

ETP: LEARNING TRANSFERABLE ECG REPRESENTATIONS VIA ECG-TEXT PRE-TRAINING

Che Liu^{1,*}, Zhongwei Wan^{2,*}, Sibò Cheng¹, Mi Zhang², Rossella Arcucci¹

¹ Imperial College London, UK

² The Ohio State University, USA

{che.liu21, sibò.cheng, r.arcucci}@imperial.ac.uk, {wan.512, mizhang.1}@osu.edu

ABSTRACT

In the domain of cardiovascular healthcare, the Electrocardiogram (ECG) serves as a critical, non-invasive diagnostic tool. Although recent strides in self-supervised learning (SSL) have been promising for ECG representation learning, these techniques often require annotated samples and struggle with classes not present in the fine-tuning stages. To address these limitations, we introduce ECG-Text Pre-training (ETP), an innovative framework designed to learn cross-modal representations that link ECG signals with textual reports. For the first time, this framework leverages the zero-shot classification task in the ECG domain. ETP employs an ECG encoder along with a pre-trained language model to align ECG signals with their corresponding textual reports. The proposed framework excels in both linear evaluation and zero-shot classification tasks, as demonstrated on the PTB-XL and CPSC2018 datasets, showcasing its ability for robust and generalizable cross-modal ECG feature learning.

Index Terms— Electrocardiogram, ECG-Text Pre-training, Self-supervised Learning

1. INTRODUCTION

The Electrocardiogram (ECG) is a crucial clinical diagnostic tool for various cardiac conditions. While deep learning has shown promise in ECG classification, its effectiveness often depends on the availability of high-quality labels and expert review, making the process labor-intensive and costly.

In the quest to circumvent the pitfalls of extensive annotation, self-supervised learning (SSL) has emerged as a promising avenue, excelling particularly with datasets harboring limited annotations [1, 2, 3]. SSL paves the way for harnessing ECG representations beneficial for a spectrum of downstream tasks like abnormality detection and arrhythmia classification [4, 5]. However, a significant bottleneck remains: extant ECG SSL [6, 7, 8, 9] strategies still lean heavily on substantial annotated data for fine-tuning on downstream applications as shown in Fig 2b. Such a dependency becomes

particularly limiting for rare cardiac conditions, steering research attention to the zero-shot classification. This paradigm aims to negate the need for annotated samples of unseen categories by leveraging cross-modal representation from ECG and disease-related textual prompt and utilize the ECG-text similarity to determine the predicted disease that do not need annotated data in downstream tasks as depicted in Fig 2a.

The path to zero-shot learning for ECG isn't devoid of obstacles. Primarily, there exists a semantic disjunction between the continuous numerical nature of ECG and the discrete clinical terminologies in textual reports [10, 11, 12]. Further complications arise from domain adaptation issues and scalability concerns, with zero-shot models often requiring considerable computational resources [13]. While recent studies, such as those by [14] and [15], have made headway in ECG zero-shot classification, they remain tethered to supervised learning during pre-training, demanding extensive annotated ECG data.

Witnessing the potential of vision-language pre-training in broader contexts, as evidenced by works like CLIP [16], we introduce **ECG-Text Pre-training (ETP)**. This innovative approach seeks to leverage the 12-leads ECG and its corresponding textual reports within a cross-modal learning paradigm. ETP features a language model paired with an ECG encoder to yield text and ECG embeddings. Leveraging a priori clinical knowledge, the text is channeled through a sizeable frozen language model with 1D CNN serving as the ECG encoder's backbone. Both components possess linear projection heads, ensuring the harmonization of text and ECG dimensions. Following this, the concordance between ECG and text embeddings becomes the focal point to minimize contrastive learning loss and yield classification probabilities for diverse ECG categories.

The key contributions from our research are outlined as follows:

- We are the pioneers in delving into and unveiling the potential of ECG-Text Pre-training (ETP) specifically for ECG signals.
- Our approach not only achieves state-of-the-art (SOTA) results in the fine-tuning phase but also becomes the

* Equal Contribution.

first to demonstrate the viability of zero-shot tasks. Furthermore, compared to uni-modal SSL, our method exhibits enhanced robustness and transferability.

