# Prune or Quantize? Layer-wise Compression of Time-Series ECG Foundation Networks

# **Anonymous Author(s)**

Affiliation Address email

#### **Abstract**

Foundation models for biosignals, such as wearable ECG monitors, face challenges in resource-constrained settings due to high memory and computational demands. We propose an adaptive layer-wise compression framework that combines quantization and pruning to reduce model size while preserving predictive performance. Layer importance, estimated via parameter contribution and weight variance, guides fine-grained assignment of bit-widths and pruning thresholds, balancing efficiency and accuracy across high- and low-sensitivity layers. Experiments on Chapman and CPSC ECG datasets show that our method consistently outperforms fixed global compression schemes, achieving up to 10.44× compression with no loss in performance. Our architecture-agnostic framework scales from lightweight residual networks to large foundation models, enabling real-time, low-resource ECG monitoring. By efficiently deploying foundation models on edge devices, this work advances scalable, physiology-aware biosignal AI for mobile health and clinical applications.

## 1 Introduction and Related Work

2

3

5

6

8

9

10

11

12

13

14

Electrocardiography (ECG) is a cornerstone of cardiac health assessment, capturing the heart's electrical activity through body-surface electrodes to reveal characteristic waveforms Trobec et al. (2018). These signals enable detection of arrhythmias, ranging from asymptomatic to life-threatening conditions like sudden cardiac death Srinivasan and Schilling (2018). Traditional rule-based diagnostics struggle with the scale and complexity of physiological data, driving demand for automated, cost-effective ECG monitoring Ebrahimi et al. (2020).

Deep learning (DL) has transformed arrhythmia detection, with convolutional neural networks 22 (CNNs) and recurrent architectures achieving high accuracy Kiranyaz et al. (2015); Alzubaidi et al. (2021). Recent innovations include transforming ECGs into images de Santana et al. (2021), CNN-LSTM hybrids Tan et al. (2018), and attention-based or transformer-based models El-Ghaish and 25 26 Eldele (2024); Jin et al. (2021). However, these approaches often lack generalization across diverse populations or robustness to class imbalance Hannun et al. (2019). Foundation models, pretrained 27 on large-scale unlabeled data via self-supervision, have revolutionized NLP, vision, and audio by 28 enabling robust generalization across tasks and domains Radford et al. (2018); He et al. (2022); 29 Hsu et al. (2021). In medicine, models like CheXzero Tiu et al. (2022), MedSAM Ma et al. (2024), 30 31 and ECGFounder Li et al. (2024) leverage large-scale biosignal data for improved transferability. However, their computational complexity and reliance on supervised pretraining with limited cohorts 32 hinder deployment in resource-constrained settings, such as wearable ECG devices.

Deploying DL models on wearables is limited by memory, energy, and latency constraints Chen 34 and Ran (2019). Compression techniques like pruning Frankle and Carbin (2018), quantization 35 Hubara et al. (2018); Krishnamoorthi (2018), and binarization Courbariaux et al. (2015) enable 36 lightweight deployment. Quantization reduces parameter precision, acting as a regularizer that 37 preserves discriminative capacity in noisy biosignals Liu et al. (2022b). Recent advances, such as 38 nonuniform-to-uniform quantization (N2UQ) Liu et al. (2022b), adapt bin widths to data distributions, 39 achieving near full-precision accuracy. Knowledge distillation Hinton et al. (2015) and neural 40 architecture search Tan (2019) further optimize model efficiency, but their application to biosignal foundation models remains underexplored. 42

#### Our contributions are:

- The first adaptive compression framework for ECG foundation models, enabling edge deployment with up to  $10 \times$  size reduction for both traditional and foundational models.
- A ResNet1D achieving state-of-the-art arrhythmia classification with high compression.
- A demonstration that compressed foundation models maintain clinical accuracy while reducing computational costs by an order of magnitude, advancing scalable biosignal AI.

