# Robust and Task-Aware Training of Deep Residual Networks for Varying-Lead ECG Classification

Hansheng Ren[*1], Miao Xiong[*1], Bryan Hooi[1]
[1]National University of Singapore

## Abstract

*In PhysioNet/Computing in Cardiology Challenge 2021, we developed an ensemble model by combining different epochs of ResNet to classify cardiac abnormalities from 12, 6, 4, 3, 2 lead electrocardiogram (ECG) signals, where epochs are chosen based on validation performance on China Physiological Signal Challenge (CPSC) dataset and Georgia dataset. In order to adapt to the specially designed Challenge score, we designed a multi-task loss to combine the benefit of binary cross-entropy loss and Challenge loss. Besides, we also integrated a subsample frequency feature into the model to learn from the signals. To gain a better generalization ability, mixup and weighted loss are introduced.*

*We submitted our model in the official phase with team name **DataLA_NUS**, and our final selected model achieved a Challenge score of 0.51, 0.51, 0.51, 0.50, 0.52 (ranked 8th, 5th, 6th, 8th, 5th) on the 12-lead, 6-lead, 4-lead, 3-lead, and 2-lead setting on the final hidden test set with the Challenge evaluation metric.*

## 1. Introduction

Electrocardiography is a commonly used, non-invasive technique for recording electrical changes and examining the physiological activities of the heart. The record, named electrocardiogram (ECG), shows the series of waves that relate to the electrical impulses that occur during each heartbeat. In the PhysioNet/Computing in Cardiology Challenge 2021 [1, 2], we are required to develop models for electrocardiograms of 12, 6, 4, 3, 2 leads to identify the cardiac abnormalities present in recordings from [1–8] datasets. In this paper, we will show our proposed methods and what we have tried in this Challenge for inspiration.[1]

## 2. Methods

In this section, we first introduce how we process the data. After obtaining the preprocessed signals, we train our model using multi-task learning to align better with

---

[1]* These authors contributed equally and are co-first authors.

---

clinical realities: We adopt a normal binary cross-entropy loss and a loss based on the Challenge score metric that is later referred as Challenge loss. For every input signal, we first feed them into our modified 1D-resnet18 [9] or 1D-effientnet [10] model to get its embedding and then pass this embedding through two independent linear layers to generate two predictions. These two predictions are used to calculate binary cross-entropy loss and Challenge loss separately.

### 2.1. Data Processing

The Challenge data was sourced from eight organizations [1–8], and recordings from different sources have distinct features. As an initial step, we want to process all samples into fixed length and fixed sample frequency. Based on the observation that most recordings are 500Hz, we decide to upsample or downsample all recordings into 500Hz. By analyzing the dataset statistics, we find that most samples are around ten seconds, and the INCART [5] dataset with 30-minute-long signals only contains 74 recordings and does not appear in the test dataset. We decide to sample all recordings to samples with lengths equal to 4096 instead of splitting them into multiple patches. Zero padding is introduced for recordings shorter than 4096, and for recordings longer than 4096, we randomly select the start point to fetch a fixed-length patch. Besides the recording itself, we also utilize some auxiliary information, including age and sex.

### 2.2. Model Architecture

For Our model architecture, we tried two types of widely used deep learning architecture for generating our recording embeddings, 1D-Resnet18 and 1D-Efficientnet, with details as below.

**1D-Resnet18** Except the design of multi-task loss, We also modify the Resnet18 [9] in a way following the last year challenge's winning solution[11]:one convolutional layer followed by $N = 8$ residual blocks (ResBs) with large kernel and dropout layer, each of which contains two convolutional layers and a squeeze and excitation (SE) block introduced by [12]. After obtaininng the embedding of sig-

nals, a fully connected layer is introduced to analyze auxiliary information, including patient age and gender, and join in recording embedding into two independent dense layers for different loss functions.

**1D-Efficient net** We also try to use the modified 1D EfficientNet [10] as the main model architecture during the official phase. This model is first published on [10] and its one-dimensional PyTorch implementation originates from this repo.