- We have established the comprehensive benchmark for ETP, focusing on the confluence of ECG-Text pre-training and ECG signals.

2. METHODOLOGY

2.1. ECG-Text Pre-training

Incorporating both ECG signals and paired textual description, we employ the following modifications based on the CLIP framework:

Given the CLIP framework as a reference [17], we integrate a contrastive learning aiming to predict the associated pair $(e_{ecg,i}, t_{ecg,i})$ among the $N \times N$ probable ECG-text combinations, while strategically positioning the remaining $N^2 - N$ negative combinations at a distance. In this context, two distinct encoders for ECG signals and text, denoted as \mathcal{F}_{ecg} and \mathcal{F}_{text} respectively, transform $e_{ecg,i}$ and $t_{ecg,i}$ into a latent embedding space, represented as $[\hat{e}]_{i=1}^N$. Subsequently, two separate non-linear projectors for ECG signals and text, denoted as \mathcal{P}_{ecg} and \mathcal{P}_{text} respectively, convert $e_{ecg,i}$ and $t_{ecg,i}$ into a consistent dimension, termed d . This process can be represented as:

$$\hat{e}_{ecg,i} = \mathcal{P}_{ecg}(\mathcal{F}_{ecg}(e_{ecg,i})), \quad (1)$$

$$\hat{t}_{ecg,i} = \mathcal{P}_t(\mathcal{F}_{text}(t_{ecg,i})), \quad (2)$$

with both $\hat{e}_{ecg,i}$ and $\hat{t}_{ecg,i}$ belonging to the set \mathbb{R}^d .

From the training set, we extract ECG feature vectors denoted by $[\hat{e}_{ecg,i}]_{i=1}^N$ and text feature vectors represented by $[\hat{t}_{ecg,i}]_{i=1}^N$. Following this, we then calculate the cosine similarities as $r_{i,i}^{e2t} = \hat{e}_{ecg,i}^\top \hat{t}_{ecg,i}$ and $r_{i,i}^{t2e} = \hat{t}_{ecg,i}^\top \hat{e}_{ecg,i}$, which illustrate the ECG-text and text-ECG compatibilities respectively. The loss function, \mathcal{L}_{CE} , is then expressed as:

$$\mathcal{L}_e^{e2t} = -\log \frac{\exp(r_{i,i}^{e2t}/\sigma_1)}{\sum_{j=1}^K \exp(r_{i,j}^{e2t}/\sigma_1)}, \quad (3)$$

$$\mathcal{L}_i^{t2e} = -\log \frac{\exp(r_{i,i}^{t2e}/\sigma_1)}{\sum_{j=1}^K \exp(r_{i,j}^{t2e}/\sigma_1)} \quad (4)$$

$$\mathcal{L}_{total} = \frac{1}{2K} \sum_{i=1}^N (\mathcal{L}_e^{e2t} + \mathcal{L}_t^{t2e}), \quad (5)$$

Here, \mathcal{L}_e^{e2t} and \mathcal{L}_t^{t2e} are ECG-text and text-ECG cross-modal contrastive losses respectively. σ_1 denotes the temperature hyper-parameter, which in our research was fixed at 0.07. Meanwhile, K symbolizes the batch size per step, with K being a subset of N . Through the total loss, \mathcal{L}_{total} , our

model gets trained to maximize mutual information between the aligned ECG-text pairs that encompass cross-referential attributes in a batch.

2.2. Self-supervised Contrastive Learning

Conventional contrastive-based SSL methods [18, 19, 20, 7, 6] rely on strong augmentation to generate two distinct views for the input data, such as random segmentation and inversion, to the original ECG signals. This creates augmented views that serve as positive pairs $[(e_{ecg,i}, e'_{ecg,i})]_{i=1}^N$, with the other ECG signals in the mini-batch being considered as negative examples. The pipeline is depicted in Fig 1b. This data augmentation approach aligns with the strategy outlined in [7].

