

# ProtoEFNet: Dynamic Prototype Learning for Inherently Interpretable Ejection Fraction Estimation in Echocardiography

Yeganeh Ghamary<sup>1</sup>, Victoria Wu<sup>2</sup>, Hooman Vaseli<sup>2</sup>, Christina Luong<sup>3</sup>, Teresa Tsang<sup>3</sup>, Siavash A. Bigdeli<sup>1</sup>, Purang Abolmaesumi<sup>2</sup>

<sup>1</sup> Department of Applied Mathematics and Computer Science, Technical University of Denmark, Kongens Lyngby, Denmark

<sup>2</sup> Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, BC, Canada

<sup>3</sup> Vancouver General Hospital, Vancouver, BC, Canada  
`{s194258,sarbi}@dtu.dk,{victoriawu,hoomanv,purang}@ece.ubc.ca` \*\*

**Abstract.** Ejection fraction (EF) is a crucial metric for assessing cardiac function and diagnosing conditions such as heart failure. Traditionally, EF estimation requires manual tracing and domain expertise, making the process time-consuming and subject to inter-observer variability. Most current deep learning methods for EF prediction are black-box models with limited transparency, which reduces clinical trust. Some post-hoc explainability methods have been proposed to interpret the decision-making process after the prediction is made. However, these explanations do not guide the model’s internal reasoning and therefore offer limited reliability in clinical applications. To address this, we introduce ProtoEFNet, a novel video-based prototype-learning model for continuous EF regression. The model learns dynamic spatio-temporal prototypes that capture clinically meaningful cardiac motion patterns. Additionally, the proposed Prototype Angular Separation (PAS) loss enforces discriminative representations across the continuous EF spectrum. Our experiments on the Echonet-Dynamic dataset show that ProtoEFNet can achieve accuracy on par with its non-interpretable counterpart while providing clinically relevant insight. The ablation study shows the proposed loss boosts the performance with a 2% increase in F1 score from  $77.67 \pm 2.68$  to  $79.64 \pm 2.10$ . Our source code is available at: <https://github.com/DeepRCL/ProtoEF>.

**Keywords:** Ultrasound · Echocardiography · Ejection Fraction · Regression · Explainable AI · Prototypical Neural Networks

\*\* Y. Ghamary and V. Wu are joint first authors.

T. Tsang, S. A. Bigdeli and P. Abolmaesumi are joint senior authors.

Work done during Y. Ghamary’s external stay at UBC.

## 1 Introduction

Heart failure is a major global health issue that requires accurate and timely assessment for effective management. A key measure of cardiac function is ejection fraction (EF), which quantifies the percentage of blood pumped from the left ventricle with each heartbeat. EF plays a central role in diagnosing heart failure, guiding treatment, and predicting outcomes [7,13]. It is typically measured using echocardiography, where clinicians manually trace the left ventricle at two key points in the cardiac cycle: end-diastole (ED) and end-systole (ES) to estimate volume changes [2]. However, this process is highly operator-dependent, with inter-observer variation ranging from 7.6 to 13.9 percent, and requires significant expertise [19]. Automating ejection fraction estimation through artificial intelligence can enhance consistency, reduce clinician workload, and support large-scale screening.

Despite progress in automated EF estimation [9,12,14,15,16,17,19,20], current methods face key challenges. Many rely on black-box deep learning models, offering limited transparency and reducing clinical trust. Explainability in medical AI is crucial, as clinicians require interpretable reasoning behind model predictions to facilitate adoption in medical settings. However, the existing explainable approaches use post-hoc techniques, such as attention weights [1,16], or gradients [21,23], that interpret the decision-making process after the prediction is made. These explanations are often inconsistent and do not inform the model’s internal reasoning, limiting their reliability in clinical applications.

To address these challenges, ante-hoc XAI methods have been introduced, embedding explainability directly into model architectures. Prototype-based models [3,10,11,25] are key examples, where networks learn class-specific prototypes that capture discriminative patterns explicitly used for classification. Unlike post-hoc methods that rely on abstract signals such as attention weights or gradients, which can be diffused, multi-layered, and difficult to interpret [3,22], prototypes offer more explicit, case-based explanations by grounding decisions in example-like visual features, mimicking a clinician’s decision process. However, adapting prototype learning to ejection fraction estimation poses specific challenges. EF is a continuous value, requiring a regression-based approach, and unlike other assessments that use static images, echocardiographic assessment depends on capturing spatio-temporal information across frames.

