# Reliability-Aware Hierarchical Learning for Chagas Detection from Electrocardiogram under Expert Label Scarcity

Hao Wen[1], Jingsu Kang[2]

[1]College of Science, China Agricultural University, Beijing, China
[2]Tianjin Medical University, Tianjin, China

## Abstract

*Aim: We present a reliability-aware hierarchical learning framework for ECG-based Chagas cardiomyopathy screening in the George B. Moody PhysioNet Challenge 2025 by Team Revenger, aiming to maximize positive case retrieval under prevalence constraints.*

*Methods: The 12-lead ECGs were resampled to 400 Hz, bandpass filtered (0.5–45 Hz), and z-score normalized. We used a ResNet model integrated with squeeze-and-excitation (SE) modules for binary classification. To address severe class imbalance and the scarcity of expert-confirmed labels, we applied stratified upsampling and reliability-weighted label smoothing to prioritize expert-confirmed positives over self-reported ones. Model training used an asymmetric loss to further penalize false negatives and was optimized with AdamW and a OneCycle learning rate scheduler. Model selection was based on the Challenge score from an internal hold-out subset.*

*Results: On the official hidden test set, our method received a Challenge score of 0.163, ranking 32nd of 40 eligible teams.*

*Conclusion: The proposed method demonstrates effective performance for ECG-based Chagas screening, highlighting its potential for improving detection accuracy and reliability in resource-limited scenarios.*

## 1. Introduction

Addressing underdiagnosis of Chagas disease through scalable ECG-based screening is the focus of the 2025 George B. Moody PhysioNet Challenge [1–3]. Enabled by aggregated multi-cohort ECG datasets [4–8], the Challenge frames a multi-source learning setting with heterogeneous label reliability and severe class imbalance.

In this work, we propose a reliability-aware hierarchical framework that prioritizes expert-confirmed labels and mitigates severe class imbalance within a deep ECG model, with optimization aligned to prevalence-constrained sensitivity objectives.

## 2. Methods

### 2.1. Datasets and Preprocessing

We used three ECG datasets for model training, with substantial differences in sample size, Chagas prevalence, and label provenance as summarized in Table 1.

| Dataset | Size | Chagas rate | Label provenance |
|---|---|---|---|
| SaMi-Trop | 1 631 | 100 % | expert-confirmed |
| CODE-15% | 345 779 | 1.795 % | self-reported |
| PTB-XL | 21 799 | 0 % | N/A |

Table 1: Dataset statistics and label provenance. Chagas rate is the proportion of recordings labeled positive in each dataset. N/A indicates that confirmed Chagas cases are not expected (non-endemic population).

All ECGs were uniformly resampled to 400 Hz, bandpass filtered (0.5–45 Hz), and z-score normalized to zero mean and unit variance computed as in Eq. 1:

$$\tilde{x} = \frac{x - \mu_x}{\sigma_x} \tag{1}$$

where $x$ is the original ECG signal, $\mu_x$ and $\sigma_x$ are the mean and standard deviation of $x$, respectively. We excluded ECGs shorter than 1200 samples to ensure inputs contain enough cardiac cycles for stable model analysis.