Next, we derive the representations of these augmented views, represented as $[\hat{e}']_{i=1}^N$, using the ECG projector p_{ecg} and ECG encoder \mathcal{F}_{ecg} . This is analogous to obtaining the representations $[\hat{e}]_{i=1}^N$. Consequently, our ECG invariant learning goal is defined as:

$$\mathcal{L}_{SSL} = -\frac{1}{K} \sum_{j=1}^N \log \frac{\exp(r_{i,i}^{e2e'}/\sigma_2)}{\sum_{j=1}^N \exp(r_{i,j}^{e2e'}/\sigma_2)} \quad (6)$$

$$\hat{e}_{ecg,i} = \mathcal{F}_{ecg}(e_{ecg,i}), \hat{e}'_{ecg,i} = \mathcal{F}_{ecg}(e'_{ecg,i}) \quad (7)$$

$$r_{i,i}^{e2e'} = \hat{e}_{ecg,i}^\top \hat{e}'_{ecg,i} \quad (8)$$

In Eq 6, the temperature hyper-parameter σ_2 retains its value of 0.07 when considering the overall loss objective \mathcal{L}_{SSL} .

3. EXPERIMENTS AND ANALYSIS

3.1. Datasets

PTB-XL The ECG dataset under examination is substantial, encompassing 21,837 ECG signals that were accumulated from 18,885 patients during the period of October 1989 to June 1996. The collected data consists of 12-lead ECG, each sampled at a rate of 500 Hz with a duration of 10 seconds, where each ECG signal is paired with the corresponding ECG reports. The reports are generated by the standard protocol and only describe the ECG without final diagnosis. The original ECG reports were written in 70.89% German, 27.9% English, and 1.21% Swedish, and were converted into structured SCP-ECG statements. For downstream tasks, we follow the official split from [6] to build the train/val/test split only with single category. Furthermore, each record in downstream task setting is classified under one of five primary diagnostic categories: Normal (NORM), Myocardial Infarction (MI), ST/T Change (STTC), Conduction Disturbance (CD), and Hypertrophy (HYP).

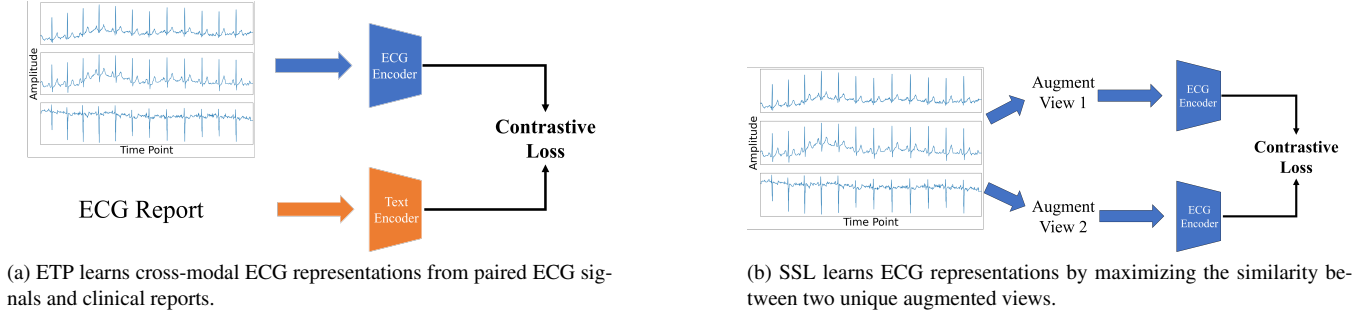


Fig. 1: Comparison between ETP and SSL.

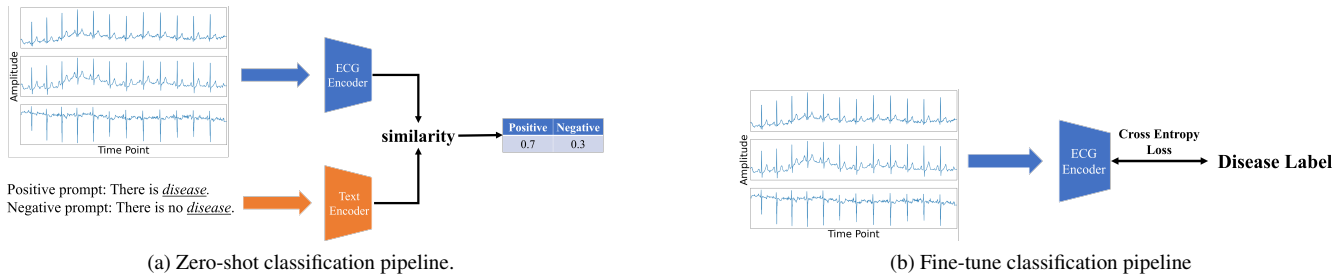


Fig. 2: The pipeline of zero-shot classification and fine-tune classification.

