

# On the Analysis of Data Augmentation Methods for Spectral Imaged Based Heart Sound Classification Using Convolutional Neural Networks

George Zhou (✉ [gez4001@med.cornell.edu](mailto:gez4001@med.cornell.edu))

Weill Cornell Medicine

Yunchan Chen

Weill Cornell Medicine

Candace Chien

Weill Cornell Medicine

---

## Research Article

**Keywords:** Machine learning, data augmentation, cardiac sound analysis, spectrograms, convolutional neural network, cardiology, healthcare automation

**Posted Date:** September 20th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-888104/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **Title:** On the Analysis of Data Augmentation Methods for Spectral Imaged Based Heart  
2 Sound Classification using Convolutional Neural Networks

3 **Authors:** George Zhou<sup>1\*</sup>, Yunchan Chen<sup>1</sup>, Candace Chien<sup>1</sup>

4 **Affiliations**

5 **1** Weill Cornell Medicine, New York, NY, 10021, USA

6 **Corresponding Author:** George Zhou, [gez4001@med.cornell.edu](mailto:gez4001@med.cornell.edu)

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26 **Abstract**

27 **Background:** The application of machine learning to cardiac auscultation has the potential to  
28 improve the accuracy and efficiency of both routine and point-of-care screenings. The use of  
29 Convolutional Neural Networks (CNN) on heart sound spectrograms in particular has defined  
30 state-of-the-art performance. However, the relative paucity of patient data remains a significant  
31 barrier to creating models that can adapt to the wide range of between-subject variability. To  
32 that end, we examined a CNN model's performance on automated heart sound classification,  
33 before and after various forms of data augmentation, and aimed to identify the most optimal  
34 augmentation methods for cardiac spectrogram analysis.

35 **Results:** We built a standard CNN model to classify cardiac sound recordings as either normal  
36 or abnormal. The baseline control model achieved an ROC AUC of  $0.945 \pm 0.016$ . Among the  
37 data augmentation techniques explored, horizontal flipping of the spectrogram image improved  
38 the model performance the most, with an ROC AUC of  $0.957 \pm 0.009$ . Principal component  
39 analysis color augmentation (PCA) and perturbations of saturation-value (SV) of the hue-  
40 saturation-value (HSV) color scale achieved an ROC AUC of  $0.949 \pm 0.014$  and  $0.946 \pm 0.019$ ,  
41 respectively. Time and frequency masking resulted in an ROC AUC of  $0.948 \pm 0.012$ . Pitch  
42 shifting, time stretching and compressing, noise injection, vertical flipping, and applying random  
43 color filters all negatively impacted model performance.

44 **Conclusion:** Data augmentation can improve classification accuracy by expanding and  
45 diversifying the dataset, which protects against overfitting to random variance. However, data  
46 augmentation is necessarily domain specific. For example, methods like noise injection have  
47 found success in other areas of automated sound classification, but in the context of cardiac  
48 sound analysis, noise injection can mimic the presence of murmurs and worsen model  
49 performance. Thus, care should be taken to ensure clinically appropriate forms of data  
50 augmentation to avoid negatively impacting model performance.

51 **Key Words:** Machine learning, data augmentation, cardiac sound analysis, spectrograms,  
52 convolutional neural network, cardiology, healthcare automation

53

54

## 55 **I. Background**

56 Cardiac auscultation has been a core element of the cardiovascular physical exam since the  
57 1800s. Sounds produced by the heart reflect its underlying biology and can cue a trained  
58 physician to different heart pathologies such as valvular defects or congenital diseases.

59 However, in recent years, cardiac auscultation has been challenged for its diagnostic utility. The  
60 decline in accurate cardiac auscultation is a well-documented phenomenon<sup>1,2,3</sup>. For example,  
61 internal medicine residents in the US made a correct assessment of auscultation findings only  
62 22% of the time<sup>2</sup>.

63

64 This has spurred an active area of research in developing suitable machine learning models to  
65 classify heart sounds based on recorded phonocardiogram (PCG) signals. Many research  
66 groups have published a wide variety of machine learning models to this end. Survey of the  
67 existing literature reveals that many different feature extraction methods (Mel-frequency cepstral  
68 coefficients<sup>4,5,6</sup>, discrete wavelet transform<sup>7,8,9</sup>, tensor decomposition<sup>10</sup>, sparse coding<sup>11</sup>) and  
69 classification methods (k-nearest neighbors<sup>7</sup>, support vector machines<sup>4,10,11,12</sup>, hidden Markov  
70 models<sup>13,14</sup>, recurrent neural networks<sup>15,16</sup>, convolution neural networks<sup>6,17,18</sup>), and their different  
71 permutations together have been extensively explored.

72

73 It is generally accepted that bigger datasets result in better machine learning models<sup>19,20</sup>.

74 However, real-world clinical applications is limited by the scarcity of labeled clinical data. This  
75 scarcity issue can be attributed to several challenges unique to the medical domain, including:  
76 the relative paucity of available clinical databases structured for machine learning research, the

77 administrative and logistical hurdles associated with collecting and working with patient data and  
78 protected health information due to Health Insurance Portability and Accountability Act (HIPAA)  
79 laws and Institutional Review Board (IRB) regulations, and finally the time-consuming and  
80 expensive nature of properly annotating health data. The gold standard for validating heart  
81 sounds is echocardiogram imaging plus the diagnosis from a cardiologist, both of which are  
82 costly to obtain. An additional challenge in creating a machine learning model to classify heart  
83 sounds is that heart sounds are not actually recorded and stored anywhere in electronic health  
84 records (EHR). Mining EHR databases is not an option, meaning heart sounds must be  
85 collected and labeled from scratch, one-by-one. Data acquisition is made even harder in times  
86 of public health crises, as we have observed with the COVID-19 pandemic, which resulted in  
87 drastic reductions in non-emergency patient volumes in clinics across the world.

88

89 Data augmentation is one solution to the legal limitations and constraints around clinical data.  
90 Data augmentation is the process of generating *synthetic* data from *real* data, while preserving  
91 the class label. In the context of developing machine learning models for heart sound  
92 classification, *real* data means heart sounds collected directly from a patient, whereas *synthetic*  
93 data means artificial heart sounds generated from *real* heart sounds via various computer-  
94 implemented methods.

95

96 The major value add of data augmentation for heart sound classification resides in its ability to  
97 significantly expand the size of available training data without the onerous task of having to  
98 actually obtain and label a large enough volume of heart sounds. An expanded dataset can  
99 improve model performance because the new data created from class-preserving  
100 transformations can help the model better learn the unique features that constitute the essence  
101 of a class, instead of the random variance that is present within each class. Data augmentation  
102 combats overfitting and can help the model make better predictions on unseen data.

103

104 Data augmentation is necessarily domain specific, as the applied transformations should reflect  
105 realistic variations and preserve the underlying features that distinguish different classes from  
106 each other. In other words, the data augmentation should 'make sense' for the task at hand.

107 Two important constraints unique to heart sound spectrograms must be considered in designing  
108 effective data augmentation strategies.

109

110 The first constraint, which we will call the "physiological constraint", is related directly to the  
111 phenomenon under study, the heart sound itself. Heart sounds naturally fall within a narrow  
112 physiological scope: heart rates are 60-100 beats per minute and the principal frequencies of  
113 heart sounds are 20 to 500 Hz. A healthy heart sound can be deconstructed into four main  
114 frequency components: S1 (mitral and tricuspid valve closing), systole (ventricles contracting),  
115 S2 (aortic and pulmonic valve closing), and diastole (ventricles relaxing). A pathological heart  
116 sound has all the same frequency components. The difference between a healthy heart sound  
117 and pathological heart sound is that a pathological heart sound will have additional frequency  
118 components such as murmurs from valve stenosis or regurgitation, rubs from pericarditis, S3  
119 gallops(from increased atrial pressure, as seen in congestive heart failure or dilated  
120 cardiomyopathy), or S4 gallops(atrium contracting against stiff ventricle caused by hypertension,  
121 pulmonary hypertension, ventricular outflow obstruction, or ischemic heart disease). Of note, an  
122 additional sound that can be produced by a healthy heart is the physiological splitting of S2 due  
123 to delayed pulmonic valve closing. Thus, the "physiologic constraint" is that any data  
124 augmentation method must reflect realistic variations of possible heart sounds and also ensure  
125 the presence or absence of additional frequency components is preserved for each individual  
126 heart sound or else the distinguishing factor between a normal and abnormal heart sound is lost  
127 and the class labels lose their meaning.