# 9 2 Method

- We aim to design scalable, interpretable biosignal foundation models that balance physiological fidelity with edge deployment efficiency. Our framework integrates morphology-aware convolutional models with self-supervised transformers, enhanced by adaptive compression to address heterogeneous biosignals.
- ResNet1D Architecture. The ResNet1D processes ECG signals  $\mathbf{X} \in \mathbb{R}^{C \times L}$ , where C is the number of leads and L is the sequence length. The initial convolution maps  $\mathbf{X}$  to a feature space:  $\mathbf{H}_0 = \mathrm{BN}(\mathrm{Conv1d}(\mathbf{X}; W_0))$ . Each residual block applies:

$$\mathbf{Z}_{k} = \sigma(\mathrm{BN}(\mathrm{Conv1d}(\mathbf{H}_{k-1}; W_{k,1}))), \quad \mathbf{Z}_{k}' = \sigma(\mathrm{BN}(\mathrm{Conv1d}(\mathbf{Z}_{k}; W_{k,2}))),$$
(1)

- with output  $\mathbf{H}_k = \mathbf{Z}_k' + \mathcal{S}(\mathbf{H}_{k-1})$ . The final output is  $\hat{\mathbf{y}} = \operatorname{Softmax}(W_f \operatorname{vec}(\mathbf{H}_K) + b_f)$ . ResNet1D captures local ECG morphology (e.g., QRS complexes), complementing the global temporal modeling of foundation models.
- ECG-HuBERT Architecture. The HuBERT-ECG model Coppola et al. (2024), pretrained on large-scale unlabeled ECG data, extracts Mel-spectrogram features:  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T], \quad \mathbf{x}_t \in \mathbb{R}^F$ . Clustering assigns pseudo-labels  $c_t = \arg\min_k \|\mathbf{x}_t \mu_k\|_2^2$ . Masked frames are embedded by a convolutional encoder  $f_{\text{conv}}$ , producing  $\mathbf{z}_t$ , contextualized by a Transformer:  $\mathbf{h}_t = \mathcal{T}(\mathbf{z}_t + \mathbf{p}_t)$ . The loss is:

$$\mathcal{L}_{\text{HuBERT}} = -\frac{1}{\sum_{t} m_{t}} \sum_{t=1}^{T} m_{t} \log p_{\theta}(c_{t}|\mathbf{h}_{t}). \tag{2}$$

- This self-supervised pretraining enables robust, generalizable representations for ECG tasks, embodying biosignal foundation model principles.
- Adaptive Model Compression. The layer-wise adaptive compression capitalizes on the heterogeneous sensitivity of network layers: more critical layers are pruned and quantized conservatively, whereas less important layers undergo more aggressive compression Shinde (2024).
- Let a neural network  $\mathcal{M}$  have L layers, each with a weight tensor  $W_l$ . The goal is to determine the bit-width  $b_l$  and pruning factor  $p_l$  for each layer, minimizing model size while ensuring minimal accuracy loss:  $\min_{\{b_l,p_l\}_{l=1}^L} \operatorname{Size}(\mathcal{M})$  s.t. Accuracy $(\mathcal{M}_{\operatorname{comp}}) \geq A_0 \Delta$ , where  $\mathcal{M}_{\operatorname{comp}}$  is the compressed model,  $A_0$  is the baseline accuracy, and  $\Delta$  is the allowable accuracy degradation.
- Layer Importance Estimation. To guide the compression process, we assign an importance score to each layer based on two factors: Parameter Density Index (PDI) reflects the proportion of parameters in layer l relative to the total parameters:  $PDI_l = \dim(W_l) / \sum_{k=1}^L \dim(W_k)$ . Parameter

Variability Index (PVI) captures the variability of weights within a layer, which influences its sensitivity to quantization. Specifically, the PVI  $(PVI_l)$  is computed based on the variance of the weights in each layer, normalized relative to the maximum variance across all layers. This helps in assessing how much a layer's weight distribution varies, affecting its ability to maintain accuracy after quantization.

