

Contrastive Multi-modal Training with Electrocardiography and Natural Language Echocardiography Reports for Zero-shot Prediction of Structural Heart Disease

Wai-Chak WONG, BSc;¹ Che LIU, BEng MSc;^{2,3} Pierre ELIAS, MD;^{4,5} John Weston HUGHES, PhD;⁴ Chun-Yu LEUNG, MBBS;⁶ Xiao-Yan QIAN, PhD;⁷ Hang-Long LI, MBBS;¹ Yuk-Ming LAU, MBBS;¹ Chao-Fan TAO, PhD;⁷ Ali CHOO, MPhil;¹ Chi-Hang YUNG, MBChB;¹ Chi-Hong FONG, BSc;¹ Wai-Kwok CHOI, MBBS;¹ Chak-Kong CHENG, MBBS;¹ Lok-Lam CHENG, MBBS;¹ Lik-Man LAU, BSc(AC);¹ Roshan RELWANI, MBBS;¹ Jing QIN, PhD;⁸ Lequan YU, PhD;⁹ Hin-Wai LUI, PhD;¹⁰ Ho-On Alston Conrad CHIU, MBBS;^{1,11,12} Hung-Fat TSE, MD PhD;^{1,11,12,13} Chung-Wah SIU, MD MS;¹ Rossella ARCUCCI, PhD;^{2,3} Joshua Wing-Kei HO, PhD;¹⁴ Chun-Ka WONG, MBBS.^{1,11,12,13,15}

¹Department of Medicine, School of Clinical Medicine, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China.

²Data Science Institute, Imperial College London, United Kingdom.

¹⁶ ³Department of Earth Science and Engineering, Faculty of Engineering, Imperial College London,
¹⁷ United Kingdom.

¹⁸ Seymour, Paul, and Gloria Milstein Division of Cardiology, Department of Medicine, Columbia
¹⁹ University Irving Medical Center, USA.

20 ⁵Department of Biomedical Informatics, Columbia University, USA.

21 ⁶Tseung Kwan O Hospital, Hong Kong SAR, China.

22 ⁷Department of Electrical and Electronic Engineering, Faculty of Engineering, The University of
23 Hong Kong, Hong Kong SAR, China.

²⁴School of Nursing, The Hong Kong Polytechnic University, Hong Kong SAR, China.

²⁵School of Computing and Data Science, The University of Hong Kong, Hong Kong SAR, China.

²⁶ ¹⁰Department of Computer Science, University of California Irvine, USA.

27 ¹¹Department of Medicine, Queen Mary Hospital, Hong Kong SAR, China.

28 ¹²Department of Medicine, Tung Wah Hospital, Hong Kong SAR, Chi

29 ¹³Cardiac and Vascular Center, The University of Hong Kong - Shenzhen Hospital, China.

30 ¹⁴School of Biomedical Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong,
31 Hong Kong SAR, China.

³² ¹⁵Cambridge Stem Cell Institute, University of Cambridge, United Kingdom.

33

34 Short title: MERL-ECHO for Structural Heart Disease Prediction

35 Word count: 4896

36 **Disclosure:** The authors report no conflict of interest.

37

38 Address for Correspondence:

39 Chun-Ka WONG
40 Clinical Assistant Professor
41 Department of Medicine, The University of Hong Kong
42 Hong Kong SAR, China
43 Tel: +852 22553597
44 Email: wongeck@hku.hk
45

46

ABSTRACT

47 **Background:** Machine learning models for predicting **structural heart disease (SHD)** from
48 electrocardiography (ECG) traditionally required structured echocardiographic data. The
49 potential of **echocardiography (ECHO)** natural language reports remains underused. We
50 describe **MERL-ECHO**, a multimodal model using **contrastive language-image pre-training (CLIP)**
51 that aligns ECG with ECHO natural language reports for zero-shot SHD prediction.

52

53 **Methods:** We conducted a multi-center retrospective study using paired ECG and ECHO natural
54 language reports from Queen Mary Hospital and Tung Wah Hospital in Hong Kong. MERL-ECHO
55 was trained on 45,016 pairs ECG-ECHO pairs. Performance was evaluated on an **internal test set**
56 covering 10 SHDs and on an **external test set** of 5,442 ECGs with ECHO-derived labels for 6 SHDs
57 from Columbia University Irving Medical Center, USA.

58

59 **Results:** The cohort included 8,192 patients (mean age 73.7 ± 16.5 years; 55.3% male). In the
60 **internal test set**, MERL-ECHO achieved an **average AUROC of 0.69**, with strongest performance
61 for left ventricular dilation (0.78), right ventricular systolic dysfunction (0.71), and tricuspid
62 regurgitation (0.71). In the **external test set**, the **average AUROC was 0.72**, with highest
63 performance for left ventricular systolic dysfunction (0.76) and aortic stenosis (0.76). **Pre-**
64 **training improved AUROC by up to 5%**, performance scaled with larger datasets, and **ResNet18**
65 **outperformed ViT-Tiny as ECG encoder by 7%**. Saliency analysis revealed interpretable ECG
66 features, including unexpected P-wave changes in aortic stenosis, suggesting novel disease
67 markers.

68

69 **Conclusions:** MERL-ECHO leverages ECHO natural language reports for multimodal training with
70 ECG. This CLIP-based model enables accurate **zero-shot prediction of SHDs** and highlights
71 interpretable ECG features with potential clinical relevance.

72

73 **Keywords:** artificial intelligence, deep learning, echocardiography, electrocardiogra

74

ABBREVIATIONS

- 75 **AUPRC** Area under the precision-recall curve
- 76 **AUROC** Area under the receiver-operating characteristic curve
- 77 **CNN** Convolutional neural network
- 78 **DNN** Deep neural network
- 79 **ECG** Electrocardiography
- 80 **ECHO** Echocardiography
- 81 **Grad-CAM** Gradient-weighted class activation mapping
- 82 **RNN** Recurrent neural network
- 83 **SHD** Structural heart disease

84

INTRODUCTION

85 The global burden of heart failure due to valvular heart disease and impaired cardiac function
86 has steadily increased over the past decades ¹⁻³. Early detection of structural heart disease
87 (SHD) is critical for enabling timely initiation of pharmacological and interventional therapies ⁴⁻⁷.
88 Recent advances have demonstrated that deep neural networks (DNNs) applied to 12-lead
89 electrocardiograms (ECGs) can accurately predict valvular heart disease ⁸⁻¹⁰.

90

91 The predictive accuracy of DNN models typically improves with larger training datasets ¹¹.
92 However, most prior models have relied on supervised learning with convolutional neural
93 networks (CNNs) or recurrent neural networks (RNNs), which require structured datasets with
94 ground-truth SHD severity labels. In clinical practice, echocardiography (ECHO) data may be
95 captured as natural language free-text reports, leaving a vast amount of valuable information
96 underutilized for ECG classifier training.

97

98 In the field of image classification, contrastive language-image pre-training (CLIP) has emerged
99 as a groundbreaking self-supervised learning paradigm. This method enables training of image
100 classifiers from paired images and natural language text without the need for manual labeling
101 ¹². By leveraging large and diverse datasets, CLIP substantially expands the scope of model
102 development. Building on this concept, we recently introduced Multimodal ECG Representation
103 Learning (MERL), a model that applies the CLIP architecture for self-supervised training of ECG
104 classifiers using paired ECG signals and natural language ECG reports ¹³. We hypothesize that

105 MERL can be extended to incorporate paired ECG signals and natural language ECHO reports for
106 the prediction of SHD.

107

108 In this multi-center study, we present MERL-ECHO, a CLIP-based model trained using
109 contrastive multimodal learning on paired ECG signals and natural language ECHO reports. We
110 demonstrate that MERL-ECHO enables zero-shot prediction of SHD across internal and external
111 test sets. Furthermore, we systematically evaluate different architectural designs and report
112 comparative results to guide future research in multimodal ECG-based SHD prediction.

113

114

115

116

METHODS

117 **Study Design**

118 This was a multi-center retrospective study involving patients from Queen Mary Hospital and
119 Tung Wah Hospital, Hong Kong. This study conformed to the Declaration of Helsinki. Follows
120 CONSORT-AI (Consolidated Standards of Reporting Trials-Artificial Intelligence) checklist for
121 reporting¹⁴. Consent from individual patients was waived as only anonymized data from the
122 hospital registries was involved. Ethical approval was granted by HKU/HA HKW Institutional
123 Review Board (HKU/HA HKW IRB; UW 25-042).

124 **2.2**

125 **Data Source**

126 Adult patients aged 18 years or above having undergone 12-lead ECGs and echocardiograms at
127 Queen Mary Hospital and Tung Wah Hospital, Hong Kong, from January 2008 to May 2025 were
128 identified. ECGs were grayscale for the period 2008-2023 and coloured for the period 2023-
129 2025. Natural language ECHO text reports for 2008-2025 were acquired. ECGs and natural
130 language ECHO text reports were paired up. Only data pairs with the time difference between
131 the ECG acquisition date and the ECHO text report date less than or equal to three years were
132 included in the final study cohort. Each paired data was randomly assigned to either the
133 training, validation, or test set (Figure 1).