CPSC2018 This dataset, which is publicly accessible, comprises 6,877 standard 12-lead ECG records, each sampled at a rate of 500 Hz, and the duration of these records ranges from 6 to 60 seconds. The dataset is annotated with nine distinct labels, which include Atrial fibrillation (AF), First-degree atrioventricular block (I-AVB), Left bundle branch block (LBBB), Right bundle branch block (RBBB), Premature atrial contraction (PAC), Premature ventricular contraction (PVC), ST-segment depression (STD), ST-segment elevation (STE), and normal (Normal).

For both datasets, we adhere to the official split as outlined in [6] and only select samples that belong to a single category.

3.2. Implementation

The ECG encoder we utilize is ResNet18-1D. This is adapted from its two-dimensional version, ResNet18-2D [21], by transitioning to 1D convolutional layers. For text encoding, we employ BioClinicalBERT [22], pre-trained on clinical notes and bio-clinical articles. Our model integrates two linear projection heads: one for the ECG encoder and another for the text encoder. Both are characterized by an output dimension of 512 and utilize a temperature parameter τ initialized to 0.07. The ECG encoder’s optimization is handled using the Adam optimizer, set with a learning rate and weight decay of $2 \times e^{-3}$ and $1 \times e^{-5}$. During pre-training, we operate over 50 epochs with a batch size of 128, while all subsequent downstream tasks are processed with a batch size of 32. All experimental procedures are executed using PyTorch 2.0 on 1 NVIDIA A100-40GB GPU.

3.3. Results on Linear Evaluation

In the task of linear evaluation, we rigorously test the quality and robustness of the ECG representations generated by our ETP framework. To do this, we keep the pre-trained ECG encoder fixed and only update a linear classifier that is initialized randomly. This evaluation methodology is applied to two large-scale public ECG datasets with disease-level annotation, PTB-XL and CPSC2018, using Area Under the Curve (AUC) score and F1-score as the primary metrics for performance assessment. As clearly indicated in Table 1, the ETP framework sets new performance standards, outclassing all existing baseline methods. Specifically, it achieves an AUC of 83.5 and an F1-score of 61.3 on the PTB-XL dataset. Similarly, on the CPSC2018 dataset, ETP registers an AUC of 86.1 and an F1-score of 63.4. These compelling results not only validate the effectiveness of ETP but also firmly establish it as the leading methodology for learning ECG representations.

3.4. Results on Zero-shot Classification

To delve deeper into the capabilities of learn cross-modal representation from proposed ETP framework, we conducted zero-shot classification tasks on both PTB-XL and CPSC2018 datasets. The results are presented in Tab 2 and 3. Our zero-shot classification pipeline is inspired by the CLIP framework [17]. We employ the phrase ‘this ECG indicates disease name’ as a positive prompt and calculate the cosine similarity between the ECG and prompt embeddings. The

Table 1: Linear evaluation results on PTB-XL and CPSC2018. Best results are in bold.

Method	PTB-XL		CPSC2018	
	AUC	F1	AUC	F1
Random init	71.5	52.3	72.1	59.9
CPC [23]	70.3	54.2	74.6	53.6
SimCLR [18]	67.5	55.5	73.2	56.8
BYOL [19]	76.1	56.8	77.4	61.3
SimSiam [20]	71.4	56.8	75.5	62.0
TS-TCC [7]	81.8	56.4	83.5	62.2
CLOCS [24]	81.7	55.8	82.0	61.3
ASTCL [6]	82.0	57.4	84.2	62.8
ETP	83.5	61.3	86.1	63.4

prompt with the highest similarity is selected as the predicted category.