82 Combined Importance. The final importance score for layer l is a weighted sum of the PDI and the 83 PVI:  $I_l = \alpha \cdot PDI_l + \beta \cdot PVI_l$ , where  $\alpha$  and  $\beta$  are hyperparameters controlling the emphasis on 84 parameter density and sensitivity to quantization.

**Quantization.** Quantization reduces the model's memory and computational requirements by converting continuous weights to discrete levels. In *Fixed Quantization*, all weights are uniformly quantized to a global bit-width  $b_{\text{fixed}}$ :  $\hat{w} = \text{Quantize}(w, b_{\text{fixed}})$ . Layer-wise Adaptive Quantization (LAQ). Bit-widths are assigned to layers based on importance scores. For layer l, the optimal bit-width  $b_l^*$  is selected by:

$$b_l^* = \min\{b_l : \text{Accuracy}(\mathcal{M}_{\text{quant}}) \ge A_0 - \gamma \cdot I_l\},\tag{3}$$

where  $\gamma$  is a global accuracy tolerance, and  $b_l$  is selected greedily to minimize accuracy degradation.

Pruning. Pruning removes low-magnitude weights, inducing sparsity and reducing model size. In Fixed Pruning, a uniform pruning factor  $p_{\rm fixed}$  is applied across the network. Layer-wise Adaptive Pruning (LAP). Pruning is performed based on layer-specific importance, with a pruning factor  $p_l$  for each layer. For each weight:

$$\hat{W}_{l,ij} = \{ W_{l,ij}, \text{ if } |W_{l,ij}| > \theta_l, \quad 0, \text{ otherwise.}$$
 (4)

The pruning threshold  $\theta_l$  is computed as the  $p_l$ -th percentile of  $|W_l|$ , where layers with higher importance are pruned conservatively. This adaptive, importance-guided compression achieves a trade-off between model size and performance, enabling efficient deployment.

# 98 3 Experimental Setup

85

86

87

89

112

All experiments were conducted on the Kaggle platform equipped with an NVIDIA Tesla P100 GPU, leveraging PyTorch for deep learning operations.

Datasets: CPSC 2018 Challenge Dataset. Contains 6,877 twelve-lead ECG recordings (6–60 seconds, 500 Hz) with nine rhythm categories Liu et al. (2018). Chapman Clinical Dataset. Includes ~10,000 subjects with 10-second, twelve-lead ECGs (500 Hz), aggregated into four rhythm groups Zheng et al. (2020); Murat et al. (2021). CPSC and Chapman enable rigorous evaluation of benchmark performance and cross-domain generalization under class imbalance and noise.

Model Architectures and Training. We evaluate: (1) ResNet1D, capturing local ECG morphology, and (2) HuBERT-ECG, a self-supervised transformer for global temporal dependencies. This pairing probes the trade-off between interpretable CNNs and scalable foundation models. Inputs are zero-padded to 5,120 samples. Training uses Adam ( $lr=10^{-3}$  for ResNet1D,  $10^{-4}$  for HuBERT, following Kiranyaz et al. (2015)), weight decay  $10^{-3}$ , dropout (0.2–0.3), and ReduceLROnPlateau. Categorical cross-entropy loss ensures robust rhythm classification.

## 4 Results and Discussion

We conduct an ablation study to evaluate the effects of quantization and pruning on the proposed ResNet1D and HuBERT-ECG models. Table 1 reports classification metrics (Accuracy, Precision, Recall, F1) and compression ratio (CR) across various settings.

Fixed Quantization. Uniform quantization (8-bit to 1-bit) reveals distinct sensitivities. For ResNet1D,
4-bit quantization achieves the highest accuracy (0.9688) and F1 (0.9657) with 7.95 × compression,
likely due to quantization noise acting as implicit regularization. Analysis of weight distributions
shows reduced variance in feature activations, mitigating overfitting on ECG waveforms. Performance
collapses at ≤3 bits due to excessive information loss.