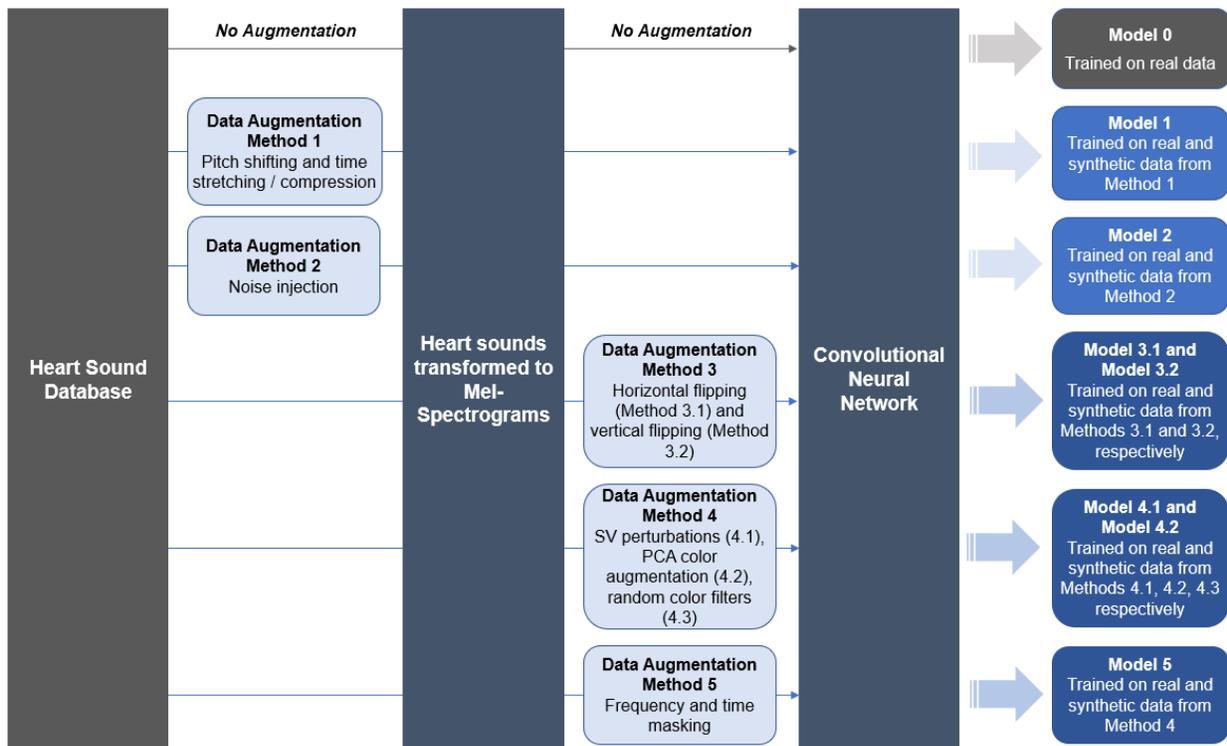
128

129 The second constraint, which we will call the “spectrogram constraint”, is related to the  
130 spectrogram image and what it represents. One advantage for using CNN to classify heart  
131 sounds is that this converts an audio classification problem into a computer vision problem,  
132 which opens the door to the extensive library of data augmentation techniques developed for  
133 images. *Shorten et al.*<sup>21</sup> published a review article surveying the gamut of image data  
134 augmentation techniques that have been researched including flipping, cropping, rotation,  
135 translations, color space transformations, kernel filters to sharpen or blur images, mixing  
136 images, and random erasing. However, not all image data augmentation techniques will  
137 translate appropriately. Although spectrograms are images from a data structure point of view,  
138 spectrograms and traditional images have a fundamental difference in terms of what information  
139 is conveyed along the x- and y- axis. For a traditional image, the axes represent physical  
140 distances, while for spectrograms the x-axis represents time and the y-axis represents  
141 frequency. Moreover, color also carries a different meaning for traditional images vs  
142 spectrogram images. The meaning of color is self-evident for traditional images. For  
143 spectrograms, color is an additional dimension that represents decibels, or the loudness and  
144 intensity of the heart sound. Thus, the “spectrogram constraint” is that any data augmentation  
145 method that operates on the spectrogram as a simple image should correlate with a real-world,  
146 physical transformation of the sound.

147

148 With these constraints in mind, we evaluate common data augmentation techniques at the audio  
149 level, including pitch shifting and time stretching/compressing and noise injection, and at the  
150 image level, including horizontal flips, vertical flips, hue/brightness transformations, principal  
151 component analysis (PCA) color augmentation, random color filters, and time/frequency  
152 masking, for classification of heart sounds based on their spectral image. We include  
153 augmentation methods that are consistent with and contradict what would be an effective data  
154 augmentation method as predicted by our theoretical considerations discussed above to 1)

155 examine the individual effectiveness of each augmentation technique on heart sound  
 156 classification and 2) assess the validity of our theoretical framework.  
 157  
 158 To study the effects of these data augmentation methods on heart sound classification, we  
 159 separate our experiments into two phases. The first phase is to establish the baseline  
 160 performance of our CNN on spectral images of heart sounds. In the second phase, the same  
 161 CNN is trained on both real and synthetically generated heart sounds. Model performance with  
 162 and without data augmentation on the same binary classification task is compared. Each  
 163 individual data augmentation scheme is carried out in a one-to-one correspondence, meaning  
 164 for every real heart sound, one synthetic heart sound is generated from it. This doubles the size  
 165 of the dataset available for training, from  $N$  to  $2N$ . Figure 1 below shows our study design.  
 166  
 167



168  
 169 **Figure 1: Overview of Study Design**

170 To study the effects of data augmentation on heart sound classification, we established the  
171 baseline performance of a machine learning algorithm trained on real heart sound data only  
172 (Model 0). We then compared this baseline performance to various models as delineated in the  
173 above diagram.  
174

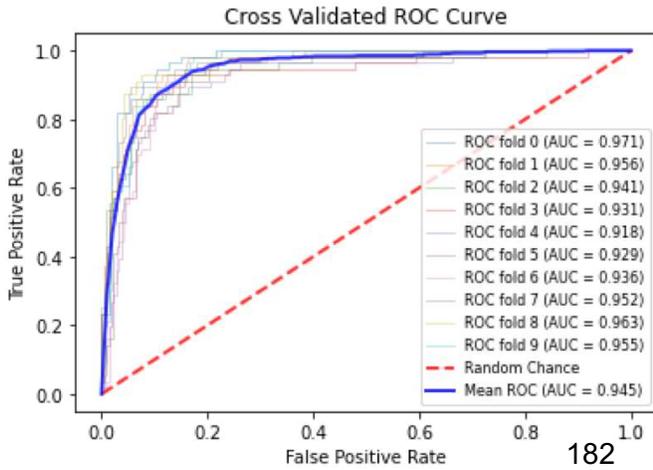
## 175 **II. Results**

176 Reported metrics are based on a stratified 10-fold cross validation. The folds are created in a  
177 consistent way across the different models. This serves to limit any potential variability in model  
178 performance that would be due to the underlying variability in the data itself. Test folds only  
179 contain real heart sounds.

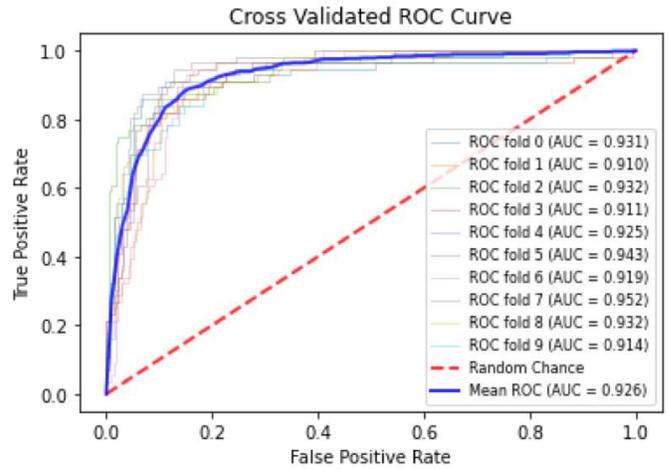
180

181 Figure 2 shows the cross validated ROC curves for the different models.

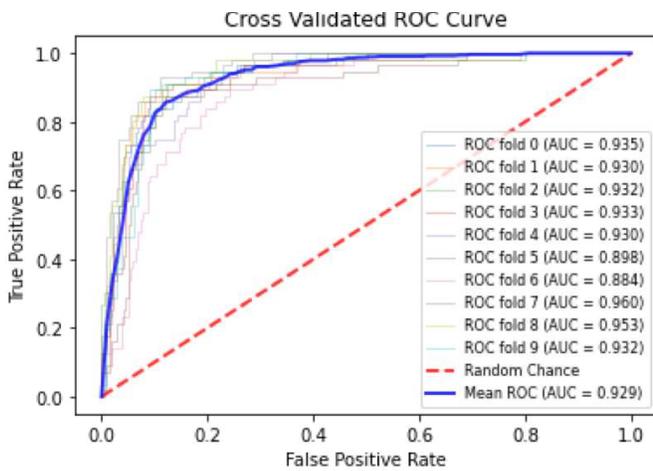
**Fig. 2a. Model 0**  
**Baseline**



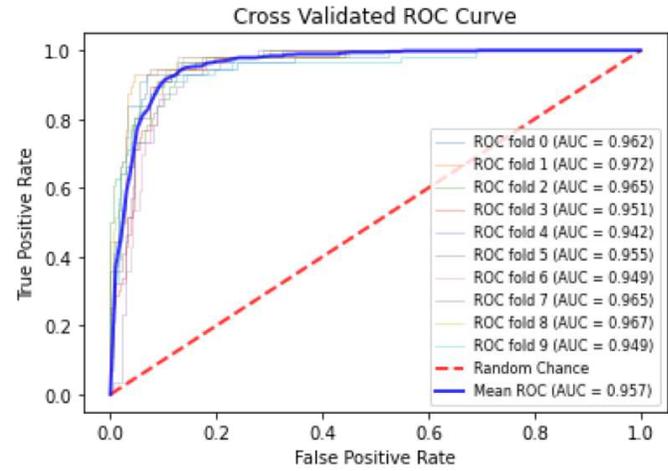
**Fig. 2b. Model 1**  
**Pitch Shifting / Time Compression**



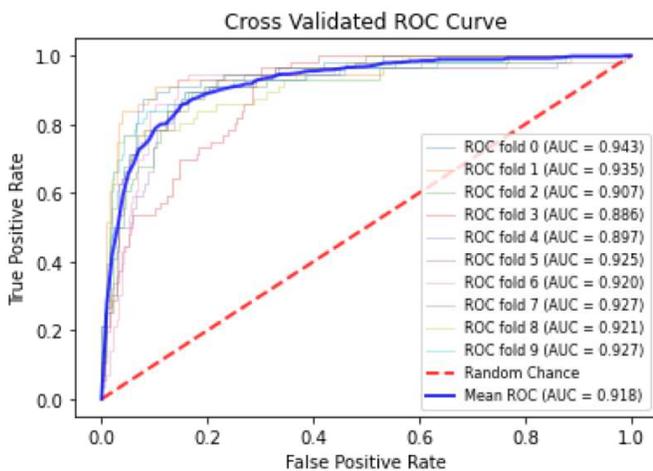
**Fig. 2c. Model 2**  
**Noise Injection**