**Challenge loss** Since our classification task considers the clinical reality that making some misdiagnoses is more harmful than others and some misdiagnoses are partially acceptable, we design a special challenge loss to adapt our model. Besides, we design a special task-aware loss function to teach our model different severity of making misdiagnoses during training. Specially speaking, we define continuous relaxation of the scoring metric and optimizing the resulting loss function directly in our end-to-end framework so the model can align better with clinical realities as captured in the scoring metric. Specifically, denote the weight of score metric as $w_{ij}$, the model prediction as x, and the label as y. The challenge loss is calculated by:

$$norm = \max(1, \text{sum}(\max(x, y))) \quad (1)$$

$$z_{ij} = y_i x_j w_{ij} \quad (2)$$

$$\text{Challenge Loss} = \frac{\text{sum}(z)}{norm} \quad (3)$$

This challenge loss is indeed the challenge scoring metric generalized to non-binary input cases. So after feeding the recordings to the model during training, we can get a generalized challenge score between model scalar outputs and labels without any threshold and then use it as a loss to supervise backward. We also implemented a batch-style computing function of this challenge loss to make its speed acceptable for training our model.

**BCEloss** We also adopt the binary cross entropy loss since this Challenge task is basically a multi-label classification problem.

**Weight Loss for Imbalanced Class** We also tried a modified BCE loss for the data imbalance problem. Figure 1 shows the positive labels of the ground truth for every class. It demonstrates that the data imbalance problem is significant in most classes. In order to solve this problem, we have tried two solutions: one is to use the weighted binary cross-entropy loss; another is to use an imbalanced data sampler. The former approach slightly increases the final performance while the latter does not work (we leave it in the discussion part). We define the weight manually, different from the standard way of using a positive/negative ratio to compute the weights. Similar to Figure 1, we calculate the validation set accuracy for every class and try to increase the weight for classes in which our proposed model performs poorly.
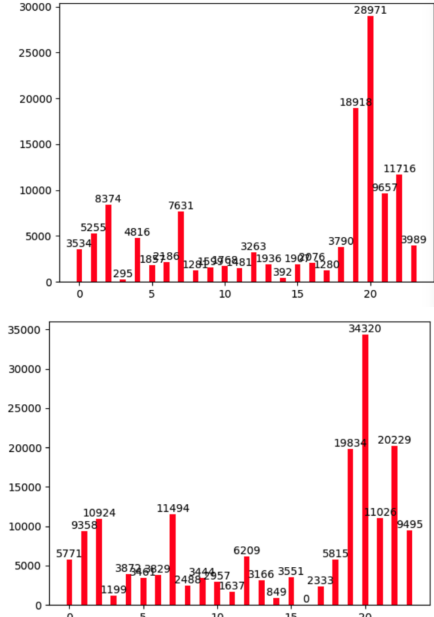


Figure 1. Comparison of the number of every class's positive samples in ground truth labels(**Upper**) and prediction results(**Bottom**) on the training data set. The result is given by the main model (resnet18 with local threshold optimization and multi-loss).

**Mixup** In order to enhance the generalization ability, we add mixup [13] to the training process. Firstly, we sample two datapoints $(x_i, y_i)$ and $(x_j, y_j)$ from the training dataset. Then we contruct a new mixed datapoint $x_{ij} = \alpha x_i + (1 - \alpha)x_j$ with new label $y_{ij} = \alpha y_i + (1 - \alpha)y_j$. Model trained on these mixed datapoints tends to learn soft decision boundary between classes.

## 2.3. Implementation Setup

In this section, we introduce three training techniques we used for our model: loading parameters of 12-lead trained model to model with fewer leads, threshold optimization and model ensemble.

**loading parameters of 12-lead trained model to model with fewer leads** We observe that if we adopt modified ResNet model architecture when we train fewer leads model, we can reuse almost all the parameters from 12 leads except the first convolution layer since the following layers' parameter layout is the same. So we first train a 12 lead model. And when we are going to train a model with fewer than 12 leads, we initialize the model's parameters by the trained 12 lead model's parameter, except the first convolutional layer.

**Threshold optimization** For multi-label classification task, we perform binary prediction for every class based on

scalar output generated by challenge loss and the learned threshold. Since we did 80/20 data split after data pre-processing, we can use the validation data to conduct a grid search on the threshold. We search a global threshold (the same threshold for all classes) first and then do a local threshold search for each class. It is worth mentioning that if we directly use the divided 20% validation set to search the model threshold and make the model selection, it seems not optimal. So we do an 80/20 data split for the whole training data first, and then in the validation set we duplicate recordings from the Georgia dataset by 21 times and recordings from the CPSC dataset by seven times to make a validation set with a data distribution similar to the hidden validation set for threshold search and model selection.

