

Deep Learning for Early Chagas Disease Diagnosis: A Comparative Analysis of 12-Lead ECG and Derived VCG

Alejandro Pascual-Mellado¹, Vicent Tores-Sastre¹, Cristina Albert¹, Alejandro Pérez-González¹, Raúl Alós¹, Elisa Ramírez¹, Francisco Castells¹, José Millet¹

¹ EP Analytics Lab - ITACA Institute, Universitat Politècnica de València

Abstract

Chagas disease (ChD) is a chronic parasitic condition that can lead to severe cardiac complications. The use of electrocardiographic (ECG) analysis has emerged as a promising tool for early, non-invasive detection. This work, developed by the EPBandoleroLab team for the PhysioNet Challenge 2025, presents a deep learning approach for ChD classification using the CODE-15, SaMi-Trop, and PTB-XL databases. Our methodology explores the effectiveness of different signal representations, comparing the standard 12-lead ECG with the derived Vectorcardiogram (VCG). Furthermore, we address the significant class imbalance through a controlled sampling strategy. Our findings indicate that the model performs best when trained on the full 12-lead ECG representation with a moderately imbalanced dataset. This configuration achieved a Challenge Score of 0.259 in the official phase, placing our team in the top half of all competitors.

1. Introduction

ChD, or American trypanosomiasis, is a neglected tropical disease caused by the protozoan *Trypanosoma cruzi*. It affects over 7 million people worldwide, mainly in Latin America, and leads to severe cardiac complications, including heart failure and sudden death, especially in its chronic phase [1]. Transmission occurs predominantly via triatomine insects ("kissing bugs"), but also through congenital routes, blood transfusion, and ingestion of contaminated food [1, 2].

Although the acute phase is often asymptomatic, a significant proportion of patients develop chronic ChD cardiomyopathy, characterized by ventricular dysfunction, thromboembolism, arrhythmias and dysautonomia [3]. Early diagnosis is crucial but is often hindered by limited access to serological testing in rural or under-resourced settings [4].

ECG, as a low-cost and widely available diagnostic tool, holds promise for detecting early signs of chronic ChD car-

diomyopathy. Certain alterations in ECG, such as right bundle branch block, premature ventricular beats, ST-T changes, abnormal Q waves, various degrees of AV block, sick sinus syndrome and low QRS voltage, may suggest ChD even in asymptomatic individuals [3, 5].

The 2025 PhysioNet Challenge focuses precisely on this problem: detecting ChD disease using standard 12-lead ECG recordings via machine learning and deep learning methods[6].

Recent studies have demonstrated the significant promise of applying artificial intelligence (AI) to ECG-based disease detection. Notably, deep neural networks have been shown to outperform medical residents in classifying various cardiac conditions [7] and have successfully automated the classification of numerous arrhythmias with high accuracy [8]. These findings underscore the growing potential of deep learning to support clinical diagnosis in cardiology.

This challenge aims to develop scalable, AI-powered tools for the early, preclinical diagnosis of ChD from a 10-second ECG, enabling timely intervention for at-risk individuals before irreversible cardiac damage occurs.

2. Materials

The dataset provided for the PhysioNet Challenge 2025 comprises three distinct ECG databases, each offering unique characteristics relevant to the Chagas detection task [6]:

- **CODE-15%:** A subset of the larger CODE cohort, it includes over 300,000 ECGs collected in Brazil between 2010 and 2016. The dataset is provided in 18 distinct partitions. Each ECG lasts approximately 7.3 to 10.2 seconds and was recorded at around 400 Hz. The Chagas labels in this set are self-reported (i.e., weak labels) with unknown accuracy and low prevalence, adding real-world noise and variability to the training data [9].

- **SaMi-Trop:** This cohort includes 1,631 12-lead ECGs collected between 2011 and 2012 from Brazilian patients in endemic areas. Importantly, Chagas diagnoses are sero-

logically confirmed, making this the only strongly labeled positive dataset [10].

- **PTB-XL:** This European dataset consists of 21,799 ECGs, recorded at 500 Hz with 10-second duration. Due to its geographical origin (Germany), it is assumed to contain only Chagas-negative cases, offering a high-quality negative class with low noise [11].

To create a robust **external test** set for evaluating generalization, we randomly held out two of the eighteen CODE-15% partitions. This held-out set was reserved exclusively for final testing and was not used during training or hyperparameter tuning.