### 2.2. Reliability-Aware Hierarchical Supervision

We introduce a hierarchical supervision scheme that encodes source reliability through stratified label smoothing and adaptive upsampling. Three reliability levels are defined: (1) expert-confirmed positives (SaMi-Trop, maximal trust), (2) self-reported samples (CODE-15%, both positives and negatives, higher uncertainty), and (3) non-endemic negatives (PTB-XL, very low true prevalence but still mildly regularized). Given a one-hot label $y$ ($[0, 1]$ for positives, $[1, 0]$ for negatives) and number of classes
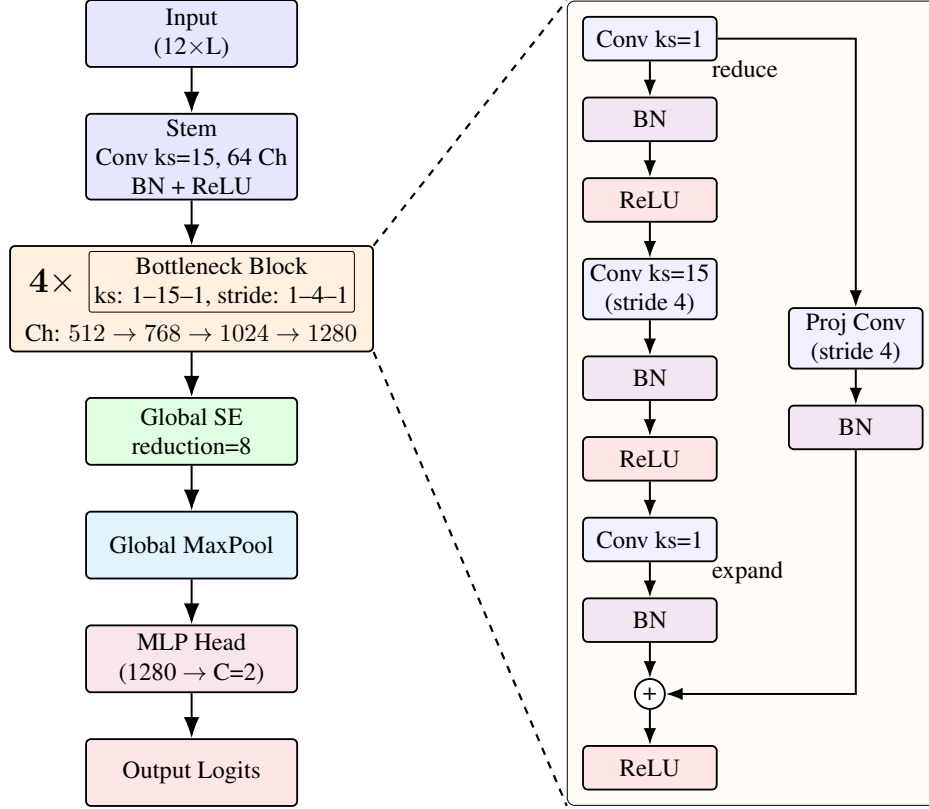
Figure 1: Model architecture. Left: overall network: a stem (Conv1d, kernel size 15, 64 output channels, Batch Normalization, ReLU) followed by four bottleneck residual blocks, a global SE module, global max pooling, and an MLP head producing $C = 2$ logits. Channel widths shown ($512 \rightarrow 768 \rightarrow 1024 \rightarrow 1280$) are the expanded channels. Right: internal bottleneck structure. The middle convolution of kernel size 15 uses a stride of 4 for temporal downsampling; kernel size 1 convolutions reduce and then expand channels, and a projection convolution (kernel size 1, stride 4) aligns resolution and width for the residual path. For clarity, dropout layers present in the implementation are omitted. Abbreviations: ks kernel size; Ch channels; BN Batch Normalization; SE squeeze-and-excitation; MLP multi-layer perceptron.

$C = 2$, the smoothed target is computed as in Eq. 2:

$$\tilde{\mathbf{y}} = (1 - \varepsilon) \cdot \mathbf{y} + \frac{\varepsilon}{C} \cdot \mathbf{1}, \qquad (2)$$

where $\varepsilon$ is the smoothing factor which depends on the reliability level: 0.0 (SaMi positives), 0.6 (CODE-15% positives & negatives), 0.2 (PTB-XL negatives). This attenuates overconfident gradients for noisier or potentially misreported labels while preserving sharp supervision on expert-confirmed cases.

Severe class imbalance was mitigated by upsampling positives during training: positive samples from CODE-15% were upsampled by a factor of 3, and those from SaMi-Trop by 12. No upsampling was applied to PTB-XL, which contains no positives. We chose these factors after reviewing hidden validation scores from multiple submissions, as shown in Table 2.

Smoothed labels and upsampling strategies for each dataset are summarized in Table 3.

| Upsample factor | | Challenge score |
|---|---|---|
| CODE-15% | SaMi-Trop | |
| - | - | 0.239[†] |
| 3 | 7 | 0.212 |
| 3 | 12 | **0.245** |
| 10 | 120 | 0.210 |
| 6 | 36 | 0.221 |

Table 2: Representative upsampling schemes and corresponding Challenge scores on the hidden validation set. The model and training strategies used were the same. "-" indicates no upsampling.
[†] obtained during the unofficial phase.

## 2.3. Model Architecture

We build upon the 1D ResNet ECG classifier of Ribeiro et al. [4] and introduce three modifications.

**(1) Bottleneck residual blocks.** We replace basic

| Dataset | Upsampling | Smoothed labels | |
| | factor | negative | positive |
| --- | --- | --- | --- |
| SaMi-Trop | 12 | N/A | [0, 1] |
| CODE-15% | 3 | [0.7, 0.3] | [0.3, 0.7] |
| PTB-XL | 1 | [0.9, 0.1] | N/A |

Table 3: Smoothed labels (computed from Eq. 2) and up-sampling strategies for each dataset.

ResNet blocks with bottleneck blocks of kernel sizes 1–15–1 (pointwise–temporal–pointwise). The middle temporal convolution applies a stride of 4 for downsampling; the two convolutions (kernel size 1) first reduce the number of channels and then expand them with an expansion factor of 4. A projection convolution (kernel size 1, stride 4) is used in the residual branch whenever temporal resolution or channel width changes. Across the four blocks, the reduced (bottleneck) channel widths are $128 \rightarrow 192 \rightarrow 256 \rightarrow 320$, yielding expanded output widths $512 \rightarrow 768 \rightarrow 1024 \rightarrow 1280$.

