

# ECG Classification Report

Omenetti Matteo<sup>1</sup> and Srinivasan Ravi<sup>2</sup>

<sup>1</sup>ETH Zurich - momenetti@ethz.ch

<sup>2</sup>ETH Zurich - rsrinivasan@ethz.ch

**Summary** This report is part of the first project for the course "Machine Learning for Health Care" hosted at ETH Zürich. The scope of this project was to build different models for ECG classification for two different datasets PTB and MIT\_BIH. All of the models described in this report can be found in "src/" and they have been trained (in a reasonable amount of time) using the GPUs/TPUs offered by Google COLAB-Pro. All of the images that summarize the architecture of the different models have been omitted due to the length constraints of the report, but can be found in the respective notebooks together the graphs for the training loss and accuracy and the confusion matrices for multiclass classification.

**Convolutional Neural Network (CNN)** For the vanilla CNN, a smaller model compared to the one provided as a baseline was built. The reason behind this shallower architecture is that we were curious to understand whether the decreased training cost would also result in a statistically relevant decrease in performance. However, as it will be reported in the two subsections below this is not the case.

**MIT\_BIH** The model on the Mit\_Bih dataset achieves a test accuracy score of 0.982, taking around 6 minutes to train. This model was trained for 25 epochs with the option to reduce the learning rate on plateau, early stopping and validation split 0.1. The reported accuracy show the performance of this shallower model are very close to the baseline model (accuracy 0.985) that instead takes more than 7 minutes to train on the same GPU.

**PTB** The model on the PTB dataset achieves a test accuracy score of 0.975, taking around 1 minute to train. This model was trained for 33 epochs with the option to reduce the learning rate on plateau, early stopping and validation split 0.1. The reported accuracy show, the performance of this shallower model is very close to

the baseline model (accuracy 0.991). The model achieved a AUROC score of 0.974 and a AUPRC score of 0.983.

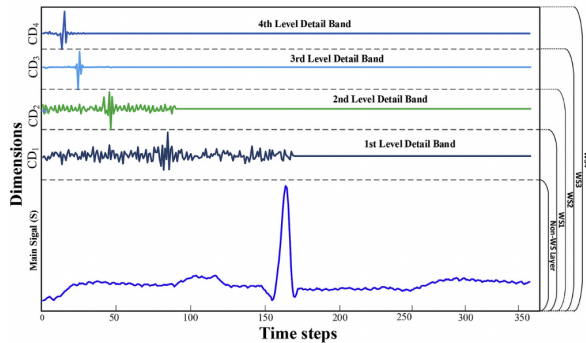
**Recurrent Neural Network (RNN)** This model consists of three sequential SimpleRNN blocks (respectively with dimension 128, 256 and 100), followed by a dense layer to perform the classification. The model is trained by unrolling the RNN blocks to speed up the training with GPUs.

**MIT\_BIH** The training completed in 40 epochs (40s/epochs) with early stopping and validation split set at 0.2. This model on the Mit\_Bih dataset achieves an accuracy of 0.888.

**PTB** The training completed in 89 epochs (5s/epochs) with early stopping and validation split set at 0.2. This model on the Mit\_Bih dataset achieves an accuracy of 0.795. The model achieved a AUROC score of 0.857 and a AUPRC score of 0.940.

**DBLSTM-WS** This network consists of 2 Bidirectional LSTM blocks (dimension 128 and 64) both with return\_sequences=True and dropout rate at 0.01. A dense layer follows with dimension 128 and dropout rate 0.3 and ReLu as activation function. A final dense layer follows for classification. The network takes as input a matrix where the first layer represents the original signal, while the other rows each represent the DWT (discrete wavelet decomposition) of the previous row. To achieve this we use the pywavelet library and the Daubechies db6 wavelet member (the rows are padded to 0 to obtain same-length signals). In both cases we found that using 3 levels of discrete wavelet transforms. The main idea of using DWT of the signal is to isolate sub-bands of the signal to improve the classification performance of DLSTMs. In the WS layer, frequency sub-bands belonging to the signals were obtained for use in the ECG signals classification phase. These sub-bands are presented in sequence to

other layers without any additional calculation. Using the Bidirectional LSTM layer instead of the standard LSTM layer (which already improves performance compared to the standard RNN), both past and future situations of sequential inputs in a time frame are evaluated without delay. Adding more levels for the DWT transform caused overfitting in the training process (this could be related to the fact that adding more levels increases the complexity of the model).



**Figure 1: WS layer outputs for LSTM inputs of the DBLSTM-WS**

**MIT\_BIH** The model trained for 25 epochs (22sec/epochs) using early stopping, reduce LR on plateau and validation split 0.2. This model on the Mit\_Bih dataset achieves an accuracy of 0.985.