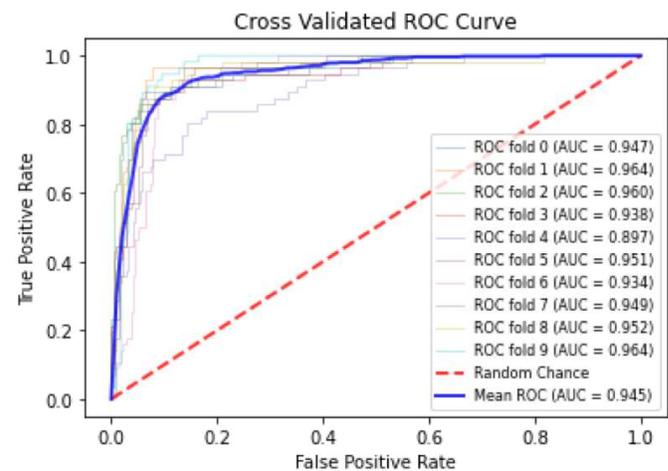
**Fig. 2d. Model 3.1**  
**Horizontal Flip**



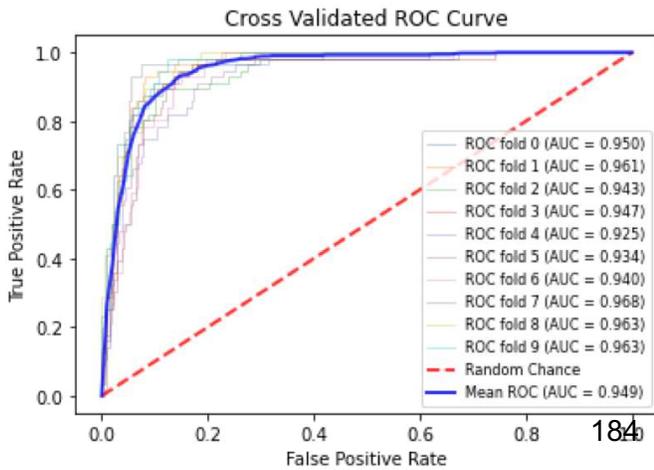
**Fig. 2e. Model 3.2**  
**Vertical Flip**



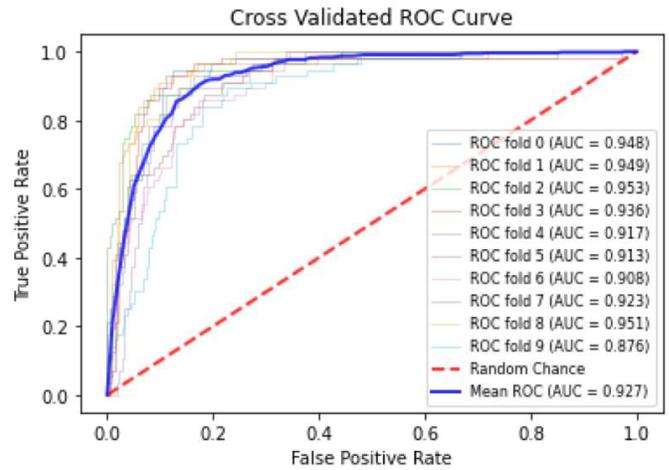
**Fig. 2f. Model 4.1**  
**SV Perturbations**



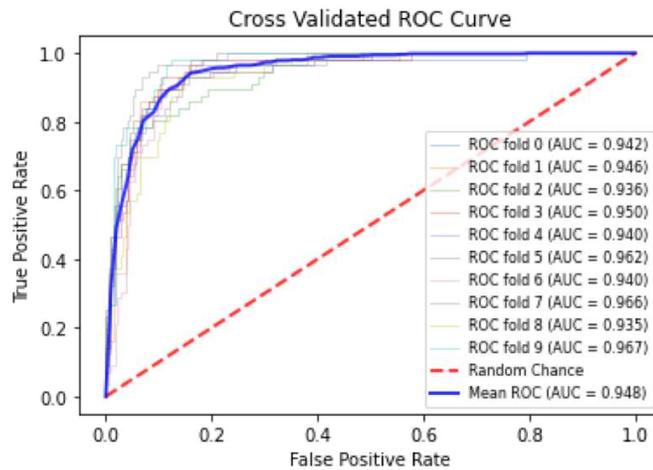
**Fig. 2g. Model 4.2**  
**PCA Color Augmentation**



**Fig. 2h. Model 4.3**  
**Random Color Filters**



**Fig. 2i. Model 5**  
**Time and Frequency Masking**



**Figure 2: ROC curves for Model 0 (a), Model 1 (b), Model 2 (c), Model 3 (d,e), Model 4 (f,g,h), Model 5 (i)**

Comparison of the ROC curve for Model 0, trained on real data only (2a); ROC curve for Model 1, trained on Mel-Spectrograms of real plus pitch shifted and time stretched/compressed heart sounds (2b); ROC curve for Model 2, trained on Mel-Spectrograms of real plus noise injected heart sounds (2c); ROC curves for Models 3.1 and 3.2, trained on real and horizontally flipped Mel-Spectrograms (2d), and real and vertically flipped Mel-Spectrograms (2e); the ROC curves for Model 4.1, 4.2, and 4.3, trained on real and saturation/value transformed images (2f), real and multi-color transformed Mel-Spectrograms (2g), real and PCA color augmented Mel-Spectrograms (2h); and the ROC curve for Model 5, trained on real and frequency/time masked Mel-spectrograms (2i). The dotted red line represents the no-discrimination line.

Table 1 is a numerical summary of the performance of each model.

207  
208  
209

**Table 1. Average performance of each model according to accuracy, sensitivity, specificity and the ROC AUC**

	<b>Accuracy (<math>\pm Stdev</math>)</b>	<b>Specificity (<math>\pm Stdev</math>) (at 90% Sensitivity)</b>	<b>ROC AUC (<math>\pm Stdev</math>)</b>
<b>Model 0</b> <i>Baseline</i>	89.7% (1.7)	86.6% (3.8)	0.945 (0.016)
<b>Model 1</b> Pitch/time alterations	88.2% (2.4) ↓	82.3% (4.7) ↓	0.926 (0.013) ↓
<b>Model 2</b> Noise Injection	88.6% (2.1) ↓	82.2% (6.2) ↓	0.929 (0.021) ↓
<b>Model 3.1</b> Horizontal Flip	90.2% (1.8) ↑	90.8% (2.7) ↑	0.957 (0.009) ↑
<b>Model 3.2</b> Vertical Flip	89.2% (2.7) ↓	79.5% (6.9) ↓	0.919 (0.017) ↓
<b>Model 4.1</b> SV Perturbations	90.6% (1.7) ↑	80.3% (26.9) ↓	0.946 (0.019) ↑
<b>Model 4.2</b> PCA Color Augmentation	89.2% (2.2) ↓	87.8% (4.3) ↑	0.949 (0.014) ↑
<b>Model 4.3</b> Random Color Filters	87.4% (3.0) ↓	81.4% (7.0) ↓	0.927 (0.024) ↓
<b>Model 5</b> Time/Frequency Masking	89.5% (1.7) ↓	86.2% (5.1) ↓	0.948 (0.012) ↑

210  
211

### 212 III. Discussion

213 In summary, our objective was to identify the optimal forms of data augmentation for the binary  
214 classification of PCG signals using their spectral image representation. Our baseline CNN  
215 model achieved specificity of 86.6% at 90% sensitivity, and a ROC AUC of 0.95, which makes it  
216 comparable to state-of-the-art<sup>30, 31</sup>. As previously discussed, one of the unique challenges of  
217 heart sound augmentation is that the generated samples must fulfill certain “physiological  
218 constraints” to remain meaningful. More explicitly, the rate, rhythm, and pitch of cardiac sounds  
219 are bounded within a narrow range. Values that fall outside of these limits would be unrealistic,

220 and hence detract from the classification. Additionally, the original spectral components of the  
221 heart sounds must be maintained to ensure that a normal sound does not become pathological.  
222 The presence or absence of frequency components like murmurs, rubs, S3, or S4 gallops  
223 should be preserved through these transformations. Secondly, the “spectrogram constraint”  
224 stems from the fact that spectrograms and photographs fundamentally convey different  
225 information along their respective dimensions. Image data augmentation methods can work for  
226 spectral images only if they correlate with realistic physical variations in the sound.

227

228 The data augmentation method that satisfied both the “physiological constraint” and the  
229 “spectrogram constraint” improved model performance, while all the data augmentation  
230 methods that failed to satisfy at least one of the constraints worsened model performance in  
231 some respect, experimentally supporting our theoretical framework. We provide a rationale for  
232 why each data augmentation method either improved, did not effect, or worsened model  
233 performance using our framework below.

234

235 The first augmentation method was pitch shifting and time stretching/compressing. Since this  
236 augmentation is done at the audio level, the “spectrogram constraint” does not apply. Natural  
237 pitch variations reflect different anatomical variations of the heart including differing myocardium  
238 wall thickness, body fat/water composition, patient bone/rib structure, and the actual heart size,  
239 all of which may lead to variabilities in heart sound attenuation. The data augmentation  
240 technique of pitch shifting aims to capture these natural variations. There is also variability in  
241 how fast the heart beats. Time stretching and compressing represents heart sounds at different  
242 heart rates, such as in tachycardia or bradycardia. Although pitch shifting and time  
243 stretching/compressing as data augmentation techniques reflects possible physiological  
244 variations, experimentally we see worsening model performance when these data augmentation  
245 techniques are applied. At first this seems to contradict our theoretical framework because the

246 “physiological constraint” is supposedly satisfied. However, if we considered that the natural  
247 heart sound exists within a very narrow physiological range, it is likely that the upper and lower  
248 limits of our pitch shifting, and time stretching/ compressing may have pushed the audio outside  
249 the normal physiological range. Thus, the “physiological constraint” was not actually satisfied  
250 because our augmentation techniques created sounds that would never exist clinically, which is  
251 consistent with the worsening model performance.