PTB-XL As shown in Table 2, the ETP pre-trained model consistently surpasses models with random initialization across various metrics, including AUC, ACC, and F1-score. For example, in the ‘NORM’ category, ETP achieves an AUC of 71.8, compared to 52.7 for random initialization. It also attains an ACC of 87.4 in the ‘HYP’ category, as opposed to 10.6 with random initialization. However, it’s important to note that the AUC score for ETP is lower in the ‘MI’ category, indicating potential areas for improvement in specific disease classifications. Overall, ETP demonstrates significant enhancements, with average scores of 54.6 for AUC, 60.8 for ACC, and 33.1 for F1-score.

CPSC2018 Table 3 shows similar trends. The ETP pre-trained model consistently outperforms models initialized randomly across various metrics, such as AUC, ACC, and F1-score. Specifically, in categories like ‘LBBB,’ ETP achieves an AUC of 81.3, compared to 33.3 from a randomly initialized model. Additionally, ETP attains an ACC of 72.3 in the ‘PAC’ category, as opposed to 9.6 from a randomly initialized model. The average scores across all categories further underscore ETP’s superiority, with an average AUC of 57.1, ACC of 60.9, and F1-score of 27.1. These substantial improvements across all disease categories highlight the effectiveness of the cross-modal representation learned by ETP.

The results affirm the efficacy of ETP in learning robust cross-modal representations for ECG and paired reports. While ETP shows promising results in most categories, there are specific areas, such as the ‘MI’ category in the PTB-XL dataset, where further refinement could be beneficial. Overall, the ETP framework demonstrates a compelling advantage over random initialization in zero-shot classification tasks, thereby validating its potential for practical applications in cardiovascular healthcare.

Table 2: Zero-shot classification Results on PTB-XL. Best results are in bold.

Category	Method	PTB-XL		
		AUC	ACC	F1
NORM	Random init	52.7	58.8	72.8
	ETP	71.8	56.8	73.4
MI	Random init	57.6	54.4	28.7
	ETP	46.4	15.5	26.6
STTC	Random init	55.3	43.7	26.4
	ETP	56.3	57.8	24.8
CD	Random init	35.5	10.6	19.3
	ETP Pre-trained	52.6	87.4	28.1
HYP	Random init	25.4	3.5	6.3
	ETP	45.8	86.4	12.2
Average	Random init	45.3	34.2	30.7
	ETP	54.6	60.8	33.1

Table 3: Zero-shot classification Results on CPSC2018. Best results are in bold.

Category	Method	CPSC2018		
		AUC	ACC	F1
Normal	Random init	51.6	23.1	21.9
	ETP	55.1	60.6	26.8
AF	Random init	49.9	15.9	27.4
	ETP	50.9	52.7	32.5
I-AVB	Random init	46.3	47.4	22.7
	ETP	50.8	51.8	23.4
LBBB	Random init	33.3	4.9	6.0
	ETP	81.3	94.0	35.1
RBBB	Random init	45.3	24.3	39.3
	ETP	55.3	24.3	39.3
PAC	Random init	39.9	9.6	16.4
	ETP Pre-trained	46.3	72.3	18.8
PVC	Random init	43.8	10.8	19.4
	ETP	65.9	78.3	27.0
STD	Random init	39.5	35.1	22.7
	ETP	47.0	19.6	24.3
STE	Random init	45.0	27.0	8.1
	ETP	61.2	95.0	16.2
Average	Random init	43.8	24.2	18.2
	ETP	57.1	60.9	27.1

4. CONCLUSION

In this work, we propose ETP, the novel framework to learn cross-modal representation from unannotated ECG and associated report. We also first build the comprehensive benchmark on linear evaluation and zero-shot classification with ECG cross-modal learning and SSL. ETP surpass all SSL methods on linear evaluation task and endow the zero-shot ability to ECG community via the proposed framework and evaluated on two large-scale public datasets, PTB-XL and CPSC2018. Overall, this work establishes the first comprehensive benchmark for ECG zero-shot classification and cross-modal learning, demonstrating the capability and potential of jointly learning ECG signals and paired reports.