Table 1: Comparison of performance and compression ratio (CR) of the Proposed ResNet1D Model and HuBERT ECG Model under different quantization and pruning settings. Best, second-best, and third-best values per column are highlighted.

| Method            | Proposed ResNet1D Model |        |        |        |        | HuBERT ECG Model |        |        |        |        |
|-------------------|-------------------------|--------|--------|--------|--------|------------------|--------|--------|--------|--------|
|                   | Acc                     | Prec   | Rec    | F1     | CR     | Acc              | Prec   | Rec    | F1     | CR     |
| Full-Precision    | 0.9660                  | 0.9632 | 0.9618 | 0.9624 | 1.00x  | 0.9712           | 0.9685 | 0.9686 | 0.9685 | 1.00x  |
| Quantized (8-bit) | 0.9660                  | 0.9630 | 0.9618 | 0.9623 | 3.99x  | 0.9707           | 0.9678 | 0.9680 | 0.9678 | 3.98x  |
| Quantized (7-bit) | 0.9665                  | 0.9635 | 0.9623 | 0.9629 | 4.56x  | 0.9703           | 0.9670 | 0.9677 | 0.9673 | 4.55x  |
| Quantized (6-bit) | 0.9655                  | 0.9624 | 0.9612 | 0.9617 | 5.31x  | 0.9698           | 0.9668 | 0.9669 | 0.9668 | 5.30x  |
| Quantized (5-bit) | 0.9669                  | 0.9641 | 0.9629 | 0.9634 | 6.37x  | 0.9693           | 0.9662 | 0.9662 | 0.9662 | 6.35x  |
| Quantized (4-bit) | 0.9688                  | 0.9665 | 0.9651 | 0.9657 | 7.95x  | 0.9566           | 0.9523 | 0.9535 | 0.9524 | 7.91x  |
| Quantized (3-bit) | 0.9410                  | 0.9358 | 0.9344 | 0.9341 | 10.58x | 0.5038           | 0.6033 | 0.5508 | 0.4917 | 10.50x |
| Quantized (2-bit) | 0.3555                  | 0.2672 | 0.4175 | 0.2967 | 15.79x | 0.2101           | 0.0525 | 0.2500 | 0.0868 | 15.63x |
| Quantized (1-bit) | 0.3551                  | 0.0888 | 0.2500 | 0.1310 | 31.16x | 0.2101           | 0.0525 | 0.2500 | 0.0868 | 30.49x |
| Pruned (10%)      | 0.9504                  | 0.9478 | 0.9450 | 0.9456 | 1.11x  | 0.9660           | 0.9627 | 0.9619 | 0.9620 | 1.11x  |
| Pruned (30%)      | 0.4396                  | 0.5505 | 0.4879 | 0.3964 | 1.43x  | 0.9363           | 0.9323 | 0.9274 | 0.9288 | 1.43x  |
| Pruned (50%)      | 0.2158                  | 0.0540 | 0.2500 | 0.0888 | 2.00x  | 0.3617           | 0.5014 | 0.3959 | 0.3130 | 2.00x  |
| Pruned (70%)      | 0.2101                  | 0.0525 | 0.2500 | 0.0868 | 3.33x  | 0.2342           | 0.1094 | 0.2704 | 0.1497 | 3.32x  |
| Pruned (90%)      | 0.3551                  | 0.0888 | 0.2500 | 0.1310 | 9.92x  | 0.2158           | 0.0539 | 0.2500 | 0.0887 | 9.86x  |
| Proposed LAQ      | 0.9660                  | 0.9626 | 0.9624 | 0.9625 | 10.44x | 0.9703           | 0.9675 | 0.9667 | 0.9670 | 9.43x  |
| Proposed LAP      | 0.9665                  | 0.9639 | 0.9626 | 0.9631 | 1.67x  | 0.9717           | 0.9690 | 0.9677 | 0.9683 | 2.41x  |

In contrast, HuBERT-ECG is more sensitive to low-precision quantization, as its self-attention layers require fine-grained weight resolution to capture global temporal dependencies Vaswani et al. (2017). While 8–7 bits maintain accuracy  $\approx$ 0.97, performance drops sharply at 4-bit (0.9566), unlike the robust ResNet1D.