252

253 The second augmentation method was noise injection. Noise injection has a regularization  
254 effect that can improve model performance by reducing overfitting and is a widely used audio  
255 data augmentation method for improving model performance. This augmentation is also done at  
256 the audio level, so again the “spectrogram constraint” does not apply. Despite the known ability  
257 of noise injection for improving model performance, we observe that noise injection actually  
258 worsens model performance for heart sound spectral image classification. This can be  
259 understood from the fact that the fundamental difference between normal and abnormal heart  
260 sounds is that the latter has additional frequency components (murmurs, rubs, S3 gallops, S4  
261 gallops). By definition, noise injection is the act of introducing new frequency components to an  
262 audio file. Thus, noise injection is essentially converting normal heart sounds into abnormal  
263 heart sounds. Noise injection fails to satisfy the “physiological constraint” because it ruins the  
264 distinction that separates normal and abnormal heart sounds.

265

266 The third augmentation method is flipping the spectrogram image. Horizontal flipping improved  
267 model performance on all three counts, while vertical flipping worsened model performance on  
268 all three counts. This is explained by the fact that information conveyed by sound is encoded in  
269 the frequency domain, which is represented on the y-axis of spectrogram images. This is an  
270 important distinction from traditional images, where the y-axis represents a physical distance.  
271 Although vertical flipping has been shown to be an effective augmentation technique for

272 improving model performance on many image datasets such as ImageNet and CIFAR-10<sup>32</sup>  
273 (which consist of images of commonplace objects like dogs, cats, cars, etc.), a vertical flip is not  
274 appropriate for a spectrogram image. Transformations of the y-axis of spectrograms would  
275 scramble the frequency content of the sound, rendering any meaningful information that was  
276 encoded in the sound to be lost. A vertical flip has no physical correlation, and so does not  
277 satisfy the “spectrogram constraint.” In fact, the vertical flip worsened model performance the  
278 most out of all the data augmentation techniques explored, underscoring the importance of not  
279 distorting the y-axis of spectrogram images. Horizontal flipping leaves the frequency axis intact,  
280 so it satisfies the “spectrogram constraint”. A horizontal flip alters the temporal relationships of  
281 the frequency components, but as discussed above, a normal and pathological heart sound  
282 mostly contain the same frequency components (S1, S2, systole, diastole). The major difference  
283 is the presence or absence of other frequency components such as murmurs. It is not so much  
284 the temporal relationship of these frequency components with each other that help discern a  
285 normal heart sound from a pathological one. Thus, horizontal flips satisfy the “physiological  
286 constraint” as well, and experimentally we observe that horizontal flips improve model  
287 performance the most out of all data augmentation methods explored. Horizontal flipping as a  
288 data augmentation technique is most likely unique to heart sound spectral images compared to  
289 many other audio classification problems that represent sound as spectral images, owing to the  
290 rhythmic nature of heart sounds. In other audio classification tasks such as speech recognition,  
291 the temporary relationship of the different frequency components is important, and thus a  
292 horizontal flip would most likely hinder model performance.

293

294 The next set of data augmentation methods (methods 4.1, 4.2, and 4.3) are various color space  
295 transformations. Although these transformations do not distort the frequency axis of the  
296 spectrogram, it is important to keep in mind the role of color as an additional dimension in  
297 spectrogram images. In a regular photo, color represents the wavelength of light reflecting off an

298 object. In a spectrogram, color represents the loudness/intensity of the signal measured in  
299 decibels. Factors that contribute to the natural variation in heart sound amplitudes (i.e. how loud  
300 the heart sound is) include the size and position of the heart in the mediastinum, the presence  
301 of fluid within or fibrous thickening of the pericardium, and the position and extent of aeration of  
302 the lungs. For example, heart sounds are usually loudest at the apex where the heart is in direct  
303 contact with the anterior wall of the thorax. Younger patients tend to have louder heart sounds  
304 due to elastic and thin chest walls, whereas older patients tend to have quieter heart sounds  
305 due to stiffer and thicker chest walls. Heart sounds are louder when the patient is in full  
306 expiration, and quieter when the patient is in full inspiration. The data augmentation technique of  
307 color space transformations aims to capture these variations. Experimentally, we observe that  
308 SV (method 4.1) and PCA (method 4.2) did not lead to statistically significant improvements in  
309 model performance, while adding random color filters (method 4.3) unequivocally worsened  
310 model performance. Neither SV (method 4.1) and PCA (method 4.2) introduces temporal or  
311 spectral distortions to the underlying image, thus satisfying the “spectrogram constraint.”  
312 However, specificity post-SV augmentation worsened significantly, likely due to the  
313 unconstrained shading changes to the spectrogram, which translates to drastic alterations of  
314 loudness/intensity at the audio level. The model is less able to identify “normal” heart sounds  
315 due to the sheer amount of unnatural variations in the training set that were labeled as normal  
316 based on the lack of murmurs. In contrast, incorporation of PCA data in the training set  
317 improved sensitivity and ROC AUC at the expense of a minor decrease in accuracy, and overall  
318 appears to be the second-best data augmentation method for cardiac analysis next to horizontal  
319 flip. At root, PCA establishes new features, known as “principal components,” from the original  
320 dataset. The goal is to compress the initial input dimensionality without compromising the most  
321 valuable information that were conveyed. Alterations along these “principal components”  
322 accomplish two objectives. First, they enrich the image along the axes of natural variation,  
323 which are by definition where the maximum between-sample variabilities exist. Second, since

324 changes are made at the color level, the underlying object invariance is maintained, which  
325 preserves the temporal and spectral properties of the original spectrograms. While intensity  
326 changes are unpredictable in SV because they are randomly generated, PCA's perturbations  
327 were derived mathematically, though still unconstrained by human physiological limits.  
328 Therefore, PCA suffers a similar pitfall as SV, though the detrimental effects are arguably much  
329 more blunted because the "physiologic constraint" is satisfied to a greater extent.

330

331 In contrast to the previous two techniques, random color filters entirely shift the hues outside the  
332 scope of our predetermined color-axis (i.e. orange). This may work for images of commonplace  
333 objects like cars, which can be observed in a wide variety of colors, but these augmentations  
334 are nonsensical for our heart sound spectrograms as they have no associated physical  
335 meaning. The spectrogram constraint is severely violated, and experimentally we observe that  
336 multicolor filters worsen model performance to the largest degree on all three counts. It is also  
337 important to note that in addition to the natural variations in heart sounds amplitudes, changes  
338 in amplitude may also reflect clinically relevant information. Pathological conditions such as  
339 cardiac tamponade classically lead to diminished heart sounds. Pleural effusions, subcutaneous  
340 edema, pneumothorax, and chronic obstructive pulmonary diseases (COPD) such as  
341 emphysema would also muffle heart sounds, although in these conditions the heart itself would  
342 be considered healthy. Similar to noise injection, alterations in heart sound amplitude could  
343 potentially blur the distinction between normal and abnormal heart sounds, which would worsen  
344 model performance. Epidemiologically, distant heart sounds from tamponade, pneumothorax, or  
345 COPD that is severe enough to muffle heart sounds are much rarer than murmurs. The majority  
346 of abnormal heart sounds in our data set are characterized by murmurs rather than distant heart  
347 sounds, explaining why amplitude perturbations did not have as much of a deleterious effect  
348 compared to noise injections.

349

350 The fifth augmentation method is time and frequency masking. Masking induces partial  
351 information loss at random points in the time and frequency domain. We surmise that masking  
352 has a similar effect to the regularization technique of dropout, where randomly selected neurons  
353 are ignored during training. However, in clinical practice, sudden quiescent periods occur in  
354 diseases such as AV heart block, cardiac arrest, or sick sinus syndrome. The original labels are  
355 preserved, so images that sprung from masking of normal spectrograms are still labeled as  
356 normal, despite the introduction of sudden pauses. Hence, masking does not satisfy the  
357 “physiologic constraint” and we observe model performance is not improved. Unlike noise  
358 injection and similar to amplitude changes, this type of pathological heart sound is relatively  
359 rare, thus there is no drastic reduction in performance. This stands in contrast to the state-of-the  
360 art results that masking has achieved in automated speech recognition<sup>33</sup>, further illustrating the  
361 distinction between clinical sound analysis and traditional audio processing.

362

### 363 **III. Conclusions**

364 Our experimental results corroborate our theoretical framework for thinking about heart sound  
365 spectrogram classification. Methods that violated the “spectrogram constraint”, such as vertical  
366 flipping and applying random color filters, worsened model performance by the greatest extent.  
367 Among the methods that did not violate the “spectrogram constraint”, the degree to which the  
368 “physiological constraint” was adhered to correlated with how much model performance  
369 improved or worsened. Noise injection is not a safe operation because the fundamental  
370 distinction between normal and abnormal heart sounds is blurred since the majority of abnormal  
371 heart sounds (murmurs, gallops, rubs) are just normal heart sounds with additional frequency  
372 components. Amplitude variation (via sensible color space transformations) and masking are  
373 also limited by fact that the distinction between normal and abnormal heart sounds are blurred:  
374 heart sounds with decreased amplitudes can be found in diseases such as cardiac tamponade,  
375 and heart sounds with quiescent periods can be found in disease such as AV block. However,

376 these augmentation methods are less fatal compared to noise injection because  
377 epidemiologically these heart sounds are much rarer, explaining why we did not observe a  
378 drastic reduction in model performance compared to noise injection. Pitch shifting and time  
379 stretching/compressing worsened model performance most likely because the alterations were  
380 outside physiological ranges. There is potential for this augmentation method to work but given  
381 that heart sounds naturally exist within a narrow physiologic range, future work includes  
382 precisely defining these boundaries. Interestingly, horizontal flipping is not actually rooted in any  
383 true physiological variation but has proven to be the superior data augmentation method.  
384 Horizontal flipping is able to create variation in the data without unnatural variations (such as at  
385 the extreme ends of pitch and time alterations) or run the risk of transforming normal sounds  
386 into abnormal sounds (such as with amplitude variations or masking). The “physiological  
387 constraint” and “spectrogram constraint” can be used as a guide for theory crafting future data  
388 augmentation methods for heart sound classification based on their spectral image. Moreover,  
389 the ideas behind the “physiological constraint” can be extended to related works seeking to  
390 classify heart sounds, while the ideas behind the “spectrogram constraint” can be extended to  
391 related work using spectrograms to classify audio.