Some prior work has extended prototype-based models to regression [5,6] and video classification [25]. InsightRNet [6] applies prototype learning to regression, but uses discrete ordinal labels rather than truly continuous targets. Furthermore, [5,6] rely on fixed-size patch-based prototypes, which are insufficient for capturing the complex and spatially diverse visual features of ejection fraction (EF), particularly the broad and varying motion of the LV wall. We introduce ProtoEFNet (Fig. 1), the first video-based prototype model for continuous regression and the first inherently interpretable approach to EF estimation. Our key contributions are: (1) learning dynamic spatio-temporal prototypes that capture clinically relevant motion patterns; (2) proposing a prototype angular separation (PAS) loss to enforce discriminative representations

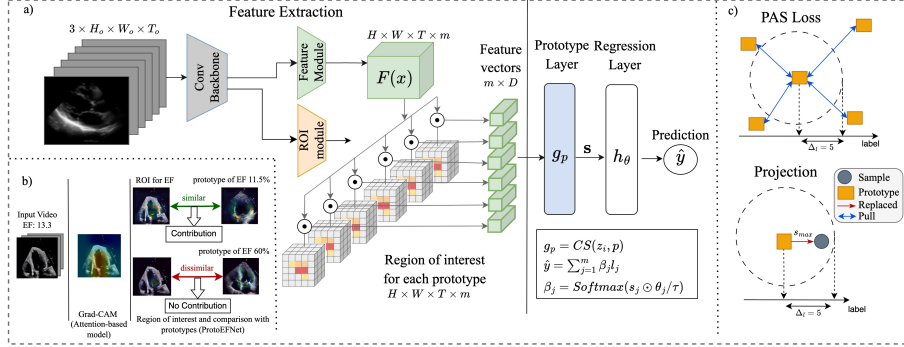


Fig. 1: (a) An overview of the architecture of ProtoEFNet. The feature extractor uses an ROI module to focus on clinically relevant spatio-temporal regions, and the final prediction is the weighted sum of the prototype labels, (b) Grad-CAM attention of CoReEcho [14], and the decision process of ProtoEFNet. The activation maps on the input data and prototypes show where the model "looks at" when calculating the cosine similarity, (c) The Prototype Angular Separation (PAS) loss increases separation between prototypes with different EF ranges, and prototypes are projected to/replaced with the closest training sample within its EF range.

across the continuous EF spectrum; (3) achieving state-of-the-art explainability on the EchoNet-Dynamic dataset, and accuracy on par with leading black-box models; (4) demonstrating through visualizations that ProtoEFNet uniquely attends to essential cardiac features—such as LV wall motion and reduced mitral valve movement in  $EF < 40\%$  cases—whereas leading black-box methods produce diffuse or clinically irrelevant attention.

## 2 Methods

**Problem Definition** We define a fixed number of  $m$  learnable prototype vectors  $\mathcal{P} = \{p_1, p_2, p_3, \dots, p_m\}$ , each associated with a continuous EF label  $\mathbf{l} = \{l_1, \dots, l_m\}$  and an *importance score*  $\theta = \{\theta_1, \dots, \theta_m\}$  that is learned during training. Based on the *similarity scores* and the learned importance of each prototype, we generate a score sheet with scores indicating *prototype contributions*. Unlike in conventional classification tasks, where multiple prototypes are associated with each discrete class and scores are calculated at the class level [3], our approach assigns probabilities to individual prototypes, capturing the continuous nature of the label space.

**Feature Extractor** has a similar architecture to [25], consisting of a pre-trained R(2+1)D-18 backbone, a feature module  $F(\cdot)$ , and a Region of Interest (ROI) module  $M_{pk}(\cdot)$ . Given an input video  $x \in \mathbb{R}^{H_o \times W_o \times T_o \times 3}$  with  $T_o$  frames, the

learned features are  $F(x) \in \mathbb{R}^{H \times W \times T \times D}$ , and the ROI module generates P occurrence maps  $M_{p_k}(x) \in \mathbb{R}^{H \times W \times T}$  highlighting which regions of  $F(x)$  are relevant when compared to  $p_k$ . Occurrence maps highlight the regions in space and time where it is likely to observe relevant features for EF prediction. We perform a weighted average pooling of the spatiotemporal features using the learned occurrence maps as weights.