Layer-wise Adaptive Quantization (LAQ). Our LAQ strategy achieves near-baseline accuracy with high compression. By allocating precision based on layer importance, LAQ preserves critical layers (e.g., convolutional filters capturing QRS complexes) while aggressively compressing redundant ones, optimizing for noisy biosignals. ResNet1D reaches 0.9660 accuracy at 10.44× CR, while HuBERT-ECG attains 0.9703 at 9.43×, consistently outperforming fixed quantization schemes.

Fixed Pruning. Light pruning (10%) has minimal impact, but pruning beyond 30% rapidly degrades accuracy by disrupting structural connectivity essential for ECG pattern recognition. ResNet1D accuracy falls to 0.4396 at 30% pruning, and HuBERT-ECG shows similar degradation, confirming that indiscriminate weight removal is detrimental.

Layer-wise Adaptive Pruning (LAP). LAP offers stable trade-offs. For ResNet1D, LAP retains 0.9665 accuracy at 1.67× CR. HuBERT-ECG achieves 0.9717 accuracy with 2.41× CR, surpassing fixed pruning. By prioritizing layers with high importance (e.g., self-attention heads in HuBERT), LAP preserves expressive capacity, crucial for generalizable biosignal modeling.

Discussion. Adaptive, layer-aware compression (LAQ/LAP) achieves Pareto-optimal trade-offs between accuracy and efficiency, enabling real-time ECG monitoring on edge devices. ResNet1D's robustness to quantization and pruning makes it ideal for lightweight applications, while HuBERT-ECG benefits from adaptive strategies to preserve self-supervised features. The framework's ability to generalize across datasets and handle noisy biosignals aligns with the scalability and robustness goals of foundation models, advancing clinical deployment of AI-driven health monitoring.

# 144 5 Conclusion

We present an adaptive compression framework for biosignal foundation models, enabling efficient ECG monitoring on edge devices with up to 10.44× compression without any loss. Layer importance guides conservative compression of critical layers and aggressive optimization of redundant ones. The framework's architecture-agnostic design generalizes across datasets and modalities, supporting real-time health monitoring. Future work will explore multimodal biosignal integration (e.g., EEG, EMG), dynamic inference, and ethical considerations for clinical adoption, enhancing the framework's impact on scalable, reliable biosignal AI.

Table 2: Classification performance comparisons on Chapman and CPSC 2018 datasets. Best, second-best, and third-best values per column are highlighted.