The remaining data, including the other 16 CODE-15% partitions, SaMi-Trop, and PTB-XL constituted our development set. From this set, various experimental subsets were generated using a data sampling strategy, which is further detailed in the Methods section. Each of these subsets was then split at the patient level into training (70%) and validation (30%) partitions.

The final test set, used for the official evaluation and ranking of competitors, remains hidden and is exclusively accessed by the event organizers.

3. Methods

3.1. Data Selection and Preprocessing

Due to a class imbalance favoring negative cases, a data selection strategy was implemented. While all positive records from the three databases were included, negative samples were filtered based on demographic characteristics (age and gender). This allowed us to control the ratio between negative and positive records using a balancing parameter R to experimentally investigate whether a controlled imbalance is beneficial for the model’s generalization.

The signals were preprocessed by resampling to 400 Hz, followed by a two-stage filtering process (median and wavelet) to remove noise. Each lead was then standardized using Z-score normalization. To create fixed-length inputs, signals were segmented into 1024-sample windows. Patient-level data splits were enforced to prevent data leakage, and the final inference for a record is the average of its segment probabilities.

3.2. Signal Representation: ECG vs. VCG

To determine the most effective input representation, an experimental comparison was conducted between two modalities. The first is the standard 12-lead ECG, which provides detailed temporal information of the cardiac vector voltage from multiple anatomical perspectives.

As an alternative, the VCG was evaluated, a three-dimensional representation of the heart’s electrical activity mathematically derived from the 12 leads using the inverse Dower transformation matrix. The VCG projects the information onto three orthogonal axes (X, Y, Z), offering a spatial view of the cardiac vector.

The underlying hypothesis is that the VCG, by eliminating the inherent redundancy among ECG leads, could allow the model to learn global diagnostic features more efficiently [12]. However, it is important to acknowledge that the Dower transformation is an estimate of the true cardiac vector, not a direct measurement, which entails the risk of introducing slight signal distortion. This experiment therefore investigates whether the benefit of a compact and non-redundant representation outweighs the potential loss of fidelity, building on previous studies indicating that key diagnostic information is largely preserved.

3.3. Hybrid Model Architecture

A hybrid architecture was designed to use a CNN for feature extraction and a Transformer for contextual modeling. The Transformer’s output is fed into a Multilayer Perceptron (MLP) for the final classification. Figure 1 shows a diagram of this architecture.

The stack of 1D convolutional layers processes the input signal. This part acts as a local feature extractor, learning to identify morphological patterns in the signal. Through layers of convolution, normalization, and pooling, the CNN transforms the signal into a shorter, denser sequence of feature vectors.

The feature sequence generated by the CNN is fed into a Transformer encoder. This component, through its multi-head self-attention mechanism, models long-term temporal dependencies in the signal. A special classification token ($[CLS]$) is prepended to the sequence to aggregate the contextual information of the entire segment into a single representation vector.

The feature vector corresponding to the $[CLS]$ token is then passed to a final MLP. After an initial normalization, the MLP projects the input through dense layers with Dropout to produce a single logit, which is mapped to a probability using a sigmoid function.

3.4. Training and Evaluation Strategy

For the experimentation, multiple development sets were generated by varying the balance ratio and the signal representation (ECG vs. VCG). Each configuration was trained and its hyperparameters optimized using its own training and validation subsets. The final performance of each optimized configuration was evaluated on the external test set to ensure a fair and rigorous comparison.