**Prototype Layer.** We use cosine similarity (CS) score as a similarity function between the features  $f_{p_k}(x)$  and prototypes  $p_k$ , which are both D-dimensional vectors  $s_k = g_{p_k}(f_{p_k}(x)) = CS(f_{p_k}(x), p_k)$ .

**Regression Layer** is a linear layer with m weights denoted by  $\theta$ . Unlike class-based prototype models, the final prediction is calculated as a weighted average of the prototype labels:

$$\hat{y} = \sum_{k=1}^m \beta_k \cdot l_k, \quad \beta_k(x) = \frac{e^{(s_k * \theta_k)/\tau}}{\sum_{i=1}^m e^{(s_i * \theta_i)/\tau}}, \quad (1)$$

where  $\beta_k(\cdot)$  indicates the *contribution* of each prototype  $p_k$  in the final prediction and is calculated using the softmax function of scaled similarity scores. To encourage *sparse* explanations, we use a small temperature parameter of  $\tau = 0.2$ . This ensures that dissimilar prototypes have 0 contribution to the final prediction (see Figure 1(b)). The overview of ProtoEFNet can be seen in Figure 1(a).

**Training Algorithm.** The feature extractor, ROI module, regression layer, and prototype vectors are trained jointly. In the last epoch, the prototypes (and their labels) are *projected* (replaced) onto the closest training feature (see Figure 1(c)). This allows us to visualize the prototypes, enhancing human-level explainability. We adapt the notion of class-based projection to the regression task using a threshold of  $\Delta_l$ . It is mathematically defined as:

$$p_j \leftarrow \operatorname{argmax}_z CS(z, p_j), \quad z = f_{p_j}(x_i), |y_i - l_j| < \Delta_l. \quad (2)$$

**Latent Space Regularization.** To learn representative prototypes, the latent space is learned using two distance-based losses. A regression-based cluster loss and prototype sample distance loss (PSD) [6]. Cluster loss creates clusters around prototypes by pulling the samples with similar EF labels using a threshold of  $\Delta_l$ , and PSD loss ensures that there is at least one sample close to each prototype:

$$\mathcal{L}_{Clst} = -\frac{1}{n} \sum_{i=1}^n \operatorname{kmax}_{p \in \mathcal{P}^c}(s_p(i)), \quad |y_i - l_j| < \Delta_l, \quad (3)$$

$$\mathcal{L}_{PSD} = -\frac{1}{m} \sum_{j=1}^m \log(1 - \min_{i \in [1, n]} \frac{d_{i,j}}{d_{max}}), \quad (4)$$

where  $n$  is the batch size,  $kmax$  is the maximum operation selecting the  $k$  closest prototypes to each sample, and  $s_p(i)$  is the cosine similarity score between the prototype  $p$  and sample  $i$ .  $\mathcal{P}^c$  visualizes the set of prototypes in the vicinity of the sample in the label space.  $d_{max}$  is the maximum possible distance in the latent space, which is 2 for the cosine distance score, and  $d_{i,j}$  is the cosine distance between sample  $i$  and prototype  $j$ . Prototype models [3,10] have been effective in classification by using separation losses to enforce clear boundaries between class-specific prototypes. However, in regression tasks, where outputs are continuous and ordered, preserving the ordinality and smooth transitions between clusters is essential. In the meantime, prototype vectors representing different EF regions should be further apart than those prototypes belonging to the same EF range. This ensures ordinality in the embedding space and that prototypes are semantically meaningful. However, in practice, we observed that even the most distant prototypes with EF labels of 15% and 86% have similar embedding vectors with a cosine similarity of around 0.7, suggesting that the embeddings lack meaningful distinction. To address this, and inspired by [11], we propose a novel prototype angular separation loss (PAS) that promotes greater inter-prototype distinction while respecting the continuity between the cluster of the data points required for regression tasks. The PAS loss addresses this by pulling prototypes from different EF regions further apart (see Figure 1(c)). To measure the distance between prototypes, we utilize angular similarity—a monotonic transformation of cosine similarity that ranges from 0 to 1. The PAS loss can be defined as follows:

$$\mathcal{L}_{PAS} = -\frac{1}{m} \sum_{i=1}^m \left[ \frac{1}{|\mathcal{P}^c|} \sum_{j \in \mathcal{P}^c} \log(1 - AS(p_i, p_j)) \right], \quad \mathcal{P}^c = \{p_j \mid |l_i - l_j| > \Delta_l\}, \quad (5)$$

where  $AS$  indicates angular similarity defined as:

$$AS(\mathbf{p}_i, \mathbf{p}_j) = 1 - \frac{1}{\pi} \arccos(CS(\mathbf{p}_i, \mathbf{p}_j)). \quad (6)$$

$\mathcal{P}^c$  is the set of prototypes that have EF label larger than  $\Delta_l$  of the given prototype. Inside the bracket is the average of the log of the pairwise angular similarity. We chose the average due to its superiority over the max function.