134

135 For the external test data set at Columbia University Irving Medical Center, all data were
136 collected from adult patients aged 18 or above who had a digitally stored 12-lead ECG and a
137 transthoracic echocardiogram within a 1-year interval between 2008 and 2022 (Figure 1). ECG

138 waveform data were extracted from the GE MUSE ECG management system at a sampling
139 frequency of 250 Hz across all 12 leads. Echocardiographic data were extracted from the Syngo
140 Dynamics (Siemens) and Xcelera (Philips) systems. Each ECG was accompanied by the echo-
141 derived binary labels, which were derived from structured echocardiogram reports and
142 binarized using clinically relevant thresholds to indicate at least moderate disease severity.

143 2.3

144 **ECG pre-processing**

145 Our training data were stored as images rather than raw digital signals. We utilized a simplified
146 version of our previously developed DigitHeart pipeline to extract voltage-time series data for
147 downstream use¹⁵. First, a Yolo v11 object detection model was used to locate and label each
148 of the 12 leads in ECGs¹⁶ (Supplemental Figure S1), 99.5% of ECGs were successfully detected
149 and cropped to the 12-lead region; the remaining ECGs (0.479%) were manually annotated to
150 identify the 12 leads. Second, binary thresholding using OpenCV v4.10.0 with a cut-off value of
151 80 was used to grossly remove grid lines, followed by using connected components size filtering
152 to further remove residual smaller grid lines¹⁷. Third, two-dimensional binary images were
153 converted into one-dimensional digital signals. The conversion algorithm raster-scanned the
154 binary image matrix along its horizontal axis, which serves as the temporal domain, while the
155 vertical axis represents the uncalibrated amplitude. For each pixel column (a discrete time
156 step), the vertical centroid of the white pixels (representing the ECG waveform) was computed
157 by averaging their y-coordinates. This centroid value represented a single sample point in the
158 raw time series. If no white pixels were found in a column, the vertical position was
159 interpolated by carrying forward the value from the preceding time step to maintain signal

160 continuity. To reduce noise introduced during the image-to-signal conversion, the extracted
161 signal was smoothed using a convolution-based moving average filter with a fixed window size
162 of 10 applied. Finally, the extracted ECG waveform data in the internal dataset and ECG
163 waveform data in the external dataset were up-sampled to a frequency of 500 Hz and stored
164 in .mat format. Upsampling the digital signal to 500 Hz ensures that the fine-tuned data inputs
165 are consistent with the pre-training ECG data during the pre-training stage of MERL-ECHO,
166 thereby minimizing potential distribution shift arising from disparate sampling rates. To meet
167 the model's input requirement of 10-second ECG recordings sampled at 500 Hz, the up-sampled
168 ECG waveforms were repeated as needed to achieve the required duration.

169 24

170 **ECHO natural language text report pre-processing**

171 First, we used the Qwen3-4B model, a general-purpose large language model, to analyze the
172 natural language text of each ECHO report¹⁸. Our objective was to exclude ECG-ECHO text
173 pairs whose reports documented any cardiac valve interventions listed in Supplemental Table
174 S1; only pairs linked to reports without these interventions were included. Second, the cardiac
175 abbreviation terms were expanded to the full terms, such as "AS" to "aortic stenosis" and "MR"
176 to "mitral regurgitation". Details are provided in the Supplemental Table S2. Empty or ECHO
177 reports with less than 4 words were excluded. This text pre-processing preserves the ECHO text
178 in a form that is as close as possible to its original state. We will also explore whether imposing
179 a more structured representation of ECHO text yields different results; the structured pre-
180 processing approach will be discussed in a subsequent section.

181

2.5

182 Diagnosis annotation

183 To validate the performance of our MERL-ECHO model, we developed a web application for
184 annotating a subset of our internal data sets to serve as ground truth (Supplemental Figure S2).
185 A committee consisting of 3 cardiologists and trained personnel was responsible for annotating
186 10 types of labels for each ECHO report, including aortic, mitral, and tricuspid stenosis and
187 regurgitation, left and right ventricular dilation and systolic impairment, and left ventricular
188 hypertrophy based on the same thresholding criteria. The first 500 ECHO reports were manually
189 annotated. Thereafter, the remaining 4,016 ECHO reports were annotated using a semi-
190 automated, human-in-the-loop workflow. In this pipeline, locally deployed large language
191 model Qwen3-4B was used to generate initial predictions for each report (average F1 score
192 0.77; Supplemental Figure S3; Supplemental Table S3)¹⁸. Subsequently, all reports were
193 reviewed and corrected by human annotators to ensure the final accuracy and fidelity of all
194 labels. This hybrid approach significantly accelerated the data curation process while
195 maintaining expert-level quality.

196 2.6

197 CLIP model development and training

198 Our methodology is predicated on a multimodal learning framework designed to train a zero-
199 shot diagnosis predictor, MERL-ECHO, to predict SHDs from ECGs by integrating ECG signals
200 with their corresponding natural language ECHO text reports. This approach, inspired by the
201 principles of CLIP and adapted for the medical signal-text domain¹², leverages the rich semantic
202 information contained within clinical narratives to inform and structure the learned ECG
203 representations. We previously described the use of CLIP architecture to train a zero-shot ECG

204 diagnosis predictor (MERL) by using ECG signal and ECG text report pairs¹³. In this study, we
205 want to detect SHDs from ECGs by training a model (MERL-ECHO) using paired ECGs and natural
206 language ECHO text reports.

207

208 MERL-ECHO is a dual-encoder architecture, featuring two distinct encoders for ECG signals and
209 ECHO text reports, which are used to process the corresponding modalities, respectively.
210 During the training process, all parameters in the ECG encoder will be updated, while only
211 those in the last 6 layers of the text encoder will be updated. In our whole dataset \mathcal{X} , we
212 represent each ECG-ECHO text pair as (e_i, r_i) , where as $e_i \in \mathcal{E}$ denotes the raw ECG signal and
213 $r_i \in \mathcal{R}$ denotes the associated natural language ECHO text report, respectively, with $i =$
214 $1, 2, 3, \dots, N$. Two contrastive learning strategies are utilized and described as follows.

215

216 The first and most important is the Cross-Modal Alignment (CMA). It applied a contrastive
217 learning approach on ECG-ECHO text pairs (Figure 2A) and enables MERL-ECHO to align the
218 feature representation of an ECG signal with the feature representation of its corresponding
219 natural language ECHO text report. ECG encoder with 1D-ResNet18¹⁹ or ViT-Tiny²⁰ as the
220 backbone embeds the ECG signal, and Med-CPT²¹ as the text encoder embeds the natural
221 language ECHO text report. Each ECG-ECHO text pair (e_i, r_i) will be embedded into the shared,
222 high-dimensional latent embedding space, denoted as $(z_{e,i}, z_{r,i})$, by the corresponding
223 encoders. Subsequently, the ECG embedding and ECHO-text embedding will be mapped into
224 the same dimensionality d by two different non-linear projectors (P_e and P_r), they are
225 represented as $\hat{z}_{e,i}$ and $\hat{z}_{r,i}$, respectively. Third, cosine similarities for ECG-ECHO text and ECHO

226 Text-ECG will be computed as $s_{i,i}^{e2r} = \hat{e}_i^\top \hat{r}_i$ and $s_{i,i}^{r2e} = \hat{r}_i^\top \hat{e}_i$, respectively. $\mathcal{L}_{i,j}^{e2r}$ and $\mathcal{L}_{i,j}^{r2e}$
227 represent ECG-ECHO text and ECHO text-ECG cross-modal contrastive losses, the final loss
228 function (\mathcal{L}_{CMA}) in cross-modal alignment is the average of these two losses (Formula 1),
229 allowing MERL-ECHO to be trained to maximize the cosine similarity of embeddings from true
230 pairs while minimizing the cosine similarity of non-paired embeddings, and making the
231 alignment learning robust and bidirectional:

232

233 **Formula 1:**

234
$$\mathcal{L}_{i,j}^{e2r} = -\log \frac{\exp(s_{i,j}^{e2r}/\tau)}{\sum_{k=1}^L \mathbb{1}_{[k \neq i]} \exp(s_{i,k}^{e2r}/\tau)}$$

235
$$\mathcal{L}_{i,j}^{r2e} = -\log \frac{\exp(s_{i,j}^{r2e}/\tau)}{\sum_{k=1}^L \mathbb{1}_{[k \neq i]} \exp(s_{i,k}^{r2e}/\tau)}$$

236
$$\mathcal{L}_{CMA} = \frac{1}{2L} \sum_{i=1}^N \sum_{j=1}^N (\mathcal{L}_{i,j}^{e2r} + \mathcal{L}_{i,j}^{r2e})$$

237

238 Cross-Modal Alignment strategy not only aligns ECG signal features and semantic features from
239 natural language ECHO text report, but it also enables MERL-ECHO the capability of zero-shot
240 classification (Figure 2B). As ECG and ECHO-text features are aligned in the same high-
241 dimensional latent embedding space during the training process, ECG signal features are
242 injected with high-level, human-understandable semantic information from ECHO-text. We can
243 test the model by simply providing a text prompt and an ECG signal. The model can then find
244 matching ECGs by calculating feature similarity, without requiring any labelled data or fine-
245 tuning for the downstream task.