**Model ensemble** During the official phase, we submit several ensemble versions of training models to the challenge server. We have two types of ensemble: the first one is the majority vote of models, that is, each model produces a binary prediction based on their own saved threshold searched by our self-divided validation set for a test recording; another is that we calculate the mean of these model's corresponding searched thresholds and apply it to the model's scalar outputs.

## 3. Results

The 2 lead Challenge scores of our submitted models on validation set are shown in table 1. We denote the default version as ResNet18 trained using BCEloss plus 0.1 Challenge loss, without model ensemble and mixup technique. Then, sub version 1 denotes ResNet18 with only global threshold search; sub 2 is ResNet18 with local threshold search; sub3 is ResNet18 trained by mixup technique with 0.3 weight and global threshold search. Sub 4 is ResNet18 trained by weighted loss mention above and with local threshold search. Sub 5 is 1D EffientNet with local threshold search. Sub 6 is ResNet18 trained by mixup technique with 0.15 weight and local threshold search. Sub 7 is the ResNet18 with no validation set, trained on all given data and use hand-set threshold 0.1, ensemble models from 20 epoch to 30 epoch with majority vote strategy. Sub 8 is ResNet18 with local threshold search on validation set while ensemble models from 20 epoch to 30 epoch with mean threshold ensemble strategy mentioned above.

Table 2 shows the results across all leads where version 1 denotes Resnet18 trained by 0.1 Challenge loss and BCE loss with trainable parameters from 12-lead applying to to fewer leads of model as pretrain, and version 2 indicates Resnet18 trained by 0.1 Challenge loss and BCE loss and using mixup training technique with weight 0.15. version 3 means Resnet18 trained by 0.1 Challenge loss and BCE loss, ensemble models from 20 epoch to 30 epochs using mean thresholds.

| Methods | 2-Lead Score |
|---|---|
| sub1(global threshold search) | 0.612 |
| sub2(local threshold search) | 0.621 |
| sub3(0.3 mixup training) | 0.598 |
| sub4(weighted loss) | 0.602 |
| sub5(1D EfficientNet) | 0.578 |
| sub6(0.15 mixup training) | 0.623 |
| sub7(predefined threshold 0.1 ensemble) | 0.162 |
| sub8(mean threshold ensemble) | 0.625 |

Table 1. 2 lead Challenge scores on the validation set for different methods.

| Methods | 12 leads | 6 leads | 4 leads | 3 leads | 2 leads |
|---|---|---|---|---|---|
| version 1 | 0.631 | 0.603 | 0.597 | 0.574 | 0.560 |
| version 2 | 0.620 | 0.615 | 0.615 | 0.623 | 0.620 |
| version 3 | 0.637 | 0.622 | 0.623 | 0.621 | 0.625 |

Table 2. Challenge scores of submitted models on 12/6/4/3/2 leads on validation set.

The final test results are listed below:

| Leads | Validation | Test | Ranking |
|---|---|---|---|
| 12 | 0.64 | 0.51 | 8 |
| 6 | 0.62 | 0.51 | 5 |
| 4 | 0.62 | 0.51 | 6 |
| 3 | 0.62 | 0.50 | 8 |
| 2 | 0.63 | 0.52 | 5 |

Table 3. Challenge scores for our final selected entry (team DataLA_NUS) using **version 3** on the hidden validation set and test set as well as the ranking on the hidden test set.

## 4. Discussions

During the exploration of finding a suitable classifier for this Challenge, we have tried many different approaches. Some do increase the performance, while some are not. Even though some methods do not work, we believe it might shed light on improving the performance.

**Rule-based Method** We did this exploration by starting from the 'Bradycardia' disease, whose definition is simple enough for a non-expert to identify. Figure 1 shows our default model's performance on every class, indicating that our model fails to perform well at identifying Bradycardia. Based on Wiki, Bradycardia is a condition typically defined wherein an individual has a resting heart rate of under 60 beats per minute (BPM) in adults, although some studies use a heart rate of less than 50 BPM [14]. Inspired by this, we compute every recording's BPM and assign positive labels to samples with BPM lower than 60. Unfortunately, we find that many recordings with BPM less
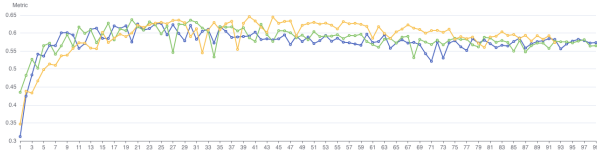
Figure 2. Validation accuracy during every training epoch. The blue plot represents the model trained on resnet18, the green plot denotes the model trained no resnet34, while the yellow one represents the model trained on resnet152. It can be seen that after 20-30 epochs the validation set performance has been decreased, **indicating the occurance of overfitting**.

than 60 and 50 are not labeled as Bradycardia, making this algorithm fail to improve performance.