**Occurrence Map Regularizer.** To further emphasize disease-specific areas, we incorporate the LV segmentation mask into the L1 regularization, penalizing activations outside the LV and reducing attention to irrelevant regions like the background. The overall cost function is defined below, where  $\lambda$  represents weight of each loss term:

$$\mathcal{L} = \lambda_{MSE} \mathcal{L}_{MSE} + \lambda_{Clst} \mathcal{L}_{Clst} + \lambda_{PSD} \mathcal{L}_{PSD} + \lambda_{PAS} \mathcal{L}_{PAS} + \lambda_{Occur} \mathcal{L}_{Occur}. \quad (7)$$

### 3 Experiment

**Dataset and Implementation.** EchoNet-Dynamic [18] is the largest public echocardiogram dataset, containing 10,036 apical four-chamber (AP4) videos (112×112 resolution) with varying lengths, each labeled with a single EF value, LV segmentation, and end-systolic/diastolic frame indices. We follow the train-val-test split from [19] and frame sampling from [16]. We augment the data using random rotation. Label balance was achieved by oversampling the minority region ( $EF < 50\%$ ), as addressing imbalance was beyond the scope of this study. The EF value outside the range of  $[10\%, 90\%]$  is clinically uncommon, and the dataset does not include any samples outside this range. Derived from the hyperparameter selection, each video is sampled as a clip of 64 frames (sampling period of 1) with the initial frame index sampled uniformly.

All experiments were conducted using an NVIDIA A100 (40 GB) with PyTorch 2.0.1 and CUDA 12.8. Following hyperparameter tuning, we used 40 prototypes, a softmax temperature of  $\tau = 0.2$  for regression layer,  $\Delta_l = 5.0\%$ , and  $k = 3$  in  $\mathcal{L}_{Clst}$ . The prototype vectors were initialized randomly. Regression layer weights  $\theta_h$  were initialized to 1, assigning equal importance to prototypes, and prototype labels were uniformly set from 10% to 90%. This is to ensure that prototypes represent all regions of the label space, including low-density regions. We used Adam optimizer with a learning rate of  $1e-4$  for the backbone and regression layer,  $1e-3$  for feature and ROI modules and  $3e-3$  for prototype vectors. We use pretrained weights of [8] and fine-tuned jointly for 30 epochs with a batch size of 16.

**Comparison with SOTA models.** In Table 1, we compare ProtoEFNet with SOTA methods on EchoNet-Dynamic. We excluded [5,6], as adapting them to the dataset required major changes that altered their core architecture. We furthermore report the F1 score for the task of indicating whether EF values are lower than 40%, which is a strong indicator of heart failure [19]. ProtoEFNet outperforms most of the SOTA models, including the post-hoc explainable model [16] and EchoCoTr with same frame size. The performance gap between ProtoEFNet and CoReEcho [14] could potentially be narrowed with further training, and more extensive hyperparameter tuning. A key advantage of ProtoEFNet is its inherent explainability, offering transparent insights into its predictions.

**Explainability Analysis.** Figure 1(b) illustrates ProtoEFNet’s decision process for a video sample, a more detailed description of the activation map visualisation can be found in [10]. ProtoEFNet is *inherently interpretable* and *transparent*, with predictions directly linked to prototype contribution and similarity scores. Its explanations are *sparse* and *faithful*, relying only on prototypes with EF values close to the true label—prototypes with distant EF values (e.g., 60%) do not influence the prediction for the sample with EF of 13%. Activation maps on the input video demonstrate both spatial (anatomical localization) and temporal (motion) explainability, highlighting *clinically relevant* features such

Table 1: The quantitative results on the EchoNet-Dynamic test set. **P** indicates the frame sampling period. Most approaches average clip-level predictions over the entire video (**EV**), while others use a single heartbeat (**SH**). Results marked with (\*) denote reproduced performance, while (<sup>1</sup>) indicates the SOTA method with the frame sampling similar to ours.