246

247 In addition to the alignment between the features of ECG and ECHO-text, we further employ
248 Uni-Modal Alignment (UMA). It is also a contrastive learning method, but it operates within a
249 single modality, specifically the ECG domain (Figure 2A). Its purpose is to learn higher-quality,
250 more discriminative feature representations for ECGs. There are three steps in this workflow:
251 First, the ECG encoder embeds the ECG signal to obtain $z_{e,i}$. Then, embedding $z_{e,i}$ is duplicated.
252 Independent and random dropout with a probability (p) 0.1 is applied to these two identical
253 embeddings for generating the positive pair $(z_{e,i}^1, z_{e,i}^2)$ ²² (Supplemental Figure S4 for notation
254 details). Third, we use standard contrastive loss on the positive pair and treat other unpaired
255 combinations as negative pairs (Formula 2). This latent augmentation technique, by using two
256 independent dropout operations to construct positive pairs, allows MERL-ECHO to learn ECG
257 signal features that are more powerfully discriminative. It works in tandem with CMA (Cross-
258 Modal Alignment) to form the training core of the MERL-ECHO framework.

259

260 **Formula 2:**

$$261 \quad \mathcal{L}_{UMA} = -\frac{1}{L} \sum_{i=1}^N \sum_{j=1}^N \log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^L \mathbb{1}_{[k \neq i]} \exp(s_{i,j}/\tau)}$$

$$262 \quad \text{where } s_{i,i} = z_{e,i}^{1\top} z_{e,i}^2$$

$$263 \quad z_{e,i}^1 = z_{e,i} \odot M^1, M^1 \sim \text{Bernoulli}(p),$$

$$264 \quad z_{e,i}^2 = z_{e,i} \odot M^2, M^2 \sim \text{Bernoulli}(p).$$

265

266 To further enhance performance, we also included a pre-training phase based on the previous
267 work from our group¹³, where a CLIP-based model was trained on 771,693 ECG digital signals
268 and ECG text report pairs initially on the publicly accessible dataset, MIMIC-ECG²³, and it
269 enabled zero-shot classification for various electrocardiographic abnormalities, we leverage
270 pre-trained model weights to facilitate knowledge transfer (Supplemental Figure S5) from ECG-
271 based representations to echocardiographic insights.

272

273 We trained the model for 15 epochs on the training set with the AdamW optimizer, setting a
274 learning rate of 1×10^{-5} and a weight decay of 1×10^{-8} . We unfroze the last six layers of the
275 text encoder. We applied a cosine annealing scheduler for learning rate adjustments and
276 maintained a batch size of 256 per GPU (two GeForce RTX 4090D). We set up early stopping
277 criteria if the evaluation loss did not improve for three consecutive epochs.

278 *L.7*

279 Prompts for zero-shot classification during inference

280 During inference, a key feature of our CLIP-based model is the capability to perform zero-shot
281 classification. Our model can predict any arbitrary diagnosis from a new ECG by inputting text
282 prompts, without the need for retraining (Figure 2B). By directly leveraging ECG waveform
283 signals along with carefully engineered text prompts for 10 distinct echocardiographic
284 prediction classes, the model can predict the binary classification for each of the
285 echocardiographic labels for the ECG by calculating the cosine similarity between the ECG
286 embedding and each text prompt embedding. The ECG-ECHO text pair with the highest
287 similarity will be the output.

288

289 To prevent model overfitting to the text prompt during the training stage, we utilized the
290 diagnostic category names as the text prompts. Prompts for zero-shot inference were explicitly
291 designed to enhance model prediction performance.

292

293 A prompt is a natural language text describing the semantic context of an image or a task. By
294 designing appropriate prompt engineering, the CLIP-based model can better align the
295 echocardiographic textual description, the prompt, with the ECG digital signal. To optimize
296 prompt design, we introduced enhanced cardiology domain prompt engineering to allow
297 prompts with variant phrases and associated features. The complete prompt list is in
298 Supplemental Table S4. We first queried a large language model, GPT-4, which has been shown
299 to achieve high accuracy in generating cardiologist-level ECHO reports²⁴⁻²⁷, to generate
300 prompts for 10 echocardiographic labels. Query details are in Supplemental Table S5. Following
301 prompt generation, experienced cardiologists reviewed and refined the prompts based on
302 accuracy, clinical relevance, and clarity criteria. This multi-step validation ensures that the
303 generated prompts accurately reflect the underlying cardiac pathophysiology and are
304 contextually appropriate for domain-specific zero-shot classification tasks.

305

2.8

306 **Saliency Mapping**

307 To better understand the ECG representations that are most important to model prediction, we
308 used the Gradient-weighted Class Activation Mapping (Grad-CAM) technique to visualize ECGs
309 with positive classes and the highest confidence model predictions²⁸. To do this, our model first

310 receives both ECG signals and text prompts simultaneously, extracting features through the
311 ECG encoder and text encoder, respectively, and calculating similarity scores via dot products,
312 which serves as the target output of Grad-CAM analysis. To generate the saliency map, we
313 targeted the final convolutional layer within the last residual block of our ECG encoder. During
314 the backward propagation, the gradients of the similarity score with respect to the feature
315 maps of this target layer were computed. These gradients were then global-average-pooled to
316 obtain a set of weights, representing the importance of each feature map. A weighted linear
317 combination of the feature maps was computed using these weights, followed by a rectified
318 linear unit (ReLU) activation to isolate features with a positive influence on the prediction. The
319 resulting saliency map is then up-sampled to the original length of 5000 time steps using linear
320 interpolation. This global saliency map is further overlaid onto each of the 12 leads of the
321 original ECG waveform. It provides a comprehensive view, where the heat intensity (color-
322 coded in red) indicates the temporal importance, as determined by the model, allowing for a
323 qualitative assessment of which segments of the cardiac cycle contribute most significantly to
324 the final similarity score.

325 2.9

326 Statistical Analysis

327 Descriptive statistics using counts, percentages, and medians were provided as appropriate.
328 Discrete variables were presented as mean \pm standard deviation. F1 score, area under the
329 receiver operating curves (AUROC), and area under the precision-recall curves (AUPRC) were
330 calculated for the model prediction on all classes in both the internal and external test sets to
331 assess their performance. The F1 score was calculated as the harmonic mean of the precision

332 and recall using a method that achieved the best balance between precision and recall; a series
333 of thresholds were calculated to get different values of precision and recall, and the F1 score
334 was chosen when the calculated precision and recall resulting the maximum F1 score based on
335 the specific threshold value that was also the cut-off value used to distinguish between positive
336 and negative class. It helps the model to find an optimal classification threshold that achieved
337 the best predictive performance (measured by the F1 score). The range of F1 score values is
338 from 0 to 1, where 1 indicates perfect classification performance, and 0 indicates the worst
339 possible performance. The AUROC was used to describe model accuracy, plotting sensitivity
340 (true positive rate) against 1-specificity (false positive rate) with values ranging from 0.5 to 1,
341 where 1 indicates perfect classification, and 0.5 indicates random guessing. The AUPRC was also
342 used to describe model accuracy, plotting precision (positive predictive value) against recall
343 (sensitivity). All statistical analyses were performed using Python 3.12.3.

344 3.

RESULTS

345 3.1

Patient Characteristics

346 A total of 45,016 unique ECG signals and ECHO text pairs were identified in 8,192 unique
347 patients from 2 public hospitals in Hong Kong, namely Queen Mary Hospital and Tung Wah
348 Hospital, respectively, of whom 4530 (55.3 %) were male, and 3662 (44.7 %) were female. Their
349 mean age was 73.7 ± 16.5 years. Patient characteristics for the training, validation, test, and
350 external test sets are described in **Table 1**.