## 5.   Conclusions

During this Challenge, we tried two types of widely used deep learning architectures for generating feature embeddings of every recordings: 1D-Resnet18 and 1D-Efficientnet. Besides, we utilized the multi-task weighted loss to adapt to the specially designed Challenge score metric, and used mixup training technique to increase the generalization ability. To better classify, we also integrated model ensemble, threshold optimization and pre-train techniques. All these tricks contribute to our final performance and make us ranked as the 7th in the official test ranking in total, while achieved 5th on the 6/2 leads.

## References

[1]  Alday EAP, Gu A, Shah AJ, Robichaux C, Wong AKI, Liu C, et al. Classification of 12-lead ECGs: the PhysioNet/Computing in Cardiology Challenge 2020. Physiological Measurement January 2021;41(12).

[2]  Reyna MA, Sadr N, Perez Alday EA, Gu A, Shah A, Robichaux C, et al. Will Two Do? Varying Dimensions in Electrocardiography: the PhysioNet/Computing in Cardiology Challenge 2021. Computing in Cardiology 2021;48:1–4.

[3]  Bousseljot R, Kreiseler D, Schnabel A. Nutzung der EKG-Signaldatenbank CARDIODAT der PTB über das Internet. Biomedizinische Technik 1995;40(S1):317–318.

[4]  Liu F, Liu C, Zhao L, Zhang X, Wu X, Xu X, et al. An Open Access Database for Evaluating the Algorithms of Electro-cardiogram Rhythm and Morphology Abnormality Detection. Journal of Medical Imaging and Health Informatics 2018;8(7):1368—1373.

[5]  Tihonenko V, Khaustov A, Ivanov S, Rivin A, Yakushenko E. St Petersburg INCART 12-lead Arrhythmia Database. PhysioBank PhysioToolkit and PhysioNet 2008;Doi: 10.13026/C2V88N.

[6]  Wagner P, Strodthoff N, Bousseljot RD, Kreiseler D, Lunze FI, Samek W, et al. PTB-XL, a Large Publicly Available Electrocardiography Dataset. Scientific Data 2020;7(1):1–15.

[7]  Zheng J, Zhang J, Danioko S, Yao H, Guo H, Rakovski C. A 12-lead Electrocardiogram Database for Arrhythmia Research Covering More Than 10,000 Patients. Scientific Data 2020;7(48):1–8.

[8]  Zheng J, Cui H, Struppa D, Zhang J, Yacoub SM, El-Askary H, et al. Optimal Multi-Stage Arrhythmia Classification Approach. Scientific Data 2020;10(2898):1–17.

[9]  He K, Zhang X, Ren S, Sun J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In Proceedings of the IEEE international conference on computer vision. 2015; 1026–1034.

[10]  Tan M, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In International Conference on Machine Learning. PMLR, 2019; 6105–6114.

[11]  Zhao Z, Fang H, Relton SD, Yan R, Liu Y, Li Z, et al. Adaptive lead weighted resnet trained with different duration signals for classifying 12-lead ecgs. In 2020 Computing in Cardiology. IEEE, 2020; 1–4.

[12]  Hu J, Shen L, Sun G. Squeeze-and-excitation networks. In Proceedings of the IEEE conference on computer vision and pattern recognition. 2018; 7132–7141.

[13]  Zhang H, Cisse M, Dauphin YN, Lopez-Paz D. Mixup: Beyond empirical risk minimization. ICLR 2017;.

[14]  Bradycardia. Bradycardia, 2010. URL https://en.wikipedia.org/wiki/Bradycardia. [Online; accessed 31-August-2021].

Address for correspondence:

Hansheng Ren
3 Research Link, Singapore 117602
hanshengren@u.nus.edu

Miao Xiong
3 Research Link, Singapore 117602
miao.xiong@u.nus.edu