| Dataset   | Author                    | Classes | #Lead | Method                       | Acc.  | Prec. | Rec.  | F1    | CR           |
|-----------|---------------------------|---------|-------|------------------------------|-------|-------|-------|-------|--------------|
| Chapman   | Yildirim et al. (2020)    | 4       | 12    | Deep neural network          | 96.13 | 95.78 | 95.43 | 95.57 | _            |
|           | Baygin et al. (2021)      | 4       | 1     | HIT pattern SVM              | 97.18 | 97.07 | 96.77 | 96.91 | _            |
|           | Murat et al. (2021)       | 4       | 1     | DNN + feature fusion         | 98.00 | 97.76 | 97.70 | 97.72 | _            |
|           | Domazetoski et al. (2022) | 3       | 12    | XGBoost                      | 89.37 | -     | -     | -     | _            |
|           | Venkatesh et al. (2024)   | 5       | 1     | 1D-CNN-BiLSTM                | 93.97 | 93.96 | 98.49 | 93.95 | _            |
|           | ResNet1D + LAQ            | 4       | 1     | Residual Network             | 96.60 | 96.26 | 96.24 | 96.25 | 10.44x       |
|           | HuBERT ECG + LAQ          | 4       | 1     | Foundational Network         | 97.03 | 96.75 | 96.67 | 96.70 | <u>9.43x</u> |
| CPSC 2018 | Zhang et al. (2020)       | 9       | 12    | CNN+Attention+BiGRU          | 86.83 | 84.18 | 82.93 | 83.51 | _            |
|           | Ge et al. (2021)          | 9       | 1     | SEBlock(CNN)                 | _     | 83.00 | 82.70 | 82.80 | _            |
|           | Liu et al. (2022a)        | 9       | 12    | CRT-Net                      | 87.20 | 87.30 | 87.20 | 86.90 | _            |
|           | Li and Zhang (2023)       | 9       | 12    | KNN+CNN                      | 88.50 | 87.77 | 87.08 | 87.37 | _            |
|           | Dhyani et al. (2023)      | 9       | 12    | ResNet+RNN                   | 93.29 | 93.38 | 93.10 | 93.09 | -            |
|           | Ji et al. (2024)          | 9       | 12    | Multi-scale grid transformer | 87.34 | 85.67 | 86.21 | 85.90 | _            |
|           | Proposed ResNet1D         | 9       | 1     | Residual Network             | 95.78 | 95.61 | 95.81 | 95.68 | 10.44x       |

# A Comparison with Existing Work.

On the Chapman dataset, ResNet1D+LAQ achieves 0.9660 accuracy with 10.44× compression, and 153 HuBERT-ECG+LAQ reaches 0.9703 with 9.43× (see Table 2). Unlike prior methods optimizing 154 solely for accuracy Murat et al. (2021), our approach sets a new state-of-the-art by balancing clinical 155 fidelity and edge deployability. On CPSC 2018, ResNet1D achieves 95.78% accuracy with 10× 156 compression, outperforming baselines Dhyani et al. (2023). These results provide the first evidence 157 of compact ECG models achieving superior performance while enabling real-time deployment on 158 resource-constrained devices. To assess generalization, a key property of biosignal foundation models, 159 we compare performance across CPSC and Chapman datasets. HuBERT-ECG's self-supervised 160 pretraining enhances robustness to Chapman's class imbalance, achieving 0.9703 accuracy despite 161 fewer training samples. ResNet1D excels on CPSC (95.78%) due to its focus on local morphology, 162 but shows slightly lower generalization on Chapman's heterogeneous clinical data. These findings 163 underscore the complementary strengths of convolutional and transformer-based foundation models 164 for biosignals. 165

## 166 References

152

- L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, Y. Duan, O. Al-Shamma, J. Santamaría, M. A. Fadhel,
   M. Al-Amidie, and L. Farhan. Review of deep learning: concepts, cnn architectures, challenges, applications,
   future directions. *Journal of big Data*, 8:1–74, 2021.
- M. Baygin, T. Tuncer, S. Dogan, R.-S. Tan, and U. R. Acharya. Automated arrhythmia detection with homeomorphically irreducible tree technique using more than 10,000 individual subject ecg records. *Information Sciences*, 575:323–337, 2021.
- 173 J. Chen and X. Ran. Deep learning with edge computing: A review. *Proceedings of the IEEE*, 107(8):1655–1674, 2019.
- 175 E. Coppola, M. Savardi, M. Massussi, M. Adamo, M. Metra, and A. Signoroni. Hubert-ecg as a self-supervised foundation model for broad and scalable cardiac applications. *medRxiv*, pages 2024–11, 2024.
- M. Courbariaux, Y. Bengio, and J.-P. David. Binaryconnect: Training deep neural networks with binary weights
   during propagations. Advances in neural information processing systems, 28, 2015.
- J. R. G. de Santana, M. G. F. Costa, and C. F. F. Costa Filho. A new approach to classify cardiac arrythmias using
   2d convolutional neural networks. In 2021 43rd Annual International Conference of the IEEE Engineering in
   Medicine & Biology Society (EMBC), pages 566–570, 2021. doi: 10.1109/EMBC46164.2021.9630938.
- S. Dhyani, A. Kumar, and S. Choudhury. Arrhythmia disease classification utilizing resrnn. *Biomedical Signal Processing and Control*, 79:104160, 2023.
- V. Domazetoski, G. Gligoric, M. Marinkovic, A. Shvilkin, J. Krsic, L. Kocarev, and M. D. Ivanovic. The
   influence of atrial flutter in automated detection of atrial arrhythmias-are we ready to go into clinical
   practice?". Computer Methods and Programs in Biomedicine, 221:106901, 2022.