351 3.1

Model Performance

354 The model performance in AUROC and AUPRC of each echocardiographic finding label is shown
355 in Figure 3 and Table 2. In the internal test set, the average AUROC is around 0.69. It varied
356 slightly across categories and diseases (Figure 3A). For ventricular function, the model achieved
357 the highest AUROC of 0.71 in right ventricular systolic function impairment (RV impaired). In the
358 category of ventricular size, the highest AUROC was found for left ventricular dilation (LV
359 dilated; 0.78). Among valvular heart diseases, moderate to severe tricuspid regurgitation had
360 the highest AUROC of 0.71. And left ventricular hypertrophy had the lowest AUROC of 0.64. The
361 AUROC values for moderate to severe aortic and mitral stenosis, as well as aortic and mitral
362 regurgitation, were 0.67. For the area under the precision-recall curves and F1 score, MERL-
363 ECHO achieved the highest AUPRC of 0.63 and the highest F1 score of 0.66 in left ventricular
364 systolic function impairment (LV impaired) (Figure 3B).

365
366 In the external test set (Columbia University Irving Medical Center), patients' characteristics
367 were similar to those in the internal data set in terms of sex and age (Table 1). Surprisingly,
368 MERL-ECHO even performs better in the external test set than in the internal test. The average
369 AUROC for 6 echocardiographic classification labels is 0.72 (Figure 3C). The AUROC for left
370 ventricular systolic function impairment and moderate to severe aortic stenosis was the highest
371 (AUROC: 0.76), and the lowest AUROC was found for moderate to severe aortic regurgitation
372 (AUROC: 0.68). For the area under the precision-recall curves and F1 score, MERL-ECHO
373 achieved the highest AUPRC of 0.48 and the highest F1 score of 0.49 in left ventricular systolic
374 function impairment (LV impaired) (Figure 3D). Although there were variations in the AUROC

375 for some echocardiographic labels, the model showed its ability to be robust and generalizable
376 to the external data.

377 3.3

378 **ECG Encoder**

379 For these two ECG feature extractor networks, the CNN-based ResNet18 and transformer-
380 based ViT-Tiny, the zero-shot performance using different data proportions and pre-training on
381 ECG-ECG text data shows that our CLIP-based model with CNN-based ResNet18 generally
382 surpasses ViT-Tiny's (Figure 4A). The ECG encoder using ResNet18 as the backbone surpasses
383 the ViT-Tiny one by almost 7% in AUROC when using the whole training dataset (Table 3).

384 3.4

385 **Effect of ECG-ECG Text Pre-training**

386 Pre-training is the initial stage of the model training; it allows the model to learn general ECG
387 features and patterns by initially training on a large amount of ECG-ECG text pair data without
388 manually labelled ground-truth, and provides a foundation for subsequent fine-tuning. With
389 pre-training on ECG-ECG text pair data, we observed that fine-tuning on the complete ECG-
390 ECHO text training set led to improved model performance for ResNet18 as ECG encoders
391 (Table 3), compared to models trained without pre-training. Specifically, for the CNN-based
392 ResNet18 model, pre-training resulted in the highest AUROC of around 69% in the internal test
393 set, representing a 5% improvement over the model without pre-training.

394 3.5

395 **Effect of Scaling Training Data**

396 We also conducted comprehensive experiments to investigate the effect of scaling training data
397 on AUROC. All the models with ResNet18 as the ECG encoder, with the pre-training and without
398 the pre-training on ECG-ECG text data, improve AUROC steadily (Figure 4A). Crucially, the
399 performance trend did not show signs of reaching a plateau, as AUROC continued to increase
400 steadily up to the full dataset size. Specifically, with pre-training, the model's performance with
401 ResNet18 as ECG encoder rises by 9% in AUROC when using the whole training dataset
402 compared to using 20% of data only (Table 3), and the ResNet-based encoder outperformed its
403 transformer-based Vit-Tiny counterpart, achieving a 6% higher AUROC. This sustained upward
404 trend suggests that there is still potential for increasing AUROC because of the continuously
405 growing trend for AUROC.

406 *3-6*

407 **Ablation Study on Text Preprocessing**

408 To evaluate the necessity of explicit feature engineering within the natural language clinical
409 text, we conducted an experiment comparing two text-processing strategies. The first strategy
410 involved training the model on natural language ECHO text reports, where only common
411 abbreviations were expanded into their full terms. The second strategy involved changing the
412 numerical values of the indicator into a specific diagnosis (Supplemental Table S6), making
413 natural language in a more structured way. For instance, “LVEF = 45%” was transformed into
414 “moderate left ventricular systolic function impaired”. Our results showed that providing
415 explicit categorical labels for LVEF yielded no performance benefit. In fact, the model trained on
416 unaltered, natural language text achieved a slightly superior AUROC. The model using natural
417 language text reached an average AUROC of 0.72 in the external test set, marginally

418 outperforming the model trained on categorized text, which scored an average AUROC of 0.71.
419 Specifically, the model trained on natural language ECHO text achieved 0.76 AUROC in left
420 ventricular systolic function impairment in the external test set, marginally outperforming the
421 model trained on categorized text, which scored 0.75 AUROC in the same category. This finding
422 suggests that the model is capable of learning the clinical significance of numerical LVEF values
423 directly from the contextual information present in the raw text, rendering manual
424 categorization unnecessary.

425 *3.7*

426 **Prompt Design Optimisation**

427 In addition to the prompts generated by using our enhanced cardiology domain prompt
428 engineering strategy, two more types of prompts, term-based prompts, and variant phrase
429 prompts, were generated by GPT-4 to investigate the effects of prompts in zero-shot
430 classification. Query details for the prompts generation are in Supplemental Table S5. Example
431 prompts for left ventricular dilation and its variations are shown below:

432

- 433 1. **Terms only:** Consisting solely of the diagnostic category name.
 - 434 • Example: "*Dilated left ventricle*."
- 435 2. **Terms with variant phrases:** Producing multiple phrasings or syntactic variants for each
 - 436 prediction class to capture linguistic diversity.
 - 437 • Example: "*Dilated left ventricle, LV dilation, LV enlargement (LVEDD > 58 mm in*
438 *males or > 52 mm in females)*."

439 3. Terms with variant phrases and associated features: Incorporating additional
440 echocardiographic descriptive phrases that complement and do not contradict the
441 definitions of the prediction classes.

- 442 • Example: "severely dilated LV with impaired LVEF at 30%, moderate mitral
443 regurgitation (MR), and dilated left atrium (LA)."

444

445 Our prompt list is presented in Supplemental Table S4, with terms with variant phrases and
446 associated features achieving the highest average AUROC of 0.69 among them (Figure 4B;
447 Supplemental Table S5). The advantage of using prompts of terms with variant phrases and
448 associated features is mainly demonstrated in the zero-shot inference of left ventricular dilation
449 and systolic function impairment, and right ventricular systolic function impairment. We
450 conclude that making the prompt more comprehensive and detailed may be helpful in
451 improving accuracy.

452 3.8

453 Saliency Mapping

454 To better understand which ECG features contributed most to model predictions, we applied
455 Grad-CAM to visualize salient regions identified by the MERL-ECHO model. This approach helps
456 address the interpretability challenges of black-box algorithms. Representative Grad-CAM
457 images are shown in Figure 5. In the illustrative ECG predicted as high risk for left ventricular
458 impairment in figure 5A, the model highlighted poor R-wave progression in the QRS complex
459 and premature ventricular complexes as important features. These findings are consistent with
460 established clinical knowledge ²⁹⁻³¹.

461

462 Moreover, MERL-ECHO consistently emphasized the P wave as the most important feature for
463 predicting aortic stenosis (Figure 5B). This observation was unexpected and prompted further
464 investigation. We therefore examined P-wave morphologies in patients with and without aortic
465 stenosis (Supplemental Method S1). Quantitative analysis revealed that patients with aortic
466 stenosis had a higher prevalence of flattened P waves in leads I, V5, and V6 (Figures 6A and 6B,
467 and Supplemental Tables S7 and S8).

468 

469

DISCUSSION

470 In this study, we developed MERL-ECHO, the first CLIP-based multimodal learning framework
471 that aligns 12-lead ECG signals with natural language ECHO reports using contrastive learning.
472 MERL-ECHO allows accurate zero-shot classification of 10 SHD conditions, achieving an average
473 AUROC of 0.69 in the internal test set and 0.72 in the external validation cohort, with the
474 strongest performance in predicting left ventricular systolic impairment and dilation.
475 Importantly, saliency mapping revealed clinically meaningful ECG features contributing to
476 predictions, highlighting the model's interpretability. Overall, MERL-ECHO shows strong
477 potential as a scalable, generalizable framework for detecting structural heart diseases from
478 ECGs without the need for labelled training data.

479

480 Previous work in ECG classification has largely relied on single-modality models using self-
481 supervised or supervised learning. Neural networks such as RNNs and CNNs³²⁻³⁴ demonstrated
482 the ability to learn ECG feature representations and reduce the need for handcrafted rules.

483 However, these models remain constrained by their design: they only utilize ECG signals,
484 require extensive manual labelling of diagnoses from electronic health records (EHRs), and are
485 limited to fixed label sets. As a result, they overlook the abundant contextual knowledge
486 embedded in natural language clinical text and cannot naturally extend beyond predefined
487 outputs.