**(2) Global squeeze-and-excitation (SE).** After the final bottleneck block, a single global SE module (reduction ratio 8) [9] performs temporal average pooling to a channel descriptor, applies a two-layer bottleneck multi-layer perceptron (MLP) $1280 \rightarrow 160 \rightarrow 1280$ with ReLU and sigmoid gating, and rescales the feature map channel-wise.

**(3) Global pooling head for variable input length.** Instead of flattening a fixed-length feature map as in the original baseline, we apply global max pooling over the remaining temporal dimension, yielding a 1280-dimensional vector irrespective of input length $L$. This vector is fed to a lightweight two-layer classification MLP: a hidden fully connected layer ($1280 \rightarrow 1024$) with non-linear activation and dropout (rate 0.2), followed by a final linear layer ($1024 \rightarrow 2$) producing class logits.

**(4) Stem and regularization.** A stem Conv1d (kernel size 15, stride 1, 64 channels) with BatchNorm and ReLU precedes the bottleneck stack. Within each bottleneck block, we apply BatchNorm+ReLU after the first two convolutions and dropout (rate 0.2) after each of those activations. All convolutions use "same" padding to preserve temporal length before downsampling operations.

The overall model architecture is illustrated in Fig. 1.

## 2.4. Training and Implementation Setups

We employed asymmetric loss (ASL) [10] to complement the reliability-aware label smoothing strategy, jointly addressing the challenges of severe class imbalance. Let $\mathbf{z} = (z_0, z_1)$ denote the logits and $p = \mathrm{softmax}(\mathbf{z})_1$ the predicted probability of the positive class. The ASL is defined in Eq. 3 with separate focusing parameters for positives and negatives and a clipped negative probability term:

$$L = -y \cdot (1 - p)^{\gamma_+} \log(p) \\ - (1 - y) \cdot (p_m)^{\gamma_-} \log(1 - p_m), \tag{3}$$

where $y$ is the (smoothed) positive-class target probability, $p_m = \max(p - m, 0)$, $(\gamma_+, \gamma_-) = (1, 4)$ and margin $m = 0.05$. We train for 30 epochs with batch size 128 using the AdamW optimizer (initial learning rate $1 \times 10^{-4}$, peak $6 \times 10^{-4}$ under a OneCycle scheduler, weight decay $1 \times 10^{-2}$). Early stopping (patience 10 epochs, monitored on a fixed 20% internal hold-out subset) selects the final model via the Challenge metric. Each training segment is a uniform random crop (or center padding if shorter) of length 4096 samples. The full implementation, including model construction, data pipeline, and optimization utilities, is based on the `torch-ECG` framework [11].

## 3. Results

The Challenge score of our team "Revenger" on the hidden test set was 0.163, ranking 32nd among 40 eligible teams. This score and ranking, along with extra scores on the internal hold-out of the public training data, and on the hidden validation set, are summarized in Table 4.

| Training | Validation | **Test** | **Ranking** |
| --- | --- | --- | --- |
| $0.451 \pm 0.005$ | 0.245 | **0.163** | **32 / 40** |

Table 4: Challenge scores for our submitted entries (team "Revenger"). Training: internal hold-out mean $\pm$ std over repeated runs. Validation: best among 10 validation submissions. Test: the unique test submission. Ranking: position on the hidden test leaderboard.

## 4. Discussion and Conclusions

The Challenge scores presented in Table 4 indicate that our proposed method is effective for Chagas screening from ECGs, albeit with substantial room for improvement. The result demonstrates our model's ability to learn diagnostically relevant features from ECGs for this task under scarce and noisy supervision. This is achieved through reliability-aware label smoothing, which incorporates both label provenance and reliability instead of treating all positive labels uniformly. Together with the asymmetric loss and strategic upsampling, these results indicate that explicitly modeling label reliability helps stabilize the learning process more effectively than introducing additional architectural complexity. Overall, our approach offers a scalable and resource-efficient solution and aligns well with the Challenge's objective of identifying high-risk individuals under limited serological testing capacity.