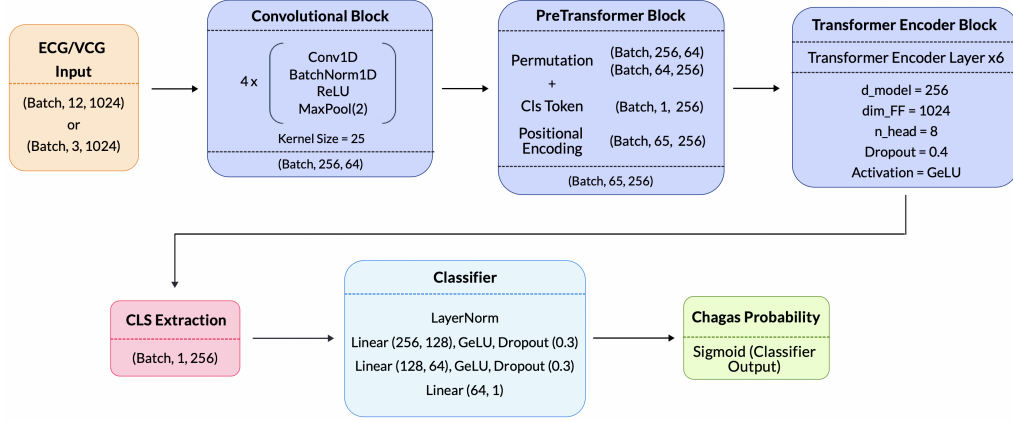


Figure 1. Model Architecture Scheme.

Each experimental configuration was trained using the Adam optimizer, a Focal Loss function to address class imbalance, and regularization techniques such as Lead Dropout and a learning rate scheduler. To optimize the model, an empirical hyperparameter search was conducted (including learning rate, weight decay, alpha and gamma), selecting the combination that maximized the AUPRC on the validation set. Early stopping was employed to prevent overfitting during this process.

4. Results and Discussions

4.1. Data Distribution and Sampling Strategy

Table 1. Distribution of classes in the datasets used for training and evaluation.

Dataset	N Positives	N Negatives	Prevalence
Ratio 1:1	7392	7392	50.0%
Ratio 3:1	7392	22176	25.0%
Ratio 5:1	7392	36960	16.7%
External test	798	39003	2%

Table 1 summarizes the composition of the datasets used. It details the three training set configurations generated by varying the balance ratio, along with the class distribution of the fixed external test set, which has a 2% prevalence of positives. This structure allows us to evaluate how different training data compositions affect performance in a realistic and consistent testing scenario.

4.2. Model Performance Analysis

The analysis of the results in Table 2 reveals two key findings. First, increasing the balance ratio during training consistently improves performance. For example, in

Table 2. Performance in the external test set of the best hyperparameters configuration found for each combination of signal representation and balance ratio.

ID	Repr.	Ratio	AUROC	Challenge Score	AUPRC
M1	ECG-12	1:1	0.805	0.363	0.139
M2	ECG-12	3:1	0.819	0.388	0.153
M3	ECG-12	5:1	0.811	0.405	0.160
M4	VCG-3	1:1	0.791	0.362	0.130
M5	ECG-3	3:1	0.805	0.377	0.146
M6	VCG-3	5:1	0.810	0.402	0.159

the ECG-12 representation, the Challenge Score rises from 0.363 (ratio 1:1) to 0.405 (ratio 5:1), demonstrating that greater exposure to the diversity of negative cases benefits model generalization.

Second, the 12-lead ECG representation slightly outperforms the VCG. Although the VCG is a compact representation, its performance is consistently lower, as observed in the Challenge Score (0.405 for ECG-12 vs. 0.402 for VCG-3 with a 5:1 ratio). A possible explanation is that the Dower transformation, being an estimate, may introduce subtle distortions that degrade diagnostic information.

Consequently, the M3 model (ECG-12, ratio 5:1) emerges as the optimal configuration, underscoring the importance of maximizing the volume of training data while preserving the full richness of the original input signal.

Table 3. Final comparison of best performance by signal representation in the external test set (2% prevalence).

Métrica	ECG-12 (M3)	VCG-3 (M6)
Challenge Score	0.405	0.402
AUPRC	0.160	0.159
AUC (AUROC)	0.811	0.810
F1-score	0.152	0.146
Precisión	0.088	0.842
Recall	0.564	0.551

Table 3 presents a direct performance comparison between our two best final configurations—one based on 12-lead ECG (M3) and the other on 3-lead VCG (M6)—both evaluated on the challenging external test set with a 2% positive prevalence.

The analysis reveals that the model using the 12-lead ECG representation demonstrates a consistent, albeit slight, superiority across most key metrics. This reinforces that while the VCG representation is more compact, the unaltered signal information present in the full 12 leads is beneficial for the model’s discriminative capacity in this task.