392

393 In conclusion, there is value in data augmentation if done correctly, particularly for binary  
394 classification of PCG signals, and most likely for other medical classification problems as well.  
395 By synthetically generating samples using simple transformations, we can expand on the  
396 existing reservoir of patient data, and further enrich the documentation of select pathological  
397 conditions, which may be rare in nature and difficult to obtain. Machine learning models are  
398 increasingly used to streamline the repetitive processes in healthcare, such as initial screening,  
399 preliminary classifications, triage, patient sorting, and specialist recommendations. Data  
400 augmentation is a method that has shown utility in improving model performance in cardiac  
401 sound analysis and should be further explored in these alternative areas as well. In addition, this

402 study corroborates the idea that models are only as good as the data from which it learns.  
403 Disease-appropriate forms of data augmentation are integral to improvements in model  
404 performance, and synthetic data is most meaningful when it lies within the scope of human  
405 physiology and can accurately mimic clinical findings. Hence, physician input should be  
406 considered when creating models, so these tools can be useful and pragmatic both empirically  
407 and at the bedside.

408

## 409 **IV. Methods**

### 410 **4.1 Data**

411 The data in this study was sourced from a publicly available database assembled from the  
412 PhysioNet/Computing in Cardiology (CinC) Challenge in 2016<sup>22,23</sup>. The directory contains 3,239  
413 recorded heart sounds that range between 5-120 seconds. The sounds were compiled by  
414 physicians and research teams across seven countries over the course of a decade<sup>22,23</sup>. Experts  
415 in cardiology labelled the heart sounds as either normal or abnormal. Normal sounds are  
416 sounds collected from patients with no underlying cardiometabolic conditions. Abnormal sounds  
417 are sounds collected from patients with an underlying cardiac pathology, including valvular  
418 defects (i.e. mitral prolapse, mitral regurgitation, aortic regurgitation, aortic stenosis and valvular  
419 surgery), as well as coronary artery disease<sup>22,23</sup>. Of the recorded heart sounds, 2575 were  
420 labeled as normal and the remaining 664 sounds were labeled as abnormal.

421

### 422 **4.2 Pre-Processing**

423 In concordance with a previous study on heart murmur identification<sup>24</sup>, the raw heart sounds  
424 were first processed by a third-order Butterworth filter with a passband of 20-500 Hz, which  
425 encapsulates the range of normal heart sound and murmur frequencies<sup>25</sup>. All sounds under 8  
426 seconds were discarded. Then, the samples were either truncated to 30-seconds if their length

427 exceeded that limit, or preserved in their entirety if the length less than 30-seconds.

428 Subsequently, the amplitudes of the signals were normalized according to equation 1:

429 
$$X_{norm} = \frac{X}{\max(|X|)} \quad (1)$$

430 where  $X$  refers to the amplitude of the signal to ensure it is standardized across all recordings.

431

### 432 **4.3 Mel-Spectrogram**

433 The samples are windowed using a Hann window of size 512 and hop length of 256. A 512-

434 point Fast Fourier Transform is applied to each window to generate a spectrogram, which

435 depicts frequency over time. The amplitude of each frequency component is encoded in color.

436 The amplitude axis is converted to the dB scale, with the maximum amplitude serving as the

437 reference point and given a value of 0 dB. The frequency axis is transformed onto the Mel scale,

438 which is characterized by equation 2,

439 
$$Mel = 2595 * \log \left( 1 + \frac{f}{500} \right) \quad (2)$$

440 where  $f$  is frequency in Hz.

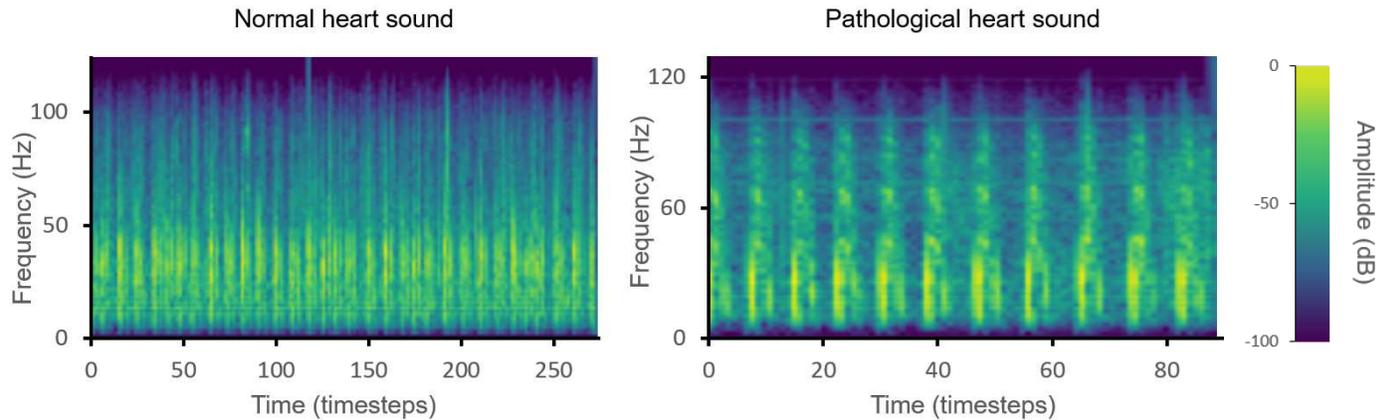
441

442 The resulting Mel-spectrogram images are standardized by rescaling each image to be of size

443 100x180 using bicubic interpolation. Figure 3 shows representative examples of the final Mel-

444 spectrogram images.

445



446 **Figure 3: Representative Mel-spectrograms of normal heart sound (left) and pathological**  
447 **heart sound (right)**

448

#### 449 **4.4 Data Augmentation**

##### 450 **4.4.1 Pitch Shifting and Time Stretching / Compression**

451 To create a synthetic heart sound under method 1, each real heart sound is first randomly pitch  
452 shifted up or down by  $p$  semitones, where  $p$  is a randomly chosen integer between 1 and 10. A  
453 semitone is defined as the interval between two adjacent notes in a 12-tone scale. For example,  
454 on a musical scale, the interval between  $C$  and  $C\#$  is one semitone. Then the pitch shifted  
455 sound is randomly time stretched/compressed by a factor of  $t$ , where  $t$  is randomly chosen from  
456 the uniform distribution  $[0.5, 2.0]$ . For example, if  $t=2.0$ , then a 30 second audio file is stretched  
457 to 60 seconds, or if  $t=0.5$ , then a 30 second audio file is compressed to 15 seconds. The pitched  
458 shifted and time stretched/compressed sounds are then converted to Mel-spectrogram images,  
459 which are used to supplement the Mel-spectrogram images derived from real heart sounds to  
460 train the convolutional neural network.

461

##### 462 **4.4.2 Noise Injection**

463 To create a synthetic heart sound under method 2, additive white Gaussian noises (AWGN) are  
464 injected element-wise into the original signal. The amplitude of AWGN is modeled as a  
465 Gaussian distribution, with  $\mu = 0$ .<sup>26</sup> The standard deviation of the noise signal is described with  
466 the following formula:

