights  $\theta$ ,  $x_m$  is input for training example  $m$  [10]. Instead of using the ratio of label, the weight ratio of 1:5:3 was determined using the ratio of competition scores. To avoid the learning from being concentrated on the Normal class, the experiment was repeated, and an appropriate value of 1.2:1 was established.

### 3. Results

Our baseline experiment is trained with ViT model which is pre-trained with ImageNet-21K after converting a full-length signal into a  $473 \times 473$  spectrogram without segmentation or augmentation and resize into a  $224 \times 224$  image. we convert each segment to spectrogram with the following conversion parameters - window size of 0.236s and the overlap ratio of 0.5. Additionally, our final model contains signal segmentation, data augmentation, and demographic data as well as pre-trained weights.

The results of an ablation study for murmur detection and outcome detection are reported in Table 1 and Table 2, respectively, along with the consequences. The majority of preprocessing techniques enhanced performance. Performance was most influenced by pre-trained, while signal segmentation and demographic data had negligible effects on both tasks. The results of the murmur detection task for the training set and hidden validation set are shown in Table 3, and the results of the clinical outcome identification task are shown in Table 4.

Methods	Weighted accuracy
Baseline	$0.708 \pm 0.045$
Final model without Segmentation	$0.751 \pm 0.025$
Final model without Pretrained	$0.661 \pm 0.016$
Final model without Augmentation	$0.704 \pm 0.042$
Final model without Demographic	$0.744 \pm 0.025$
Final model	$0.759 \pm 0.030$

Table 1. Results of our final model’s ablation study of murmur detection using 5-fold cross validation.

Methods	Challenge cost
Baseline	$15507 \pm 5485$
Final model without Segmentation	$14800 \pm 4380$
Final model without Pretrained	$23514 \pm 4976$
Final model without Augmentation	$13726 \pm 912$
Final model without Demographic	$14283 \pm 1821$
Final model	$13091 \pm 1183$

Table 2. Results for our final model’s ablation study of outcome detection using 5-fold cross validation.

Training	Validation	Test	Ranking
$0.759 \pm 0.030$	0.671		29/62

Table 3. Weighted accuracy metric scores (official Challenge score) for our final selected entry (team HCCL) for the murmur detection task, including the ranking of our team on the hidden test set. We used 5-fold cross validation on the public training set, repeated scoring on the hidden validation set, and one-time scoring on the hidden test set.

Training	Validation	Test	Ranking
$13091 \pm 1183$	9903		23/62

Table 4. Cost metric scores (official Challenge score) for our final selected entry (team HCCL) for the clinical outcome identification task, including the ranking of our team on the hidden test set. We used 5-fold cross validation on the public training set, repeated scoring on the hidden validation set, and one-time scoring on the hidden test set.

## 4. Discussions

We studied a model that can detect murmur and pathological outcome by applying the visual approach with the ViT model to PCGs modeling. Contrary to expectations, in the murmur detection task, the score of the hidden validation data was lower than the training data. Nevertheless, in the outcome identification task, the score of the hidden validation data was higher than the training data. We presume that the model overfits the murmur task and requires more generalizations for improvement, and the cause might be a small-sized dataset. In the case of signal segmentation, the training data showed a score improvement of 3.3%, but the hidden validation data showed a score improvement of 13.3%, which showed that signal segmentation was effective in the PCGs murmur detection for the ViT model.

Due to our team's limited expertise in the dataset, there might be a better optimized pre-processing technique. Although we heard PCGs, detailed information could not be obtained due to a lack of pathological understanding of the data. Potential anomalies could have not been detected in the demographic information. If exclusion criteria or re-labeling is performed through an expert on data, or additional data is utilized, the model could have performed better.

One potential benefit of the proposed method is that it allows us to examine an attention map giving opportunities what part of PCGs are informative for each pathological condition. We believe that this could open up opportunities for experts to understand, interpret, and improve the model's findings.

## Acknowledgments

TBD

## References

- [1] Goldberger AL, Amaral LA, Glass L, Hausdorff JM, Ivanov PC, Mark RG, et al. Physiobank, Physiotookit, and Physionet: Components of a new research resource for complex physiologic signals. *Circulation* 2000;101(23):e215–e220.
- [2] Reyna MA, Kiarashi Y, Elola A, Oliveira J, Renna F, Gu A, et al. Heart murmur detection from phonocardiogram recordings: The George B. Moody PhysioNet Challenge 2022. *medRxiv* 2022;URL <https://doi.org/10.1101/2022.08.11.22278688>.
- [3] Oliveira J, Renna F, Costa PD, Nogueira M, Oliveira C, Ferreira C, et al. The CirCor DigiScope dataset: from murmur detection to murmur classification. *IEEE Journal of Biomedical and Health Informatics* 2021;26(6):2524–2535.
- [4] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. *Advances in neural information processing systems* 2017;30.
- [5] Devlin J, Chang MW, Lee K, Toutanova K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv181004805* 2018;.
- [6] Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv201011929* 2020;.
- [7] Naz M, Shah JH, Khan MA, Sharif M, Raza M, Damaševičius R. From ecg signals to images: a transformation based approach for deep learning. *PeerJ Computer Science* 2021;7:e386.
- [8] Raza A, Mehmood A, Ullah S, Ahmad M, Choi GS, On BW. Heartbeat sound signal classification using deep learning. *Sensors* 2019;19(21):4819.
- [9] Park DS, Chan W, Zhang Y, Chiu CC, Zoph B, Cubuk ED, et al. SpecAugment: A simple data augmentation method for automatic speech recognition. *arXiv preprint arXiv190408779* 2019;.
- [10] Ho Y, Wooley S. The real-world-weight cross-entropy loss function: Modeling the costs of mislabeling. *IEEE Access* 2019;8:4806–4813.

Address for correspondence:

Bongwon Suh

Seoul National University, 1, Gwanak-ro, Gwanak-gu, Seoul 08826, South Korea

bongwon@snu.ac.kr