| Method                        | Frames | P | Clips   | R2 $\uparrow$ | MAE $\downarrow$ | RMSE $\downarrow$ | $F1_{<40\%} \uparrow$ | Explainable |
|-------------------------------|--------|---|---------|---------------|------------------|-------------------|-----------------------|-------------|
| Reynaud et al. [20]           | 128    | 1 | SH      | 52            | 5.59             | 8.38              | NA                    | X           |
| Bayesian [9]                  | 32     | 1 | SH      | 75            | 4.46             | NA                | 77                    | X           |
| EchoCoTr [17] <sup>1</sup>    | 36     | 2 | EV      | 79            | 4.18             | 5.59              | NA                    | X           |
| EchoCoTr [17]                 | 36     | 4 | EV      | 81            | 3.98             | 5.34              | NA                    | X           |
| CoReEcho [14]                 | 36     | 1 | 3 Clips | <b>82</b>     | <b>3.90</b>      | <b>5.13</b>       | <b>80</b>             | X           |
| GEMTrans [15] <sup>1</sup>    | 64     | 1 | EV      | 79            | 4.15             | NA                | NA                    | X           |
| Resnet2+1D [19] <sup>1*</sup> | 64     | 1 | EV      | 80            | 4.10             | 5.47              | 78                    | X           |
| EchoGNN [16] <sup>1</sup>     | 64     | 1 | EV      | 76            | 4.45             | NA                | <b>78</b>             | post-hoc    |
| ProtoEFNet (ours)             | 64     | 1 | EV      | <b>80</b>     | <b>4.07</b>      | <b>5.47</b>       | <b>78</b>             | ✓           |

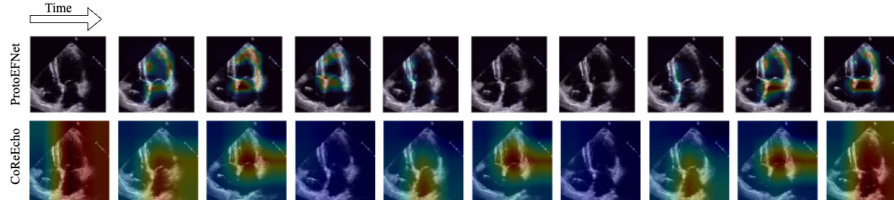


Fig. 2: Grad-CAM [21] of CoReEcho (bottom row) and the activation map of ProtoEFNet (top row). ProtoEFNet is localised on LV wall motion and mitral valve movements during systole (contraction).

as reduced LV wall motion and mitral valve movement during systole. Prototypical features shown as activation maps on top of prototypical cases show distinct EF-related patterns: the 60% EF prototype exhibits healthy LV and MV motion, while the 11% EF prototype shows reduced LV motion and a thin LV wall, both clinically meaningful. In Figure 2, CoReEcho’s spatio-temporal attention (Grad-CAM [21]) is compared to ProtoEFNet’s activation map. ProtoEFNet demonstrates superior localization, focusing on key structures like the LV wall and mitral/aortic valves, while CoReEcho and EchoCoTr highlight non-specific or clinically irrelevant regions, such as background (see [14]). ProtoEFNet also learns the periodic nature of echo videos and correctly aligns the spatio-temporal features of the input clip to those of the prototype, seen in Figure 3.

**Ablation Study.** Figure 4 shows an ablation study of different loss components. Including all loss components yields the best regression performance based on MAE and F1 scores. We evaluate prototype quality using Sparsity and Diversity metrics [6] scaled by the number of prototypes: effective explanations rely on the contribution of a few prototypes (*low Sparsity*), but different predictions

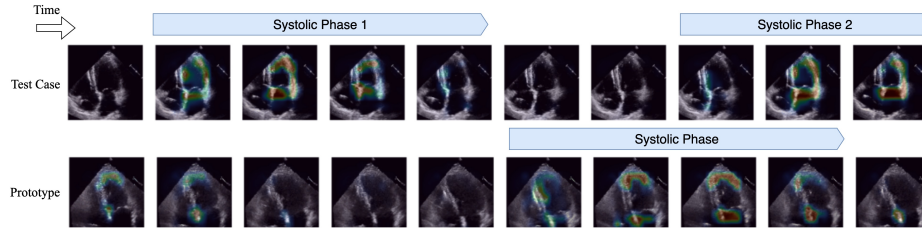


Fig. 3: Activation maps of a test case and the top contributing prototype. The model captures periodic patterns and aligns spatio-temporal features. In this example, it assigns a high similarity score as it "looks at" the LV wall during systole in the input and identifies that it "looks like" the LV wall of the prototype during the same phase.