The performance of the M3 model (ECG-12, 5:1 ratio) is representative of an effective classifier in a realistic screening scenario. It achieves a recall of 0.564, identifying more than half of the affected individuals. The consequence of this sensitivity in a low-prevalence environment is a precision of 0.088. This value, although numerically low, is more than four times higher than the 2% baseline prevalence, demonstrating that the model’s alerts are highly informative. The F1-score of 0.152 encapsulates this inherent trade-off between detecting positive cases and controlling false alarms.

5. Conclusion

This study demonstrates the potential of deep learning as a tool for the early diagnosis of ChD. Our model can function as an effective initial screening system. With a sensitivity exceeding 56%, it could alert clinicians to perform a confirmatory serological test on more than half of the affected patients.

The main limitation is its low precision, which leads to a high rate of false positives—a common challenge in low-prevalence problems. However, we propose its use as a clinical decision support system, where the model’s alerts act as a “second reader” to motivate a more thorough review of the case by a specialist. This synergy between AI and clinical expertise represents a promising avenue for improving early detection and preventing irreversible cardiac damage.

The competitiveness of this approach was validated in the official phase of the Physionet Challenge, where our model achieved a Challenge Score of 0.259, placing us in the top half of all participating teams.

Code Availability

The source code for the models and experiments presented in this paper is publicly available on GitHub at: EP-BandoleroLab Team Code

Acknowledgments

This work was supported by PID2022-142514OB-I00 (National Research Program, Ministerio de Cien-

cia e Innovación, Spanish Government) and CIBERCV CB16/11/00486 (Instituto de Salud Carlos III).

References

- [1] World Health Organization. Chagas disease (american trypanosomiasis), 2024.
- [2] Prata A. Clinical and epidemiological aspects of chagas disease. *The Lancet Infectious Diseases* 2001;1(2):92–100.
- [3] Nunes M, Dones W, Morillo C, Encina J, Ribeiro A. Chagas disease: an overview of clinical and epidemiological aspects. *Journal of the American College of Cardiology* 2013;62(9):767–776.
- [4] Ardiles-Ruesjas S, Lesmo V, González-Romero V, Cubilla Z, Chena L, Huber C, et al. Prevalence and diagnostic accuracy of different diagnostic tests for chagas disease in an indigenous community of the paraguayan chaco. *PLoS Neglected Tropical Diseases* 2025;19(2):e0012861.
- [5] Rojas L, Glisic M, Pletsch-Borba L, Echeverría L, et al. Electrocardiographic abnormalities in chagas disease in the general population: A systematic review and meta-analysis. *PLoS Neglected Tropical Diseases* 2018;12(6):e0006567.
- [6] Clifford G, Goldberger A, the Challenge Organizers. The 2025 physionet/computing in cardiology challenge, 2025.
- [7] Ribeiro A, Ribeiro A, Paixão G, Oliveira D, Gomes P, Canazart J, et al. Automatic diagnosis of the 12-lead ecg using a deep neural network. *Nature Communications* 2020; 11(1):1760.
- [8] Kim H, Sunwoo M. An automated cardiac arrhythmia classification network for 45 arrhythmia classes using 12-lead electrocardiogram. *IEEE Access* 2024;12:44527–44538.
- [9] Ribeiro A, Paixao G, Lima E, Horta Ribeiro M, Pinto Filho M, Gomes P, et al. Code-15 <https://doi.org/10.5281/zenodo.4916206>, 2021.
- [10] Ribeiro A, Ribeiro A, Paixao G, Lima E, Horta Ribeiro M, Pinto Filho M, et al. Sami-trop: 12-lead ecg traces with age and mortality annotations (1.0.0). <https://doi.org/10.5281/zenodo.4905618>, 2021.
- [11] Wagner P, Strodthoff N, Bousseljot R, Samek W, Schaeffer T. Ptb-xl, a large publicly available electrocardiography dataset (version 1.0.3). <https://doi.org/10.13026/kfzx-aw45>, 2022.
- [12] Ramirez E, Ruiperez-Campillo S, Casado-Arroyo R, Merino JL, Vogt JE, Castells F, Millet J. The art of selecting the ecg input in neural networks to classify heart diseases: a dual focus on maximizing information and reducing redundancy. *Frontiers in Physiology* 2024;Volume 15 - 2024. ISSN 1664-042X. URL <https://www.frontiersin.org/journals/physiology/articles/10.3389/fphys.2024.1452829>.

Address for correspondence:

Alejandro Pascual-Mellado
ale.pas.mel@gmail.com