**PTB** The model trained for 41 epochs (2sec/epochs) using early stopping, reduce LR on plateau and validation split 0.2. This model on the PTB dataset achieves an accuracy of 0.986. The model achieved a AUROC score of 0.995 and a AUPRC score of 0.997.

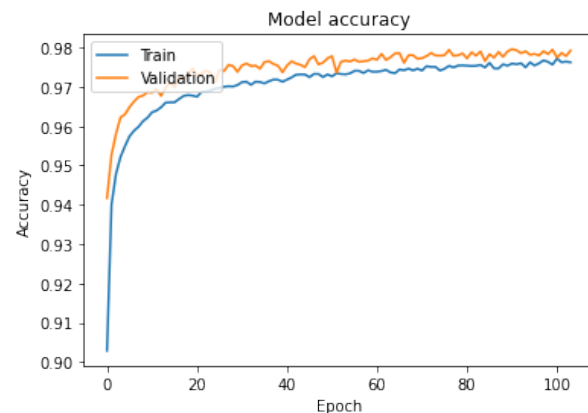
**CNN with Residual Blocks** For this task, we designed a CNN architecture similar to ResNet (but shallower). This is the model that achieves the best performance. This model compared to the vanilla CNN adds residual blocks after each BatchNormalization layer.

**MIT\_BIH** The model trained for 27 epochs (15sec/epochs) using early stopping, reduce LR on plateau and validation split 0.1. This model on the Mit\_Bih dataset achieves an accuracy of 0.988.

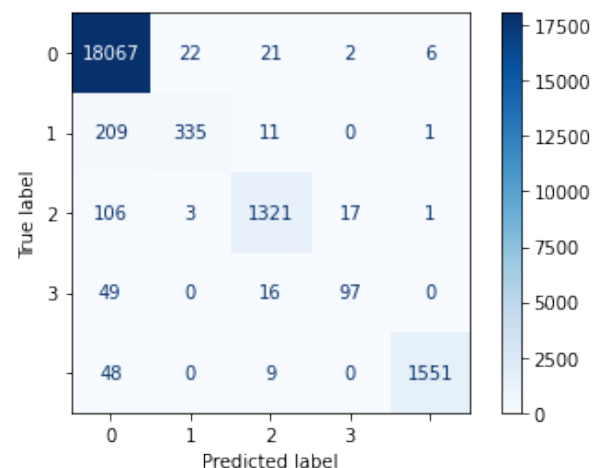
**PTB** The model trained for 44 epochs (2sec/epochs) using early stopping, reduce LR

on plateau and validation split 0.1. This model on the PTB dataset achieves an accuracy of 0.986. The model achieved a AUROC score of 0.857 and a AUPRC score of 0.94

**Transformer** This model consists of three transformers encoder blocks, stacked on top of each other with a standard CNN at the end. This model is very heavy to train. Without ColabPro the training for the Mit\_Bih dataset is around 16hours with the TPU. Using ColabPro, the amount of time required to train the model decreases drastically. Overall, its performance is not great if compared to the other much lighter models, therefore we would not recommend using this model in these particular tasks. Moreover, this model is also significantly slower than all of the other models at prediction time.



**Figure 2: Train and validation accuracy of the transformer for the Mit\_bih dataset.**

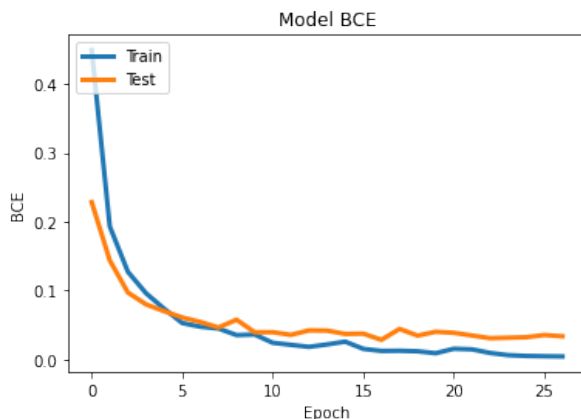


**Figure 3: Confusion matrix of the transformer for the Mit\_bih dataset.**

**MIT\_BIH** The model trained for 104 epochs (25sec/epochs) using early stopping and validation split 0.2. This model on the Mit\_Bih dataset achieves an accuracy of 0.976.

**PTB** The model trained for 146 epochs (4sec/epochs) using early stopping and validation split 0.2. This model on the PTB dataset achieves an accuracy of 0.968. The model achieved a AUROC score of 0.980 and a AUPRC score of 0.985.