- Z. Ebrahimi, M. Loni, M. Daneshtalab, and A. Gharehbaghi. A review on deep learning methods for ecg 187 arrhythmia classification. Expert Systems with Applications: X, 7:100033, 2020. 188
- 189 H. El-Ghaish and E. Eldele. Ecgtransform: Empowering adaptive ecg arrhythmia classification framework with
- bidirectional transformer. Biomedical Signal Processing and Control, 89:105714, 2024. ISSN 1746-8094. 190
- doi: https://doi.org/10.1016/j.bspc.2023.105714. URL https://www.sciencedirect.com/science/ 191 article/pii/S1746809423011473. 192
- J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. arXiv preprint 193 arXiv:1803.03635, 2018. 194
- R. Ge, T. Shen, Y. Zhou, C. Liu, L. Zhang, B. Yang, Y. Yan, J.-L. Coatrieux, and Y. Chen. Convolutional 195 squeeze-and-excitation network for ecg arrhythmia detection. Artificial Intelligence in Medicine, 121:102181, 196 197
- A. Y. Hannun, P. Raipurkar, M. Haghpanahi, G. H. Tison, C. Bourn, M. P. Turakhia, and A. Y. Ng. Cardiologist-198 level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network. 199 Nature medicine, 25(1):65-69, 2019. 200
- K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners. In 201 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 16000–16009, 202 203
- G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. arXiv preprint 204 arXiv:1503.02531, 2015. 205
- W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed. Hubert: Self-supervised 206 speech representation learning by masked prediction of hidden units. IEEE/ACM transactions on audio, 207 speech, and language processing, 29:3451-3460, 2021. 208
- I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Quantized neural networks: Training neural 209 networks with low precision weights and activations. journal of machine learning research, 18(187):1–30, 210 211 2018.
- C. Ji, L. Wang, J. Qin, L. Liu, Y. Han, and Z. Wang. Msgformer: A multi-scale grid transformer network for 212 12-lead ecg arrhythmia detection. Biomedical Signal Processing and Control, 87:105499, 2024. 213
- Y. Jin, J. Liu, Y. Liu, C. Qin, Z. Li, D. Xiao, L. Zhao, and C. Liu. A novel interpretable method based on 214 dual-level attentional deep neural network for actual multilabel arrhythmia detection. IEEE Transactions on 215 Instrumentation and Measurement, 71:1–11, 2021. 216
- S. Kiranyaz, T. Ince, and M. Gabbouj. Real-time patient-specific ecg classification by 1-d convolutional neural 217 networks. IEEE transactions on biomedical engineering, 63(3):664–675, 2015. 218
- R. Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. arXiv preprint 219 arXiv:1806.08342, 2018. 220
- J. Li, A. Aguirre, J. Moura, C. Liu, L. Zhong, C. Sun, G. Clifford, B. Westover, and S. Hong. An electrocardio-221 gram foundation model built on over 10 million recordings with external evaluation across multiple domains. 222 223 arXiv preprint arXiv:2410.04133, 2024.
- Z. Li and H. Zhang. Fusing deep metric learning with knn for 12-lead multi-labelled ecg classification. 224 Biomedical Signal Processing and Control, 85:104849, 2023. 225
- F. Liu, C. Liu, L. Zhao, X. Zhang, X. Wu, X. Xu, Y. Liu, C. Ma, S. Wei, Z. He, et al. An open access database 226 for evaluating the algorithms of electrocardiogram rhythm and morphology abnormality detection. Journal of 227 Medical Imaging and Health Informatics, 8(7):1368–1373, 2018. 228
- J. Liu, Z. Li, X. Fan, X. Hu, J. Yan, B. Li, Q. Xia, J. Zhu, and Y. Wu. Crt-net: A generalized and scalable 229 230 framework for the computer-aided diagnosis of electrocardiogram signals. Applied Soft Computing, 128: 109481, 2022a. 231
- Z. Liu, K.-T. Cheng, D. Huang, E. P. Xing, and Z. Shen. Nonuniform-to-uniform quantization: Towards accurate 232 quantization via generalized straight-through estimation. In Proceedings of the IEEE/CVF conference on 233 computer vision and pattern recognition, pages 4942-4952, 2022b. 234
- 235 J. Ma, Y. He, F. Li, L. Han, C. You, and B. Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024. 236