$$467 \quad RMS = \sqrt{\frac{\sum_i x_i^2}{n}}$$

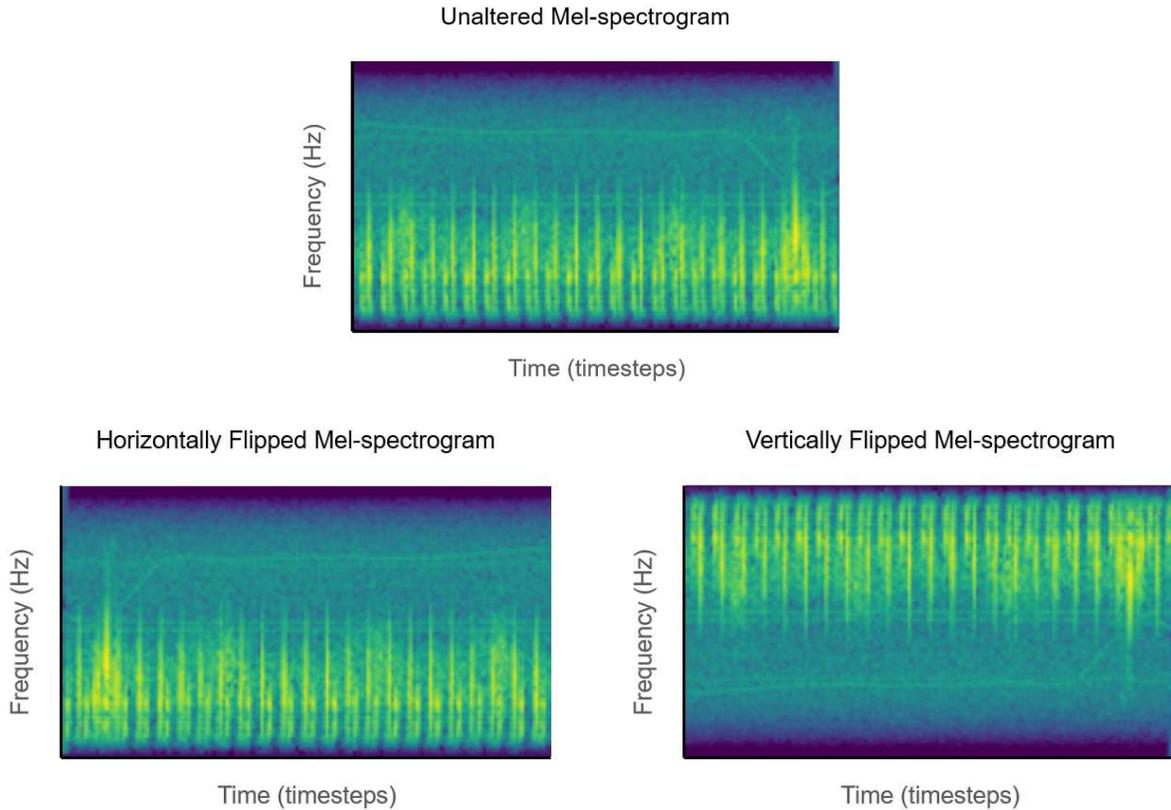
468 Assuming a signal-to-noise ratio (SNR) of 0, the required  $RMS_{noise}$  can be approximated by  
469  $RMS_{signal}$ . Each element of the noise signal is independently sampled from the distribution  
470  $X \sim N(\mu, \sigma^2)$  where  $\mu = 0$ ,  $\sigma = RMS_{signal}$ . The resulting noise signal is summed with the  
471 original sample. The synthetic samples are converted to Mel-spectrogram images and  
472 combined with the real heart sound Mel-spectrogram database to train the CNN model.

473

#### 474 **4.4.3 Image Flip**

475 To create synthetic data under method 3.1, each real heart sound is first converted to a Mel-  
476 spectrogram. The images are flipped horizontally, along an imaginary vertical axis that passes  
477 through its center, such that a given pixel with coordinate  $(x, y)$  will now be situated at  $(width -$   
478  $x - 1, y)$ . Figure 3 displays an example of the transformation. For method 3.2, the images are  
479 flipped vertically along a centered horizontal axis, such that a given pixel with coordinates  $(x, y)$   
480 will now be situated at  $(x, height - y - 1)$ . Figure 4 shows illustrative examples of a horizontally  
481 and vertically flipped spectrogram image.

482



484

485 **Figure 4: Unaltered Mel-spectrogram (top), horizontally flipped Mel-spectrogram (bottom**  
 486 **left), vertically flipped Mel-spectrogram (bottom right)**  
 487

#### 488 4.4.4 Color-Space Transformations

489 To create synthetic heart sound spectrograms under Method 4, the real heart sounds are first  
 490 converted into Mel-spectrograms. Then, each image was transformed into their RGB  
 491 representation, allowing for the extrapolation of other color-space values using pre-established  
 492 conversion factors and mathematical operations. For example, in an RBG-to-HSV  
 493 transformation, the red, green, and blue value which range from  $([0,255])$  for each pixel, is  
 494 converted into hue  $([0^\circ,360^\circ])$ , saturation  $([0-100\%])$ , and value/brightness  $([0-100\%])$  using the  
 495 following formulas<sup>27</sup>:

496

$$R' = \frac{R}{255}$$

$$497 \quad G' = \frac{G}{255}$$

$$498 \quad B' = \frac{B}{255}$$

$$499 \quad C_{max} = MAX(R', G', B')$$

$$500 \quad C_{min} = MIN(R', G', B')$$

$$501 \quad \Delta = C_{max} - C_{min}$$

$$502 \quad H = \begin{cases} 60^\circ \times \left( \frac{G' - B'}{\Delta} \bmod 6 \right), C_{max} = R' \\ 60^\circ \times \left( \frac{B' - R'}{\Delta} + 2 \right), C_{max} = G' \\ 60^\circ \times \left( \frac{R' - G'}{\Delta} + 4 \right), C_{max} = B' \end{cases}$$

$$503 \quad S = \begin{cases} 0, C_{max} = 0 \\ \frac{\Delta}{C_{max}}, C_{max} \neq 0 \end{cases}$$

$$504 \quad V = C_{max}$$

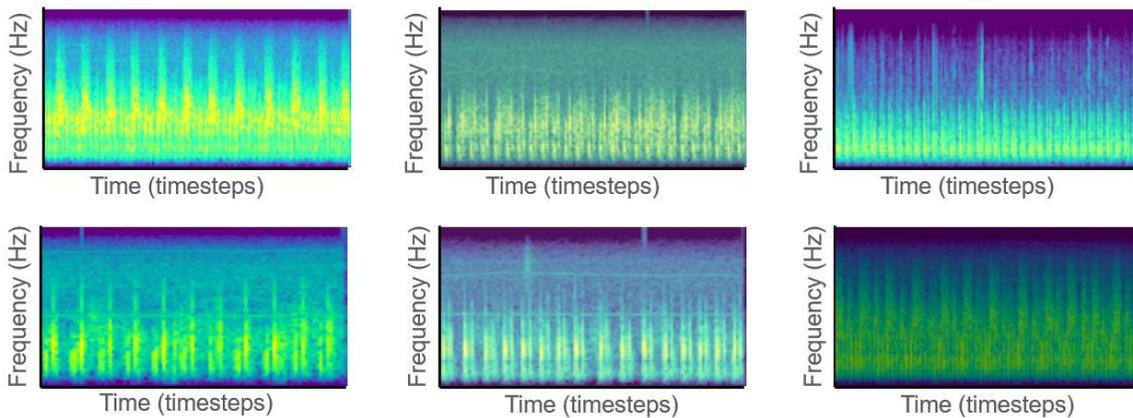
505 Within the scope of color space transformations, we explored three modalities of data  
 506 augmentation. Method 4.1 created new images from saturation and value perturbations. Method  
 507 4.2 created new images from Principal Component Analysis color augmentation, a method first  
 508 introduced in *Krizhevsky et al*<sup>28</sup>. Method 4.3 created new images from applying random color  
 509 filters.

510

#### 511 4.4.4.1 Method 4.1

512 In Method 4.1, two numbers,  $\alpha_{brightness}$  and  $\alpha_{saturation}$ , were randomly drawn from a uniform  
 513 distribution  $X \sim U(a, b)$ . Experimentally, it was determined that the  $\alpha_{brightness}$  would be bounded  
 514 by  $a=0.5$  and  $b=2$ , and  $\alpha_{saturation}$  by  $a=0.1$  and  $b=2$ .  $\alpha_{brightness}$  and  $\alpha_{saturation}$  control the  
 515 degree of brightness and saturation perturbations, respectively. The merging operation can be  
 516 described with the following formula:

517 Blending Image \* (1 -  $\alpha$ ) + Original Image \*  $\alpha$   
518 Brightness alterations were achieved by blending the original image with a pure black image of  
519 the same dimensions. Saturation alterations were achieved by blending the original image with  
520 a grey-scale image of the same dimensions. The two perturbations were applied sequentially to  
521 the original image, and the adjustment factors  $\alpha_{brightness}$  and  $\alpha_{saturation}$  were redrawn for each  
522 input spectrogram. Figure 5 shows spectrograms that have undergone saturation and  
523 brightness perturbations.



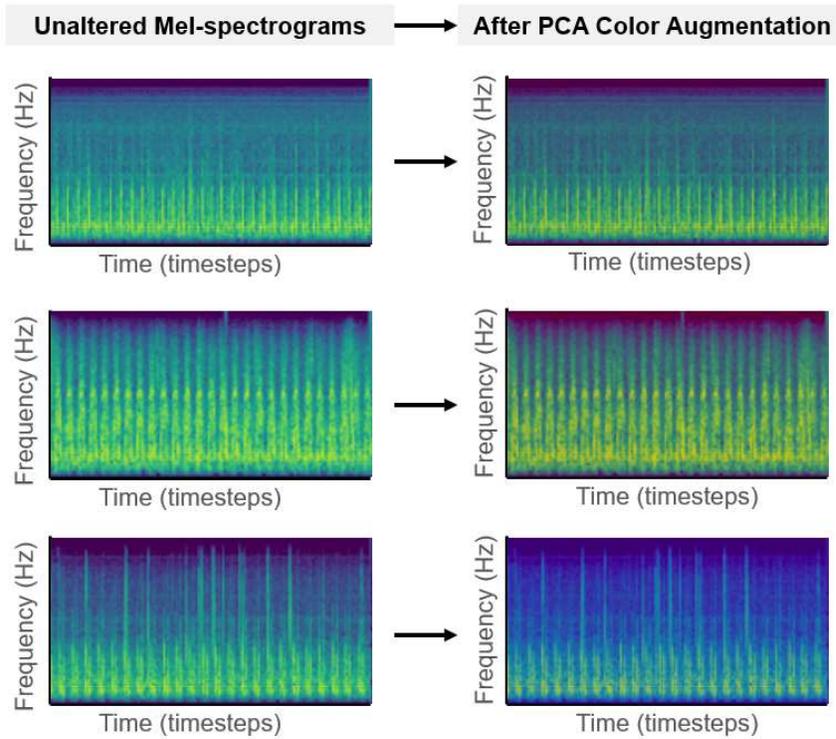
524 **Figure 5: Representative Mel-spectrograms with Saturation Brightness Perturbations**

525

#### 526 4.4.4.2 Method 4.2

527 In Method 4.2, as described in *Krizhevsky et al*<sup>28</sup>, we implemented principal component analysis  
528 on the unaltered input images, yielding a sorted set of eigenvectors and eigenvalues that are  
529 associated with the 3x3 covariance matrix of the RGB color channels. We then drew a random  
530 variable  $\alpha$  from the normal distribution  $X \sim N(\mu, \sigma^2)$ , where  $\mu = 800$ ,  $\sigma = 10$ , and multiplied it to  
531 the original eigenvalues. The principal components are scaled by the output from the previous  
532 step, and the product is added to the RGB vector of each individual pixel.  $\alpha$  is drawn once for  
533 each training image. The specific mean and standard deviation values of the perturbation were

534 chosen experimentally, to intentionally produce more pronounced differences in the output  
535 images. Figure 6 shows spectrograms that have undergone PCA color augmentation.