**Transfer Learning** For this task we used the DBLSTM-WS model trained on the MIT\_BIH dataset. We froze the Bidirectional LSTM layers and then trained the last two dense layers on the PTB dataset (adapting the last layer to fit the binary classification task), thus making these layers perform feature extraction. The model trained for 27 epochs (3sec/epoch) using early stopping, reduce LR on plateau and validation split 0.2. This model on the PTB dataset achieves an accuracy of 0.992. The model achieved a AUROC score of 0.996 and a AUPRC score of 0.997. We also tried to freeze the first dense layer and only train the last dense obtaining an accuracy of 0.988. We can see that the performance increased significantly compared to training the DBLSTM-WS directly on the PTB dataset: the model is able to extract feature much better thanks to the bigger size and variety of the MIT\_bih database compared to the PTB database.



**Figure 4: Train and validation accuracy of the transfer learning for the PTB dataset.**

**Ensemble of models** For this optional ensemble task, we decided to use two different ap-

proaches. For the first approach, we just summed up the different prediction probabilities of every model, and then we took the class with the highest one. For the second approach, we ran a logistic regression model on the output probabilities of every model, in order to better understand the relation between the models. The logistic regression ensemble achieves the best performance among all the models used for this assignment with an accuracy for the PTB dataset of 99.1%. Once again ensembles of models of different kinds prove to be more robust than each single model taken separately. The performance of each approach is summarized in their respective section.

**Average of the outputs** For this ensemble, every model is loaded using its respective h5 file, the test set is predicted and the corresponding output probabilities are summed up with the previous probabilities.

**MIT\_BIH** This ensemble reaches an accuracy of 0.984.

**PTB** This ensemble reaches an accuracy of 0.990. The model achieved a AUROC score of 0.986 and a AUPRC score of 0.990.

**Logistic regression on the outputs** For this ensemble, every model is loaded using its respective h5 file, the test set is predicted. Both the training and test set are predicted and the corresponding output probabilities are summed up with the previous probabilities. The predictions for the training set are used as input features to a logistic regression model. Finally, given the output probabilities on the test set, the corresponding class is derived using predictive capabilities of the newly trained logistic regression model.

**MIT\_BIH** This ensemble reaches an accuracy of 0.987.

**PTB** This ensemble reaches an accuracy of 0.991. The model achieved a AUROC score of 0.988 and a AUPRC score of 0.991.

**Table 1: Summary of results - Accuracy**

|                     | MIT_BIH A | PTB   |
|---------------------|-----------|-------|
| Vanilla CNN         | 0.981     | 0.982 |
| Vanilla RNN         | 0.888     | 0.795 |
| DBLSTM-WS           | 0.985     | 0.986 |
| CNN Residual Blocks | 0.988     | 0.986 |
| Transformer         | 0.976     | 0.968 |
| Transfer Learning   | –         | 0.992 |
| Ensemble Average    | 0.984     | 0.990 |
| Ensemble Logistic   | 0.987     | 0.991 |

## References

- Li, J., Si, Y., Xu, T., & Jiang, S. (2018, December). Deep Convolutional Neural Network Based ECG Classification System Using Information Fusion and One-Hot Encoding Techniques. *Mathematical Problems in Engineering*, 2018, 1–10. Retrieved 2022-03-29, from <https://www.hindawi.com/journals/mpe/2018/7354081/> doi: 10.1155/2018/7354081
- Singh, S., Pandey, S. K., Pawar, U., & Janghel, R. R. (2018). Classification of ECG Arrhythmia using Recurrent Neural Networks. *Procedia Computer Science*, 132, 1290–1297. Retrieved 2022-03-14, from <https://linkinghub.elsevier.com/retrieve/pii/S1877050918307774> doi: 10.1016/j.procs.2018.05.045
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017, December). Attention Is All You Need. *arXiv:1706.03762 [cs]*. Retrieved 2022-03-29, from <http://arxiv.org/abs/1706.03762> (arXiv: 1706.03762)
- Yildirim, (2018, May). A novel wavelet sequence based on deep bidirectional LSTM network model for ECG signal classification. *Computers in Biology and Medicine*, 96, 189–202. Retrieved 2022-03-14, from <https://linkinghub.elsevier.com/retrieve/pii/S0010482518300738> doi: 10.1016/j.compbiomed.2018.03.016
- Zhang, H., Zhao, W., & Liu, S. (2020, December). SE-ECGNet: A Multi-scale Deep Residual Network with Squeeze-and-Excitation Module for ECG Signal Classification. *arXiv:2012.05510 [cs, eess]*. Retrieved 2022-03-29, from <http://arxiv.org/abs/2012.05510> (arXiv: 2012.05510)