488

489 In contrast, CLIP provides a more flexible and scalable paradigm. Unlike conventional neural
490 networks that require fixed, manually curated label sets and labour-intensive relabelling for
491 each new task, CLIP-based architectures directly align ECG signals with free-text reports already
492 available in electronic health records. This removes the bottleneck of extensive expert
493 annotation, since the natural language text itself provides supervision. As a result, new
494 diagnostic concepts can be incorporated simply by adding text prompts, without retraining the
495 model. To our knowledge, this is the first demonstration of leveraging ECHO natural language
496 text report within a multimodal CLIP model to link ECG with echocardiographic findings.

497

498 The CLIP architecture utilized contrastive learning methods to learn the relationship between
499 ECG and ECHO text reports with two input encoders: an ECG encoder and a text encoder. Our
500 experiments revealed a stark contrast in how different ECG encoder architectures respond to
501 increased training data. The CNN-based ResNet18 model with ECG-ECG text pretraining
502 demonstrated consistent and steady performance gains as the dataset size grew, with its
503 AUROC increasing by 9% when scaling from 20% to 100% of the training data. Notably, with
504 ECG-ECG text pretraining, using ResNet18 as the ECG encoder results in a higher AUROC of 7%

505 than using ViT-Tiny as the ECG encoder; the reason behind this may stem from the different
506 processing methods. CNN-based ResNet18 processes ECG digital signals through convolutional
507 layers, which are good at extracting important local features from time-series data like ECG
508 digital signals. Transformer-based Vit-Tiny does not have such an architecture and may struggle
509 to extract local information. The continuous upward trend suggests that its performance has
510 not yet reached saturation and could be further enhanced with additional data. This highlights
511 the potential of training a large-scale foundation model on substantially larger ECG-ECHO
512 datasets, which could enable more robust zero-shot predictions across an even broader range
513 of cardiac conditions.

514

515 Saliency mapping confirmed expected predictors of left ventricular impairment, including
516 abnormal R-wave progression and ventricular ectopy, consistent with existing clinical
517 understanding. Unexpectedly, the model consistently highlighted the P wave as most important
518 for predicting aortic stenosis. Further analysis showed a higher prevalence of flattened P waves
519 in leads I, V5, and V6 among affected patients. Previous ECG studies in aortic stenosis have
520 mainly focused on QRS complex, ST segment and T wave³⁵⁻³⁸. Atrial findings have usually been
521 limited to changes in P-wave duration and dispersion^{39,40}. Subtle P-wave flattening has not
522 been systematically reported. These findings suggest that atrial remodeling in aortic stenosis
523 may be detectable through P-wave. Independent validation and mechanistic studies are needed
524 to clarify whether this feature provides additional diagnostic or prognostic value.

525

526 One possible future development of CLIP architecture for ECG classification is to expand the
527 multi-modal training to ECHO videos, as the raw spatial information may also provide useful
528 features for training ECG classifiers. However, it is well recognized that ECHO videos typically
529 require multi-frame sampling or temporal modeling^{41,42}, which increases inference time and
530 GPU memory pressure. It is therefore technical challenging to train such a model at the
531 moment.

532 

533 **Limitations**

534 Our model has several limitations that warrant consideration. First, the text encoder used in
535 this study was primarily trained on English text. ECHO reports written in other languages would
536 require a language-specific encoder to ensure reliable performance. Second, as this is a
537 retrospective study, a prospective clinical trial is necessary to assess the model's accuracy and
538 generalizability in real-world clinical settings. Third, the presence of class imbalance in our
539 datasets may have contributed to overly optimistic AUROC values for certain minority classes
540 during testing.

541  5.

542 **CONCLUSION**

543 MERL-ECHO demonstrates that aligning ECG signals with natural language echocardiography
544 reports via contrastive learning enables robust, interpretable, and generalizable zero-shot
545 detection of structural heart disease.

546

547

ARTICLE INFORMATION

548 **Author Contributions**

549 WCW, CL and CKW conceived the study design. WCW, PE, JWH, CYL, XYQ, HLL, YML, AC, CHY, CHF,
550 WKC, CKC, LLC, LML, RR, JQ, HOACC and CKW acquired internal and external data. WCW, CL, XYQ,
551 YML, CFT, RR, JQ, LY, HWL, RA, JWKH and CKW contributed to machine learning model
552 development and statistical analysis. CYL, HHL, YML, HOACC, HFT, CWS and CKW provided clinical
553 expertise for the study. WC, YML and CKW wrote first draft of the manuscript. CL, CYL, HOACC,
554 HFT and CWS revised the manuscript for intellectually critical content. All authors read the final
555 manuscript and authorized its submission.

556

557 **Data Availability Statement**

558 Pre-training data (ECG-ECG text data) are all publicly accessible datasets from MIMIC-ECG ²³.
559 Fne-tuning data (ECG-ECHO text data) at Queen Mary Hospital and Tung Wah Hospital will not
560 be publicly available owing to institutional policy. External test data from Columbia are available
561 on PhysioNet <https://physionet.org/content/echonext/1.0.0/>. Codes and model weights are
562 available at <https://github.com/ConstantineWong/MERL-ECHO>. Other anonymized data are
563 available at reasonable request from the corresponding author for 3 years from date of
564 publication.

565

566 **Sources of Funding**

567 The study was partly supported by Rosie Young Medical Fellowship for Internal Medicine, Sun
568 Chief Yeh Heart Foundation, Hong Kong.

569

570 **Acknowledgement**

571 None.

572

573 **Disclosures**

574 The authors report no conflict of interest.

575



TABLES

576

577

Table 1. Patient Characteristics

	Train [#]	Validation	Test	External test
No. of ECG-ECHO Pairs	40,500	2,250	2,266	5,442
Patients, n	7,099 ^{##}	1,693 ^{##}	1,638 ^{##}	5,442
Age, mean \pm SD, y	73.5 \pm 16.6	72.4 \pm 15.8	71.4 \pm 16.7	71.1 \pm 16.7*
Sex, %				
Female	3,206 (45.2%)	764 (45.1%)	731 (44.6%)	2,731 (50.2%)
Male	3,893 (54.8%)	929 (54.9%)	907 (55.4%)	2,711 (49.8%)
Ventricular Size and Morphology, %				
Left Ventricular Hypertrophy	25.7	18.3	19.1	-
Dilated Left Ventricle	10.3	8.6	9.1	-
Dilated Right Ventricle	13.2	10.7	10.7	-
Ventricular Function, %				
Impaired Left Ventricular Systolic Function (Moderate to Severe)	46.6	43.5	44.3	17.7

Impaired Right Ventricular Systolic Function (Moderate to Severe)	12.7	10.3	10.9	7.7
Moderate to Severe Valvular Heart Diseases, %				
Aortic Stenosis (AS)	9.6	9.9	9.3	5.3
Aortic Regurgitation (AR)	20.8	10.2	9.7	1.2
Mitral Stenosis (MS)	12.1	11.6	12.3	-
Mitral Regurgitation (MR)	35.8	26.4	26.8	6.2
Tricuspid Regurgitation (TR)	30.0	28.4	27.8	6.5

578 #Multiple ECGs and ECHO text reports per patient were included to maximize model training.

579 ##The total number of patients across the training, validation, and test sets does not equal the size of the internal dataset. The first
580 reason is in # and the second reason is that the whole internal data set was randomly split into train, validation, and test sets.

581 *Age values are capped at 90 years for de-identification purposes.

582

Table 2. Diagnostic performance by AUROC and AUPRC using MERL-ECHO for each label

Echocardiographic Finding Label**	Internal Test		External Test	
	AUROC	AUPRC	AUROC	AUPRC
LV Dilated	0.78	0.27	-	-
LV Impaired	0.70	0.63	0.76	0.48
RV Dilated	0.68	0.21	-	-
RV Impaired	0.71	0.20	0.71	0.17
AS Mod Severe	0.67	0.15	0.76	0.14
AR Mod Severe	0.67	0.15	0.69	0.02
MS Mod Severe	0.67	0.20	-	-
MR Mod Severe	0.67	0.42	0.69	0.13
TR Mod Severe	0.71	0.50	0.72	0.17
LVH	0.64	0.30	-	-

583

**LV Dilated = Left Ventricular Dilation, LV Impaired = Left Ventricular Systolic Function Impairment, RV Dilated = Right Ventricular Dilation, RV Impaired = Right Ventricular Systolic Function Impairment, AS Mod Severe = Moderate to Severe Aortic Stenosis, AR Mod Severe = Moderate to Severe Aortic Regurgitation, MS Mod Severe = Moderate to Severe Mitral Stenosis, MR Mod Severe = Moderate to Severe Mitral Regurgitation, TR Mod Severe = Moderate to Severe Tricuspid Regurgitation, LVH = Left Ventricular Hypertrophy.