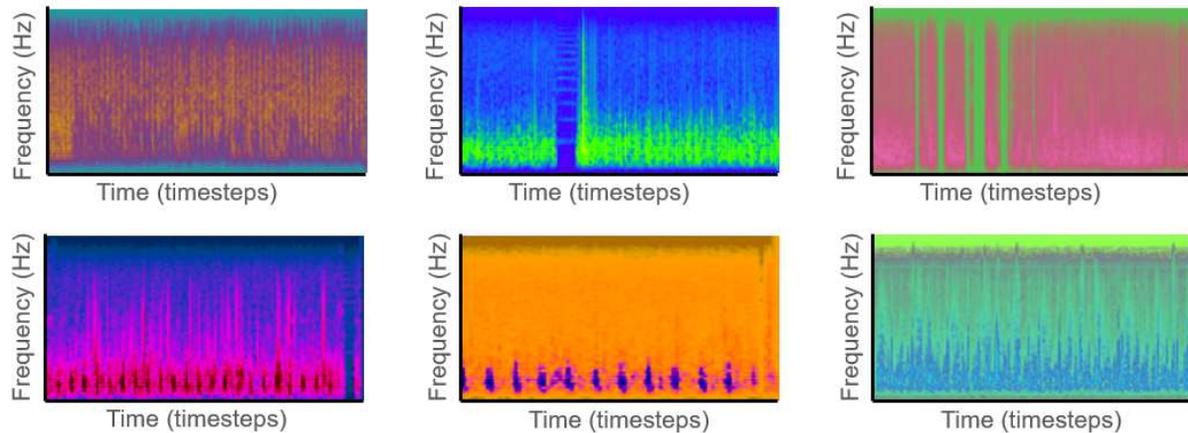
536

537 **Figure 6: Unaltered Mel-spectrograms (Left), same images after Principal Component**  
538 **Analysis (PCA) Color Augmentation (Right) (Data Augmentation Method 4.2)**  
539

#### 540 4.4.4.2 Method 4.3

541 In Method 4.3, we iterated through a library of 150 different color-space conversions using the  
542 OpenCV package, effectively generating random color balance perturbations, but preserving the  
543 underlying shapes and content of the input images. The transformed Mel-spectrograms are  
544 used to supplement the Mel-spectrograms from real heart sounds as additional training data.  
545 Figure 7 shows spectrograms with random color filters applied.

546



547

548 **Figure 7: Representative Mel-spectrograms with Random Color Filters**

549

550 **4.4.5 Time and Frequency Masks**

551 To create synthetic heart sound data under Method 5, the real heart sounds are left untouched

552 and converted to Mel-spectrogram images. To the Mel-spectrogram image, three masks are

553 randomly applied in the time domain, and three masks are randomly applied in the frequency

554 domain. In frequency masking, the frequency channels  $[f_0, f_0 + f)$  are masked, where  $f$  is

555 randomly chosen from the uniform distribution  $[0, 20]$ , and  $f_0$  is randomly chosen from  $(0, v - f)$ ,

556 where  $v$  is the total number of frequency channels. In time masking, the time steps  $[t_0, t_0 + t)$  are

557 masked, where  $t$  is randomly chosen from the uniform distribution  $[0, 20]$ , and  $t_0$  is randomly

558 chosen from  $[0, \tau - t]$ , where  $\tau$  the total number of time steps. Figure 3 illustrates an example of

559 a transformed Mel-spectrogram. The location of the masks is chosen independently, meaning it

560 is possible for masks to overlap and merge into one larger mask. The transformed Mel-

561 spectrogram images are used to supplement the Mel-spectrogram images derived from real

562 heart sounds to train the convolutional neural network. Figure 8 shows a spectrogram with time

563 and frequency masking applied.

564

565

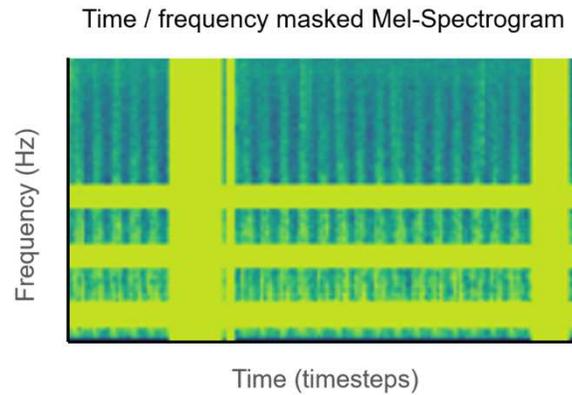
566

567

568

569

570



571 **Figure 8: Representative example of time / frequency masked Mel-spectrogram**

572 Three masks, as represented by the yellow bars, are randomly applied in the time domain, and  
573 three masks are randomly applied in the frequency domain.

574

#### 575 **4.5 Convolutional Neural Network**

576 The resulting Mel-spectrograms are treated as images and used to train a convolutional neural

577 network (CNN) for binary classification. A prior study that explored heart sound classification

578 provided an optimized CNN framework that inspired the basis of the CNN architecture used in

579 this study<sup>29</sup>. The convolutional neural network model we built consists of four layers. The first

580 layer is a convolution layer with 32 3x3 kernels, each with a stride length of one; the activation

581 function used is a rectified linear (ReLU) activation function.

582

583 This is followed by a max pooling layer with a filter of size 2x2 with a stride length of two. The

584 second layer is a convolutional layer with 64 3x3 kernels, each with a stride length of one; the

585 activation function used is a ReLU activation function. Similarly, it is followed by a max pooling

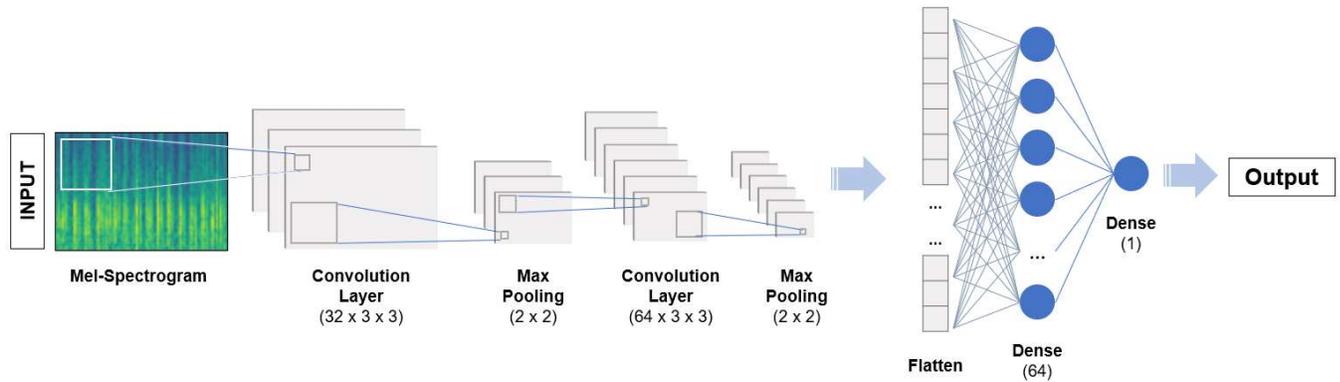
586 layer with a filter of size 2x2 with a stride length of two. Padding is not used in any layer. The

587 output from the previous operation is flattened into a one-dimensional feature vector, and then

588 passed to the third layer, a fully connected layer with 64 hidden units. The fourth and final layer

589 is a single neuron with a sigmoid activation function to make the final binary classification. We

590 used the Adaptive Moment Estimation (Adam) optimizer to iteratively improve model  
591 performance. Ten epochs are used for training. Figure 9 shows the CNN architecture.



592

### 593 **Figure 9: Convolutional Neural Network Structure**

594 Illustration of the CNN architecture employed in our study for heart sound classification.  
595

## 596 **V. List of Abbreviations**

597 **CNN:** Convolutional Neural Networks

598 **AUC:** Area under the receiver operating curve

599 **ROC:** Receiver Operating Curve

600 **PCA:** Principal component analysis

601 **SV:** Saturation/Value

602 **PCG:** phonocardiogram

603 **HIPAA:** Health Insurance Portability and Accountability Act

604 **IRB:** Institutional Review Board

605 **EHR:** Electronic health records

606 **COVID-19:** Coronavirus disease of 2019

607 **AWGN:** Additive white Gaussian noises

608 **SNR:** Signal-to-noise ratio

609 **RMS:** Root mean squared

610 **RGB:** Red/Green/Blue

611 **HSV:** Hue/Saturation/Value  
612 **ReLU:** Rectified linear activation function  
613 **Adam:** Adaptive Moment Estimation  
614 **CIFAR:** Canadian Institute For Advanced Research  
615 **COPD:** Chronic obstructive pulmonary diseases