- F. Murat, O. Yildirim, M. Talo, Y. Demir, R.-S. Tan, E. J. Ciaccio, and U. R. Acharya. Exploring deep features and ecg attributes to detect cardiac rhythm classes. *Knowledge-Based Systems*, 232:107473, 2021.
- A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- T. Shinde. Adaptive quantization and pruning of deep neural networks via layer importance estimation. In
   Workshop on Machine Learning and Compression, NeurIPS 2024, 2024.
- N. T. Srinivasan and R. J. Schilling. Sudden cardiac death and arrhythmias. *Arrhythmia & electrophysiology* review, 7(2):111, 2018.
- J. H. Tan, Y. Hagiwara, W. Pang, I. Lim, S. L. Oh, M. Adam, R. San Tan, M. Chen, and U. R. Acharya.

  Application of stacked convolutional and long short-term memory network for accurate identification of cad ecg signals. *Computers in biology and medicine*, 94:19–26, 2018.
- M. Tan. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint* arXiv:1905.11946, pages 6105–6114, 2019.
- E. Tiu, E. Talius, P. Patel, C. P. Langlotz, A. Y. Ng, and P. Rajpurkar. Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. *Nature biomedical engineering*, 6(12): 1399–1406, 2022.
- 253 R. Trobec, I. Tomašić, A. Rashkovska, M. Depolli, and V. Avbelj. *Body sensors and electrocardiography*.
  254 Springer, 2018.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- N. P. Venkatesh, R. P. Kumar, B. C. Neelapu, K. Pal, and J. Sivaraman. Automated atrial arrhythmia classification using 1d-cnn-bilstm: A deep network ensemble model. *Biomedical Signal Processing and Control*, 97:106703, 259 2024.
- O. Yildirim, M. Talo, E. J. Ciaccio, R. San Tan, and U. R. Acharya. Accurate deep neural network model to detect cardiac arrhythmia on more than 10,000 individual subject ecg records. *Computer methods and programs in biomedicine*, 197:105740, 2020.
- J. Zhang, A. Liu, M. Gao, X. Chen, X. Zhang, and X. Chen. Ecg-based multi-class arrhythmia detection using
   spatio-temporal attention-based convolutional recurrent neural network. *Artificial Intelligence in Medicine*,
   106:101856, 2020.
- J. Zheng, J. Zhang, S. Danioko, H. Yao, H. Guo, and C. Rakovski. A 12-lead electrocardiogram database for arrhythmia research covering more than 10,000 patients. *Scientific data*, 7(1):48, 2020.