5. REFERENCES

- [1] Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon, "A survey on contrastive self-supervised learning," *Technologies*, vol. 9, no. 1, pp. 2, 2020.
- [2] Yen-hsiu Chou, Shenda Hong, Yuxi Zhou, Junyuan Shang, Moxian Song, and Hongyan Li, "Knowledge-shot learning: An interpretable deep model for classifying imbalanced electrocardiography data," *Neurocomputing*, vol. 417, pp. 64–73, 2020.
- [3] Zhongwei Wan, Che Liu, Mi Zhang, Jie Fu, Benyou Wang, Sibio Cheng, Lei Ma, César Quilodrán-Casas, and Rossella Arcucci, "Med-unic: Unifying cross-lingual medical vision-language pre-training by diminishing bias," *arXiv preprint arXiv:2305.19894*, 2023.
- [4] Che Liu, Sibio Cheng, Weiping Ding, and Rossella Arcucci, "Spectral cross-domain neural network with soft-adaptive threshold spectral enhancement," *arXiv preprint arXiv:2301.10171*, 2023.
- [5] Temesgen Mehari and Nils Strodthoff, "Self-supervised representation learning from 12-lead ecg data," *Computers in Biology and Medicine*, vol. 141, pp. 105114, 2022.
- [6] Ning Wang, Panpan Feng, Zhaoyang Ge, Yanjie Zhou, Bing Zhou, and Zongmin Wang, "Adversarial spatiotemporal contrastive learning for electrocardiogram signals," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [7] Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, Chee Keong Kwok, Xiaoli Li, and Cuntai Guan, "Time-series representation learning via temporal and contextual contrasting," *arXiv preprint arXiv:2106.14112*, 2021.
- [8] Yinda Chen, Wei Huang, Xiaoyu Liu, Qi Chen, and Zhiwei Xiong, "Learning multiscale consistency for self-supervised electron microscopy instance segmentation," *arXiv preprint arXiv:2308.09917*, 2023.
- [9] Yinda Chen, Wei Huang, Shenglong Zhou, Qi Chen, and Zhiwei Xiong, "Self-supervised neuron segmentation with multi-agent reinforcement learning," .
- [10] Jun Li, Che Liu, Sibio Cheng, Rossella Arcucci, and Shenda Hong, "Frozen language model helps ecg zero-shot learning," *arXiv preprint arXiv:2303.12311*, 2023.
- [11] Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency, "Foundations and recent trends in multimodal machine learning: Principles, challenges, and open questions," *arXiv preprint arXiv:2209.03430*, 2022.
- [12] Yinda Chen, Che Liu, Wei Huang, Sibio Cheng, Rossella Arcucci, and Zhiwei Xiong, "Generative text-guided 3d vision-language pretraining for unified medical image segmentation," *arXiv preprint arXiv:2306.04811*, 2023.
- [13] Che Liu, Sibio Cheng, Chen Chen, Mengyun Qiao, Weitong Zhang, Anand Shah, Wenjia Bai, and Rossella Arcucci, "M-flag: Medical vision-language pre-training with frozen language models and latent space geometry optimization," *arXiv preprint arXiv:2307.08347*, 2023.
- [14] Mehmet Yamaç, Mert Duman, İlke Adaloğlu, Serkan Kiranyaz, and Moncef Gabbouj, "A personalized zero-shot ecg arrhythmia monitoring system: From sparse representation based domain adaptation to energy efficient abnormal beat detection for practical ecg surveillance," *arXiv preprint arXiv:2207.07089*, 2022.
- [15] Sathvik Bhaskarpanidit, Priyanka Gupta, and Manik Gupta, "Lets-gzsl: A latent embedding model for time series generalized zero shot learning," *arXiv preprint arXiv:2207.12007*, 2022.
- [16] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al., "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [17] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al., "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [18] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [19] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al., "Bootstrap your own latent—a new approach to self-supervised learning," *Advances in neural information processing systems*, vol. 33, pp. 21271–21284, 2020.
- [20] Xinlei Chen and Kaiming He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 15750–15758.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [22] Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott, "Publicly available clinical bert embeddings," *arXiv preprint arXiv:1904.03323*, 2019.
- [23] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [24] Dani Kiyasseh, Tingting Zhu, and David A Clifton, "Clocs: Contrastive learning of cardiac signals across space, time, and patients," in *International Conference on Machine Learning*. PMLR, 2021, pp. 5606–5615.