584

585

586

587

588

589

590

591

592

Table 3. AUROC comparison for two ECG encoders under different training data sizes

Metrics	Pretraining ^{###}	ECG Encoder	Training Set Size				
			8,100 (20%)	16,200 (40%)	24,300 (60%)	32,400 (80%)	40,500 (100%)
AUROC	✓	ResNet	59.9	63.1	64.3	66.9	68.9
		ViT	51.7	56.5	58.2	60.2	62.7
	✗	ResNet	49.4	56.0	57.4	59.7	63.6
		ViT	52.8	51.8	53.1	52.2	52.2

593

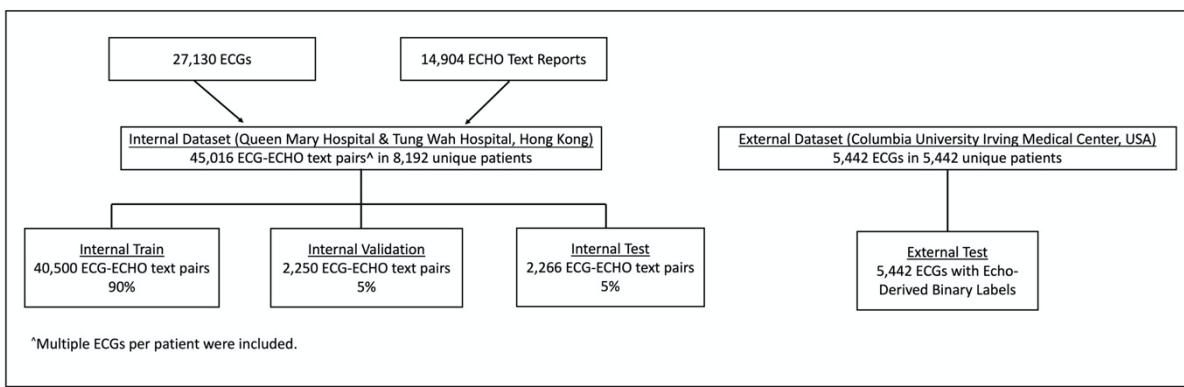
^{###}Pre-training with ECG-ECG text pair data.

594

FIGURE LEGEND

595

Figure 1. Patient Flow

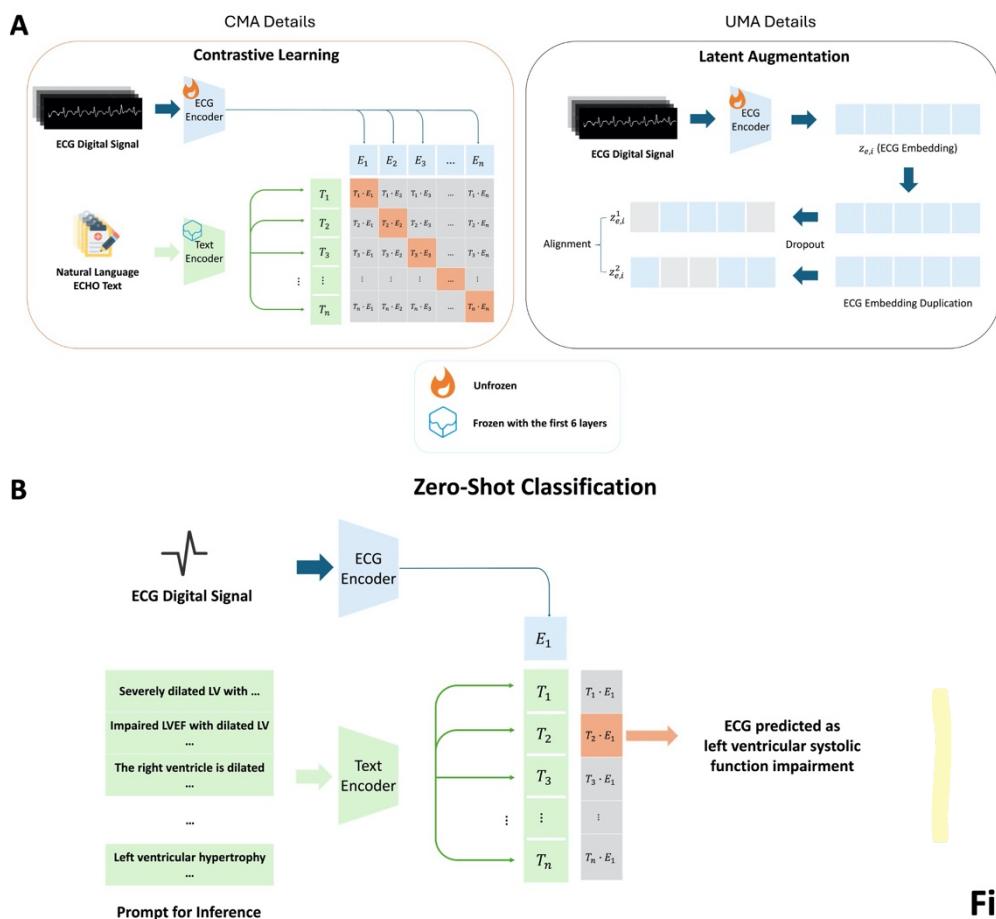


596

Figure 1

597

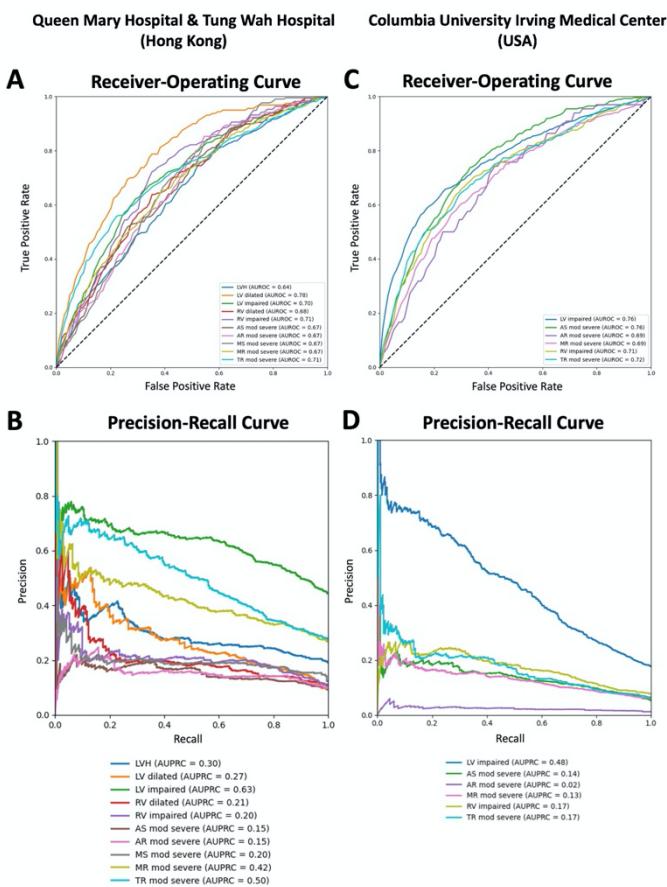
Figure 2. Model Architecture and Zero-shot Capability