616

## 617 **VI. References**

- 618 1 Nielsen T, Molgaard H, Ringsted C, Eika B. The development of a new cardiac auscultation  
619 test: how do screening and diagnostic skills differ? *Med Teach* 2009;online first, DOI:  
620 10.1080/01421590802572767.  
621
- 622 2 Mangione S. Cardiac auscultatory skills of physicians-in-training: a comparison of three  
623 English-speaking countries. *Am J Med* 2001;**110**:210–16.  
624
- 625 3 Dhuper S, Vashist S, Shah N, Sokal M. Improvement of cardiac auscultation skills in pediatric  
626 residents with training. *Clin Pediatr (Phila)* 2007;**46**:236–40.  
627
- 628 4 Nogueira D. M., Ferreira C. A., Gomes E. F., Jorge A. M. Classifying heart sounds using  
629 images of motifs, MFCC and temporal features. *Journal of Medical Systems*. 2019;43(6):p. 168.  
630 doi: 10.1007/s10916-019-1286-5  
631
- 632 5 Hamidi M., Ghassemian H., Imani M. Classification of heart sound signal using curve fitting  
633 and fractal dimension. *Biomedical Signal Processing and Control*. 2018;39:351–359. doi:  
634 10.1016/j.bspc.2017.08.002.  
635
- 636 6 V. Maknickas, A. Maknickas Recognition of normal–abnormal phonocardiographic signals  
637 using deep convolutional neural networks and mel-frequency spectral coefficients  
638 *Physiol. Meas.*, 38 (8) (2017), pp. 1671-1684  
639
- 640 7 Juniati D., Khotimah C., Wardani D. E. K., Budayasa K. Fractal dimension to classify the heart  
641 sound recordings with KNN and fuzzy c-mean clustering methods. *Journal of Physics  
642 Conference Series*. 2018;953:p. 012202. doi: 10.1088/1742-6596/953/1/012202.  
643
- 644 8 Karar M. E., El-Khafif S. H., El-Brawany M. A. Automated diagnosis of heart sounds using  
645 rule-based classification tree. *Journal of Medical Systems*. 2017;41(4):p. 60. doi:  
646 10.1007/s10916-017-0704-9  
647
- 648 9 Deng S. W., Han J. Q. Towards heart sound classification without segmentation via  
649 autocorrelation feature and diffusion maps. *Future Generation Computer Systems*. 2016;60:13–  
650 21. doi: 10.1016/j.future.2016.01.010  
651
- 652 10 Zhang W., Han J., Deng S. Heart sound classification based on scaled spectrogram and  
653 tensor decomposition. *Expert Systems with Applications*. 2017;84:220–231. doi:  
654 10.1016/j.eswa.2017.05.014

655  
656 11 Whitaker B. M., Suresha P. B., Liu C., Clifford G. D., Anderson D. V. Combining sparse  
657 coding and time-domain features for heart sound classification. *Physiological Measurement*.  
658 2017;38(8):1701–1713. doi: 10.1088/1361-6579/aa7623.  
659  
660 12 Cheng X., Zhan Q., Wang J., Ma R. A high recognition rate of feature extraction algorithm  
661 without segmentation. *IEEE 6th International Conference on Industrial Engineering and*  
662 *Applications (ICIEA)*; 2019; Tokyo, Japan. IEEE; pp. 923–927.  
663  
664 13 Wang P, Lim CS, Chauhan S, Foo JY, Anantharaman V. Phonocardiographic signal analysis  
665 method using a modified hidden Markov model. *Ann Biomed Eng*. 2007 Mar;35(3):367-74. doi:  
666 10.1007/s10439-006-9232-3. Epub 2006 Dec 14. PMID: 17171300.  
667  
668 14 Chauhan S, Wang P, Sing Lim C, Anantharaman V. A computer-aided MFCC-based HMM  
669 system for automatic auscultation. *Comput Biol Med*. 2008 Feb;38(2):221-33. doi:  
670 10.1016/j.combiomed.2007.10.006. Epub 2007 Nov 28. PMID: 18045582.  
671  
672 15 T. I. Yang and H. Hsieh, "Classification of acoustic physiological signals based on deep  
673 learning neural networks with augmented features," *2016 Computing in Cardiology Conference*  
674 *(CinC)*, Vancouver, BC, 2016, pp. 569-572.  
675  
676 16 Raza, A., Mehmood, A., Ullah, S., Ahmad, M., Choi, G. S., & On, B. W. (2019). Heartbeat  
677 Sound Signal Classification Using Deep Learning. *Sensors* (Basel, Switzerland), 19(21), 4819.  
678 <https://doi.org/10.3390/s19214819>  
679  
680 17 H. Ryu, J. Park and H. Shin, "Classification of heart sound recordings using convolution  
681 neural network," *2016 Computing in Cardiology Conference (CinC)*, Vancouver, BC, 2016, pp.  
682 1153-1156.  
683  
684 18 B. Bozkurt, I. Germanakis, Y. Stylianou A study of time-frequency features for CNN-based  
685 automatic heart sound classification for pathology detection *Comput. Biol. Med.*, 100 (August  
686 2017) (2018), pp. 132-143  
687  
688 19 Halevy A, Norvig P, Pereira F. The unreasonable effectiveness of data. *IEEE Intell Syst*.  
689 2009;24:8–12.  
690  
691 20 Chen S, Abhinav S, Saurabh S, Abhinav G. Revisiting unreasonable effectiveness of data in  
692 deep learning era. In: *ICCV*; 2017. p. 843–52.  
693  
694 21 Shorten, C., Khoshgoftaar, T.M. A survey on Image Data Augmentation for Deep Learning. *J*  
695 *Big Data* 6, 60 (2019). <https://doi.org/10.1186/s40537-019-0197-0>  
696  
697 22 Liu C, Springer D, Li Q, Moody B, Juan RA, Chorro FJ, Castells F, Roig JM, Silva I, Johnson  
698 AE, Sye Z, Schmidt SE, Papadaniil CD, Hadjileontiadis L, Naseri H, Moukadem A, Dieterlen A,  
699 Brandt C, Tang H, Samieinasab M, Samieinasab MR, Sameni R, Mark RG, Clifford GD. An  
700 open access database for the evaluation of heart sound algorithms. *Physiological Measurement*  
701 2016;37(9)  
702  
703 23 Goldberger A, Amaral L, Glass L, Hausdorff J, Ivanov PC, Mark R, Mietus JE, Moody GB,  
704 Peng CK, Stanley HE. PhysioBank, PhysioToolkit, and PhysioNet: Components of a new

- 705 research resource for complex physiologic signals. *Circulation* [Online]. 101 (23), pp. e215–  
706 e220.
- 707
- 708 24 S. Kang, R. Doroshov, J. McConnaughey and R. Shekhar, "Automated Identification of  
709 Innocent Still's Murmur in Children," in *IEEE Transactions on Biomedical Engineering*, vol. 64,  
710 no. 6, pp. 1326-1334, June 2017, doi: 10.1109/TBME.2016.2603787.
- 711
- 712 25 McGee, S. (2018). Chapter 39 - Auscultation of the Heart: General Principles. Evidence-  
713 Based Physical Diagnosis (Fourth Edition). S. McGee. Philadelphia, *Elsevier*: 327-332.e321.
- 714
- 715 26 Liu W, Lin W. Additive white Gaussian noise level estimation in SVD domain for images.  
716 *IEEE Trans Image Process* [Internet]. 2013 Mar;22(3):872–83. Available from:  
717 <http://www.ncbi.nlm.nih.gov/pubmed/23008255>
- 718
- 719 27 G. Saravanan, G. Yamuna and S. Nandhini, "Real time implementation of RGB to  
720 HSV/HSI/HSL and its reverse color space models," *2016 International Conference on*  
721 *Communication and Signal Processing (ICCSP)*, Melmaruvathur, 2016, pp. 0462-0466, doi:  
722 10.1109/ICCSP.2016.7754179.
- 723
- 724 28 Krizhevsky A, Sutskever I, Hinton GE. ImageNet classification with deep convolutional neural  
725 networks. *Commun ACM* [Internet]. 2017 May 24;60(6):84–90. Available from:  
726 <https://dl.acm.org/doi/10.1145/3065386>
- 727
- 728 29 Zhang W, Han J. Towards heart sound classification without segmentation using  
729 convolutional neural network. In: *Computing in Cardiology*. 2017.
- 730
- 731 30 M. Tschannen, T. Kramer, G. Marti, M. Heinzmann and T. Wiatowski, "Heart sound  
732 classification using deep structured features," *2016 Computing in Cardiology Conference*  
733 *(CinC)*, Vancouver, BC, 2016, pp. 565-568.
- 734
- 735 31 Bradley M Whitaker *et al* 2017 *Physiol. Meas.* 38 1701
- 736
- 737 32 J. Shijie, W. Ping, J. Peiyi and H. Siping, "Research on data augmentation for image  
738 classification based on convolution neural networks," *2017 Chinese Automation Congress*  
739 *(CAC)*, Jinan, 2017, pp. 4165-4170, doi: 10.1109/CAC.2017.8243510.
- 740
- 741 33 Park, D.S., Chan, W., Zhang, Y., Chiu, C., Zoph, B., Cubuk, E.D., Le, Q.V. (2019)  
742 SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition. *Proc.*  
743 *Interspeech 2019*, 2613-2617, DOI: 10.21437/Interspeech.2019-2680.
- 744

745

746 **Declarations**

747

748 **Ethical approval and consent to participate**

749 Not applicable

750

751 **Consent for publication**

752 Not Applicable

753

754 **Availability of data and materials**

755 The datasets used and/or analyzed during the current study are available at  
756 <https://physionet.org/content/challenge-2016/1.0.0/>

757

758 **Competing interests**

759 The authors declare no competing interests.

760

761 **Funding**

762 The authors did not receive funding for this study.

763

764 **Author's Contributions**

765 G.Z. conceived of the presented idea. G.Z. and Y.C. designed and performed the experiments.

766 G.Z. and Y.C. interpreted the results. G.Z., Y.C. and C.C. wrote the manuscript.

767

768 **Acknowledgements**

769 The authors would like to thank Dr. Patrick Flynn and Dr. George Shih for their support.

770

771 **Author's Information**

772 Affiliations

773 Weill Cornell Medicine, New York, NY, 10021, USA

774 George Zhou, Yunchan Chen, Candace Chien

775

776

777

778

779