However, the performance gap between our internal hold-out subset and the hidden test set indicates limitations

in our method's generalization capability. This is also reflected in the absolute performance, where our Challenge score is substantially behind those of the top-performing teams (e.g., 0.323, 0.283, 0.280). The leading approaches generally leveraged some form of large-scale representation learning, such as pre-trained Vision Transformer foundation models or self-supervised learning on extensive ECG datasets. These strategies focus on building robust, general-purpose feature extractors. In contrast, our work prioritized a distinct niche by designing a resource-efficient pipeline that explicitly addresses the challenges of label noise and scarcity through reliability-aware smoothing and asymmetric loss, without relying on massive pre-training. While this choice enhances practicality and stability under noisy supervision, it appears that the representational capacity of our directly-trained model is ultimately lower, limiting its ability to generalize as effectively as the foundation-model-based approaches. Furthermore, our static assignment of reliability weights, which cannot adapt to instance-specific label quality, represents another limitation compared to more dynamic or learned weighting schemes.

Building upon these insights, future research will focus on bridging the generalization gap while retaining our focus on learning under weak supervision. A primary direction is the self-adaptive supervision framework. This includes dynamic weighting schemes, moving beyond our current static factors, and adaptive sampling strategies that respond to the model's evolving confidence during training. To further improve robustness and handle class imbalance, advanced data augmentation techniques such as CutMix [12] and SMOTE [13] will be explored to diversify the limited positive samples, with a particular focus on ECG-specific transformations like lead-wise masking. Furthermore, inspired by the success of representation learning, we plan to explore self-supervised pre-training on large-scale unlabeled ECG data as a pivotal step to learn more transferable representations before fine-tuning. This direction, while computationally more demanding, addresses a key limitation identified in our current work. Additionally, a multi-task learning framework leveraging auxiliary arrhythmia labels could be integrated to impose clinically meaningful constraints and enhance the feature representation for the primary Chagas screening task.

## References

[1] Goldberger AL, Amaral LA, Glass L, Hausdorff JM, Ivanov PC, Mark RG, et al. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals. Circulation 2000;101(23):e215–e220.

[2] Reyna MA, Koscova Z, Pavlus J, Weigle J, Saghafi S, Gomes P, et al. Detection of Chagas Disease from the ECG: The George B. Moody PhysioNet Challenge 2025. In Computing in Cardiology 2025, volume 52. 2025; 1–4.

[3] Reyna MA, Koscova Z, Pavlus J, Saghafi S, Weigle J, Elola A, et al. Detection of Chagas Disease from the ECG: The George B. Moody PhysioNet Challenge 2025, 2025. URL https://arxiv.org/abs/2510.02202.

[4] Ribeiro AH, Ribeiro MH, Paixão GM, Oliveira DM, Gomes PR, Canazart JA, et al. Automatic Diagnosis of the 12-lead ECG Using a Deep Neural Network. Nature Communications 4 2020;11(1):1–9.

[5] Cardoso CS, Sabino EC, Oliveira CDL, de Oliveira LC, Ferreira AM, Cunha-Neto E, et al. Longitudinal Study of Patients with Chronic Chagas Cardiomyopathy in Brazil (SaMi-Trop Project): A Cohort Profile. BMJ Open 5 2016; 6(5):e011181. ISSN 2044-6055.

[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] Nunes MCP, Buss LF, Silva JLP, Martins LNA, Oliveira CDL, Cardoso CS, et al. Incidence and Predictors of Progression to Chagas Cardiomyopathy: Long-Term Follow-Up of Trypanosoma Cruzi –Seropositive Individuals. Circulation 11 2021;144(19):1553–1566. ISSN 1524-4539.

[8] Pinto-Filho MM, Brant LC, dos Reis RP, Giatti L, Duncan BB, Lotufo PA, et al. Prognostic Value of Electrocardiographic Abnormalities in Adults from the Brazilian Longitudinal Study of Adults' Health. Heart 12 2020; 107(19):1560–1566. ISSN 1468-201X.

[9] Hu J, Shen L, Sun G. Squeeze-and-Excitation Networks. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, 6 2018; 7132–7141.

[10] Ridnik T, Ben-Baruch E, Zamir N, Noy A, Friedman I, Protter M, et al. Asymmetric Loss for Multi-Label Classification. In Proceedings of the IEEE/CVF International Conference on Computer Vision. IEEE, 10 2021; 82–91.

[11] Wen H, Kang J. A Novel Deep Learning Package for Electrocardiography Research. Physiological Measurement 11 2022;43(11):115006. ISSN 1361-6579.

[12] Yun S, Han D, Oh SJ, Chun S, Choe J, Yoo Y. CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features. In Proceedings of the IEEE/CVF International Conference on Computer Vision. Institute of Electrical and Electronics Engineers (IEEE), 10 2019; 6022–6031.

[13] Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: Synthetic Minority Over-sampling Technique. Journal of Artificial Intelligence Research 6 2002; 16(1):321–357. ISSN 1076-9757.

Address for correspondence:

Hao Wen
No. 17, Qinghua East Road, Haidian District, Beijing, China
wenh06@cau.edu.cn,wenh06@gmail.com