598

599

Figure 3. AUROC and AUPRC for the Internal and External Test Sets



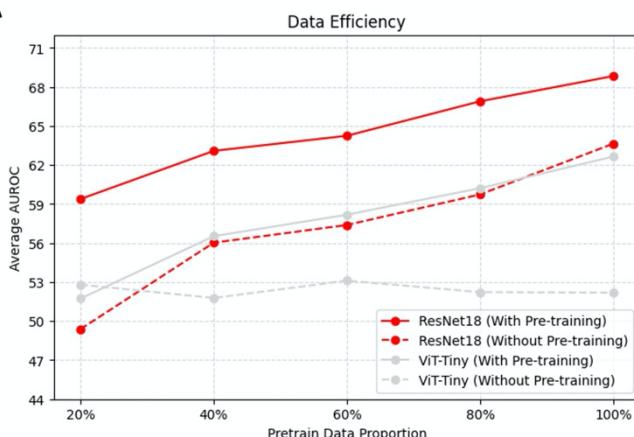
600

Figure 3

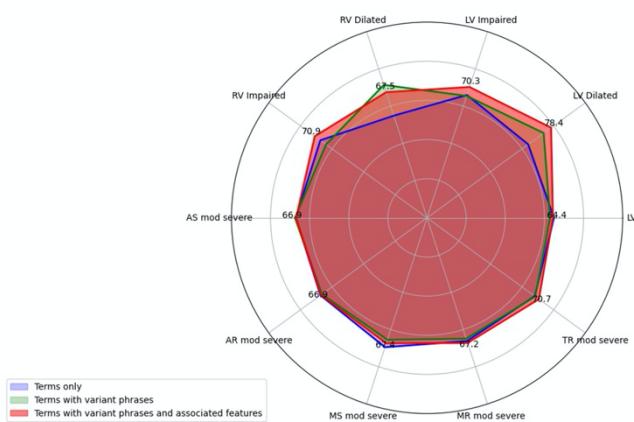
601

Figure 4. Data Efficiency Plot and Radar Plot for AUROC under Different Prompts

A



B

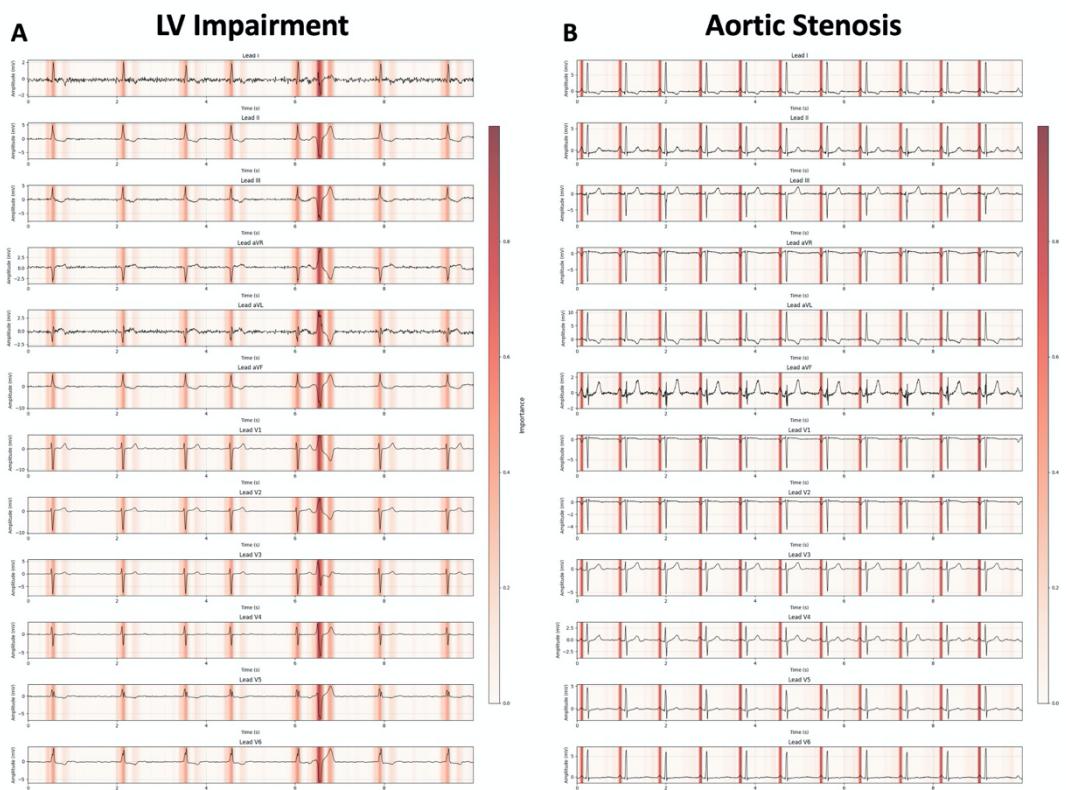


602

Figure 4

603

Figure 5. Saliency Map



604

Figure 5

605 **Figure 6. P-wave Morphology Distribution and P-wave Normalized Height Distribution for**
606 **Normal ECG and ECG with aortic stenosis (AS) labeled in the associated ECHO report**

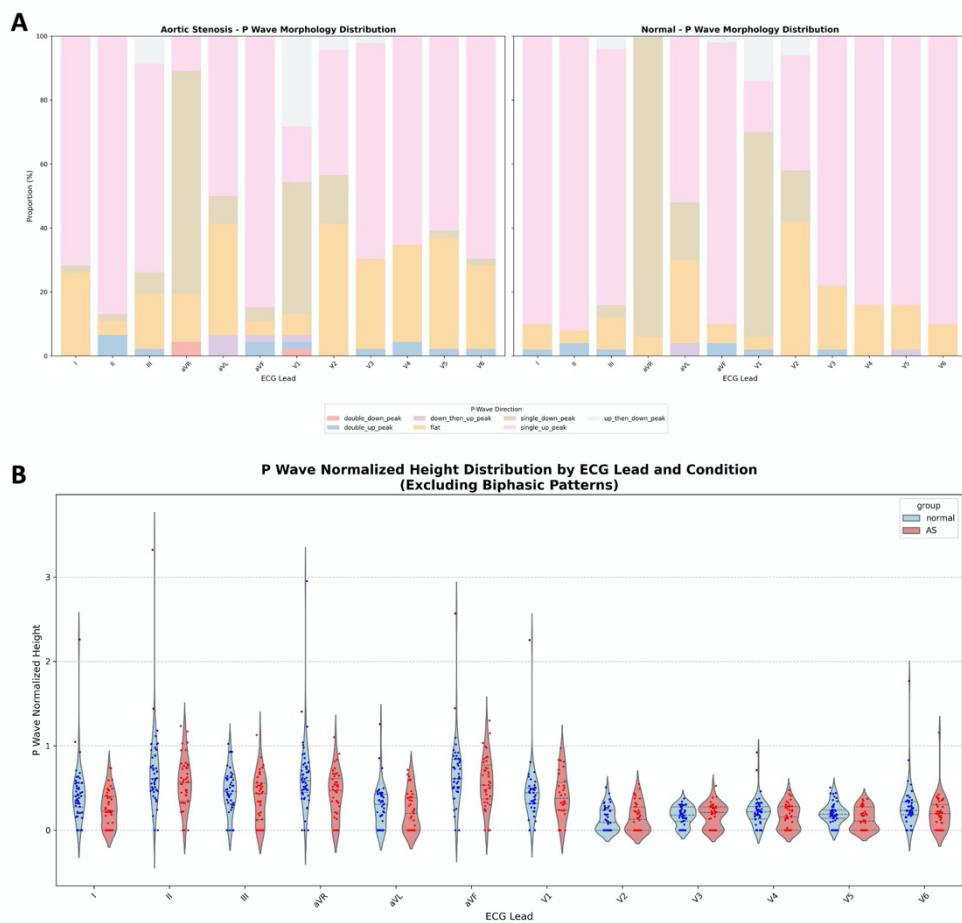


Figure 6

607
608
609
610
611
612
613
614
615

616

REFERENCES

- 617 1. Iung B, Vahanian A. Epidemiology of acquired valvular heart disease. *Canadian Journal
618 of Cardiology*. 2014;30:962-970.
- 619 2. Savarese G, Becher PM, Lund LH, Seferovic P, Rosano GM, Coats AJ. Global burden of
620 heart failure: a comprehensive and updated review of epidemiology. *Cardiovascular
621 research*. 2022;118:3272-3287.
- 622 3. Chen QF, Shi S, Wang YF, Shi J, Liu C, Xu T, Ni C, Zhou X, Lin W, Peng Y. Global,
623 regional, and national burden of valvular heart disease, 1990 to 2021. *Journal of the
624 American Heart Association*. 2024;13:e037991.
- 625 4. Packer M, Anker SD, Butler J, Filippatos G, Pocock SJ, Carson P, Januzzi J, Verma S,
626 Tsutsui H, Brueckmann M. Cardiovascular and renal outcomes with empagliflozin in heart
627 failure. *New England Journal of Medicine*. 2020;383:1413-1424.
- 628 5. McMurray JJ, Solomon SD, Inzucchi SE, Køber L, Kosiborod MN, Martinez FA,
629 Ponikowski P, Sabatine MS, Anand IS, Bělohlávek J. Dapagliflozin in patients with heart
630 failure and reduced ejection fraction. *New England Journal of Medicine*. 2019;381:1995-
631 2008.
- 632 6. Généreux P, Schwartz A, Oldemeyer JB, Pibarot P, Cohen DJ, Blanke P, Lindman BR,
633 Babaliaros V, Fearon WF, Daniels DV. Transcatheter aortic-valve replacement for
634 asymptomatic severe aortic stenosis. *New England Journal of Medicine*. 2025;392:217-
635 227.
- 636 7. Anker SD, Friede T, von Bardeleben R-S, Butler J, Khan M-S, Diek M, Heinrich J, Geyer
637 M, Placzek M, Ferrari R. Transcatheter valve repair in heart failure with moderate to severe
638 mitral regurgitation. *New England Journal of Medicine*. 2024;391:1799-1809.
- 639 8. Liang Y, Sau A, Zeidaabadi B, Barker J, Patlitzoglou K, Pastika L, Sieliwonczyk E,
640 Whinnett Z, Peters NS, Yu Z. Artificial intelligence-enhanced electrocardiography to
641 predict regurgitant valvular heart diseases: an international study. *European Heart Journal*.
642 2025:ehaf448.
- 643 9. Elias P, Poterucha TJ, Rajaram V, Moller LM, Rodriguez V, Bhave S, Hahn RT, Tison G,
644 Abreau SA, Barrios J. Deep learning electrocardiographic analysis for detection of left-
645 sided valvular heart disease. *Journal of the American College of Cardiology*. 2022;80:613-
646 626.
- 647 10. Poterucha TJ, Jing L, Ricart RP, Adjei-Mosi M, Finer J, Hartzel D, Kelsey C, Long A,
648 Rocha D, Ruhl JA. Detecting structural heart disease from electrocardiograms using AI.
649 *Nature*. 2025:1-10.
- 650 11. Kaplan J, McCandlish S, Henighan T, Brown TB, Chess B, Child R, Gray S, Radford A,
651 Wu J, Amodei D. Scaling laws for neural language models. *arXiv* 2020. *arXiv preprint
652 arXiv:200108361*. 2001.
- 653 12. Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, Sastry G, Askell A,
654 Mishkin P, Clark J. Learning transferable visual models from natural language supervision.
655 Paper/Poster presented at: International conference on machine learning; 2021;
- 656 13. Liu C, Wan Z, Ouyang C, Shah A, Bai W, Arcucci R. Zero-shot ecg classification with
657 multimodal learning and test-time clinical knowledge enhancement. *arXiv preprint
658 arXiv:240306659*. 2024.
- 659 14. Liu X, Cruz Rivera S, Moher D, Calvert MJ, Denniston AK, Chan A-W, Darzi A, Holmes
660 C, Yau C, Ashrafian H, et al. Reporting guidelines for clinical trial reports for interventions