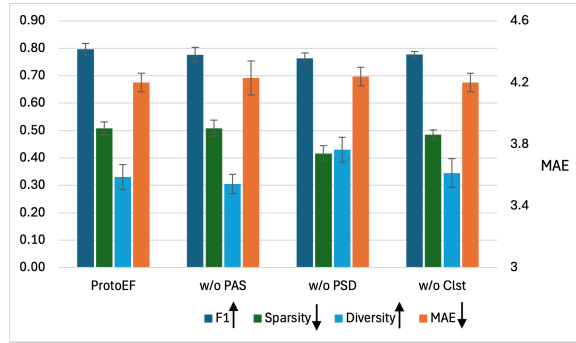


Fig. 4: Ablation Study of different loss components on validation set. Standard deviation is calculated across 5 repetitions of each experiment.

rely on different prototypes (*high Diversity*). To demonstrate the effect of these components in the embedding space, we visualize 2D PCA plot of the prototypes and the 100 closest validation features in Figure 5. Removing  $\mathcal{L}_{PSD}$  improves diversity and sparsity but degrades regression performance, as reflected by F1 and MAE scores. Moreover, some prototypes become outliers with no nearby training samples (see Figure 5). Without  $\mathcal{L}_{Clst}$ , the samples and prototypes are scattered in the embedding space without any clear ordinality.  $\mathcal{L}_{PAS}$  decreases MAE from  $4.23 \pm 0.11$  to  $4.20 \pm 0.06$ , and increases F1 score from  $77.67 \pm 2.68$  to  $79.64 \pm 2.10$ , while producing more distinct prototypes and ordinally structured clusters in the embedding space.

## 4 Conclusion

We proposed ProtoEFNet, the first prototype-based model for video-based continuous EF regression. ProtoEFNet has superior performance to the *post-hoc* explainable model [16] and a performance on par with the black-box models,



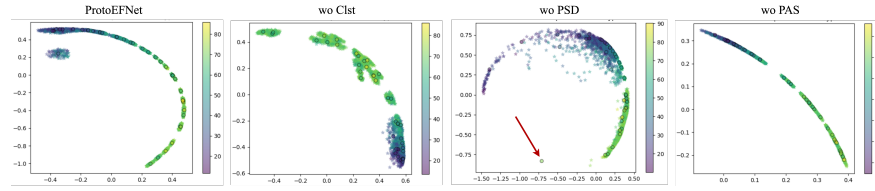


Fig. 5: The PCA plots of the prototypes (circles) and the top 100 closest latent features of the validation set to each prototype (stars). The colors indicate the ground truth EF values.

while offering inherent interpretability, transparency, and clinically meaningful explanations. Qualitative analysis shows superior focus on key cardiac structures, and ablation results confirm the effectiveness of the proposed Prototype Angular Separation loss. Future work will address prototype learning for uncommon EF ranges (minority region).

**Acknowledgements.** This work was supported in part by the Technical University of Denmark’s Travel Grant, Marie og Anders Manssons Memorial Grant, the Canadian Institutes of Health Research (CIHR), and the Natural Sciences and Engineering Research Council of Canada (NSERC). Resources were provided through DTU Computing Center at Technical University of Denmark [4] and Advanced Research Computing at the University of British Columbia [24].

## References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473 (2014)
2. Bamira, D., Picard, M.: Imaging: Echocardiology—assessment of cardiac structure and function (2018)
3. Chen, C., Li, O., Tao, D., Barnett, A., Rudin, C., Su, J.K.: This looks like that: deep learning for interpretable image recognition. *Advances in neural information processing systems* **32** (2019)
4. DTU Computing Center: DTU Computing Center resources (2024). <https://doi.org/10.48714/DTU.HPC.0001>, <https://doi.org/10.48714/DTU.HPC.0001>
5. Hesse, L.S., Dinsdale, N.K., Namburete, A.I.: Prototype learning for explainable brain age prediction. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. pp. 7903–7913 (2024)
6. Hesse, L.S., Namburete, A.I.: Insightr-net: interpretable neural network for regression using similarity-based comparisons to prototypical examples. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 502–511. Springer (2022)
7. Huang, H., Nijjar, P.S., Misialek, J.R., Blaes, A., Derrico, N.P., Kazmirczak, F., Klem, I., Farzaneh-Far, A., Shenoy, C.: Accuracy of left ventricular ejection fraction by contemporary multiple gated acquisition scanning in patients with cancer: comparison with cardiovascular magnetic resonance. *Journal of Cardiovascular Magnetic Resonance* **19**(1), 34 (2016)

8. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
9. Kazemi Esfeh, M.M., Luong, C., Behnami, D., Tsang, T., Abolmaesumi, P.: A deep bayesian video analysis framework: towards a more robust estimation of ejection fraction. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 582–590. Springer (2020)
10. Kim, E., Kim, S., Seo, M., Yoon, S.: Xprotonet: diagnosis in chest radiography with global and local explanations. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 15719–15728 (2021)
11. Kraft, S., Broelemann, K., Theissler, A., Kasneci, G., Esslingen am Neckar, G., AG, S.H.: Sparrow: Semantically coherent prototypes for image classification. In: Bmvc. p. 186 (2021)
12. Lai, S., Zhao, M., Zhao, Z., Chang, S., Yuan, X., Liu, H., Zhang, Q., Meng, G.: Echomen: Combating data imbalance in ejection fraction regression via multi-expert network. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 624–633. Springer (2024)
13. Loefer, L.R., Rosamond, W.D., Chang, P.P., Folsom, A.R., Chambless, L.E.: Heart failure incidence and survival (from the atherosclerosis risk in communities study). *The American journal of cardiology* **101**(7), 1016–1022 (2008)
14. Maani, F.A., Saeed, N., Matsun, A., Yaqub, M.: Coreecho: Continuous representation learning for 2d+ time echocardiography analysis. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 591–601. Springer (2024)
15. Mokhtari, M., Ahmadi, N., Tsang, T.S., Abolmaesumi, P., Liao, R.: Gemtrans: A general, echocardiography-based, multi-level transformer framework for cardiovascular diagnosis. In: International Workshop on Machine Learning in Medical Imaging. pp. 1–10. Springer (2023)
16. Mokhtari, M., Tsang, T., Abolmaesumi, P., Liao, R.: Echognn: explainable ejection fraction estimation with graph neural networks. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 360–369. Springer (2022)
17. Muhtaseb, R., Yaqub, M.: Echocotr: Estimation of the left ventricular ejection fraction from spatiotemporal echocardiography. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 370–379. Springer (2022)
18. Ouyang, D., He, B., Ghorbani, A., Lungren, M.P., Ashley, E.A., Liang, D.H., Zou, J.Y.: Echonet-dynamic: a large new cardiac motion video data resource for medical machine learning. In: NeurIPS ML4H Workshop. pp. 1–11 (2019)
19. Ouyang, D., He, B., Ghorbani, A., Yuan, N., Ebinger, J., Langlotz, C.P., Heidenreich, P.A., Harrington, R.A., Liang, D.H., Ashley, E.A., et al.: Video-based ai for beat-to-beat assessment of cardiac function. *Nature* **580**(7802), 252–256 (2020)
20. Reynaud, H., Vlontzos, A., Hou, B., Beqiri, A., Leeson, P., Kainz, B.: Ultrasound video transformers for cardiac ejection fraction estimation. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VI 24. pp. 495–505. Springer (2021)
21. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE international conference on computer vision* pp. 618–626 (2017)

22. Serrano, S., Smith, N.A.: Is attention interpretable? arXiv preprint arXiv:1906.03731 (2019)
23. Smilkov, D., Thorat, N., Kim, B., Viegas, F., Wattenberg, M.: Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825 (2017)
24. UBC Advanced Research Computing: UBC ARC sockeye (2019)
25. Vaseli, H., Gu, A.N., Ahmadi Amiri, S.N., Tsang, M.Y., Fung, A., Kondori, N., Saadat, A., Abolmaesumi, P., Tsang, T.S.: Protoasnet: Dynamic prototypes for inherently interpretable and uncertainty-aware aortic stenosis classification in echocardiography. In: International conference on medical image computing and computer-assisted intervention. pp. 368–378. Springer (2023)