- 661 involving artificial intelligence: the CONSORT-AI extension. *Nature Medicine*.
662 2020;26:1364-1374. doi: 10.1038/s41591-020-1034-x
- 663 15. Wong C-K, Lau YM, Lui HW, Chan WF, San WC, Zhou M, Cheng Y, Huang D, Lai WH,
664 Lau YM. Automatic detection of cardiac conditions from photos of electrocardiogram
665 captured by smartphones. *Heart*. 2024;110:1074-1082.
- 666 16. Jocher G, Qiu J, Chaurasia A. Ultralytics yolo11. *GitHub repository*. 2024.
- 667 17. Bradski G. The opencv library. *Dr Dobb's Journal: Software Tools for the Professional
668 Programmer*. 2000;25:120-123.
- 669 18. Yang A, Li A, Yang B, Zhang B, Hui B, Zheng B, Yu B, Gao C, Huang C, Lv C. Qwen3
670 technical report. *arXiv preprint arXiv:250509388*. 2025.
- 671 19. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. Paper/Poster
672 presented at: Proceedings of the IEEE conference on computer vision and pattern
673 recognition; 2016;
- 674 20. Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, Dehghani
675 M, Minderer M, Heigold G, Gelly S. An image is worth 16x16 words: Transformers for
676 image recognition at scale. *arXiv preprint arXiv:201011929*. 2020.
- 677 21. Jin Q, Kim W, Chen Q, Comeau DC, Yeganova L, Wilbur WJ, Lu Z. Medcpt: Contrastive
678 pre-trained transformers with large-scale pubmed search logs for zero-shot biomedical
679 information retrieval. *Bioinformatics*. 2023;39:btad651.
- 680 22. Gao T, Yao X, Chen D. Simcse: Simple contrastive learning of sentence embeddings. *arXiv
681 preprint arXiv:210408821*. 2021.
- 682 23. Gow B, Pollard T, Nathanson LA, Johnson A, Moody B, Fernandes C, Greenbaum N,
683 Waks JW, Eslami P, Carbonati T. MIMIC-IV-ECG: Diagnostic Electrocardiogram
684 Matched Subset. *Type: dataset*. 2023;6:13-14.
- 685 24. Lüscher TF, Wenzl FA, D'Ascenzo F, Friedman PA, Antoniades C. Artificial intelligence
686 in cardiovascular medicine: clinical applications. *European heart journal*. 2024;45:4291-
687 4304.
- 688 25. Quer G, Topol EJ. The potential for large language models to transform cardiovascular
689 medicine. *The Lancet Digital Health*. 2024;6:e767-e771.
- 690 26. Yu H, Guo P, Sano A. Ecg semantic integrator (esi): A foundation ecg model pretrained
691 with llm-enhanced cardiological text. *arXiv preprint arXiv:240519366*. 2024.
- 692 27. Zhou Y, Zhang P, Song M, Zheng A, Lu Y, Liu Z, Chen Y, Xi Z. Zodiac: A Cardiologist-
693 Level LLM Framework for Multi-Agent Diagnostics. *arXiv preprint arXiv:241002026*.
694 2024.
- 695 28. Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-cam: Visual
696 explanations from deep networks via gradient-based localization. Paper/Poster presented
697 at: Proceedings of the IEEE international conference on computer vision; 2017;
- 698 29. Thygesen K, Alpert JS, Jaffe AS, Chaitman BR, Bax JJ, Morrow DA, White HD, Infarction
699 EGobotJESoCACoCAHAWHFTFftUDoM. Fourth universal definition of myocardial
700 infarction (2018). *Journal of the American college of cardiology*. 2018;72:2231-2264.
- 701 30. Baman TS, Lange DC, Ilg KJ, Gupta SK, Liu T-Y, Alguire C, Armstrong W, Good E,
702 Chugh A, Jongnarangsin K. Relationship between burden of premature ventricular
703 complexes and left ventricular function. *Heart rhythm*. 2010;7:865-869.
- 704 31. Dukes JW, Dewland TA, Vittinghoff E, Mandyam MC, Heckbert SR, Siscovick DS, Stein
705 PK, Psaty BM, Sotoodehnia N, Gottdiener JS. Ventricular ectopy as a predictor of heart
706 failure and death. *Journal of the American College of Cardiology*. 2015;66:101-109.

- 707 32. Lai J, Tan H, Wang J, Ji L, Guo J, Han B, Shi Y, Feng Q, Yang W. Practical intelligent
708 diagnostic algorithm for wearable 12-lead ECG via self-supervised learning on large-scale
709 dataset. *Nature Communications*. 2023;14:3741.
- 710 33. Liu H, Zhao Z, She Q. Self-supervised ECG pre-training. *Biomedical Signal Processing*
711 and Control. 2021;70:103010.
- 712 34. Mehari T, Strodthoff N. Self-supervised representation learning from 12-lead ECG data.
713 *Computers in biology and medicine*. 2022;141:105114.
- 714 35. Bula K, Cmiel A, Sejud M, Sobczyk K, Ryszkiewicz S, Szydlo K, Wita M, Mizia-Stec K.
715 Electrocardiographic criteria for left ventricular hypertrophy in aortic valve stenosis:
716 Correlation with echocardiographic parameters. *Ann Noninvasive Electrocardiol*.
717 2019;24:e12645. doi: 10.1111/anec.12645
- 718 36. Shah AS, Chin CW, Vassiliou V, Cowell SJ, Doris M, Kwok TC, Semple S, Zamvar V,
719 White AC, McKillop G, et al. Left ventricular hypertrophy with strain and aortic stenosis.
720 *Circulation*. 2014;130:1607-1616. doi: 10.1161/CIRCULATIONAHA.114.011085
- 721 37. Acikgoz E, Yaman B, Acikgoz SK, Topal S, Tavil Y, Boyaci NB. Fragmented QRS can
722 predict severity of aortic stenosis. *Ann Noninvasive Electrocardiol*. 2015;20:37-42. doi:
723 10.1111/anec.12175
- 724 38. Vranic, II. Electrocardiographic appearance of aortic stenosis before and after aortic valve
725 replacement. *Ann Noninvasive Electrocardiol*. 2017;22. doi: 10.1111/anec.12457
- 726 39. Turhan H, Yetkin E, Atak R, Altinok T, Senen K, Ileri M, Sasmaz H, Cehreli S, Kutuk E.
727 Increased p-wave duration and p-wave dispersion in patients with aortic stenosis. *Ann
728 Noninvasive Electrocardiol*. 2003;8:18-21. doi: 10.1046/j.1542-474x.2003.08104.x
- 729 40. Dursun H, Tanriverdi Z, Colluoglu T, Kaya D. Effect of transcatheter aortic valve
730 replacement on P-wave duration, P-wave dispersion and left atrial size. *J Geriatr Cardiol*.
731 2015;12:613-617. doi: 10.11909/j.issn.1671-5411.2015.06.016
- 732 41. Yuan C, Yang Y, Yang Y, Cheng Z. DATE: Dynamic Absolute Time Enhancement for
733 Long Video Understanding. *arXiv preprint arXiv:250909263*. 2025.
- 734 42. Zhan Z, Wu Y, Gong Y, Meng Z, Kong Z, Yang C, Yuan G, Zhao P, Niu W, Wang Y. Fast
735 and memory-efficient video diffusion using streamlined inference. *Advances in Neural
736 Information Processing Systems*. 2024;37:13660-13684.
- 737