EchoGLAD: Hierarchical Graph Neural Networks for Left Ventricle Landmark Detection on Echocardiograms

Masoud Mokhtari^{1[0000-0001-9471-5573]}, Mobina Mahdavi¹, Hooman Vaseli^{1[0000-0002-8259-9488]}, Christina Luong², Purang Abolmaesumi^{1*}, Teresa S. M. Tsang², and Renjie Liao^{1*}

¹ Electrical and Computer Engineering, University of British Columbia, Vancouver, BC, Canada

{masoud, mobina, hoomanv, purang, rjliao}@ece.ubc.ca
² Vancouver General Hospital, Vancouver, BC, Canada
 {christina.luong, t.tsang}@ubc.ca

Abstract. The functional assessment of the left ventricle chamber of the heart requires detecting four landmark locations and measuring the internal dimension of the left ventricle and the approximate mass of the surrounding muscle. The key challenge of automating this task with machine learning is the sparsity of clinical labels, i.e., only a few landmark pixels in a high-dimensional image are annotated, leading many prior works to heavily rely on isotropic label smoothing. However, such a label smoothing strategy ignores the anatomical information of the image and induces some bias. To address this challenge, we introduce an echocardiogram-based, hierarchical graph neural network (GNN) for left ventricle landmark detection (EchoGLAD). Our main contributions are: 1) a hierarchical graph representation learning framework for multi-resolution landmark detection via GNNs; 2) induced hierarchical supervision at different levels of granularity using a multi-level loss. We evaluate our model on a public and a private dataset under the in-distribution (ID) and out-of-distribution (OOD) settings. For the ID setting, we achieve the state-of-the-art mean absolute errors (MAEs) of 1.46 mm and 1.86 mm on the two datasets. Our model also shows better OOD generalization than prior works with a testing MAE of 4.3 mm.

 ${\bf Keywords:} \ {\rm Graph \ Neural \ Networks} \cdot {\rm Landmark \ Detection} \cdot {\rm Ultrasound}.$

1 Introduction

Left Ventricular Hypertrophy (LVH), one of the leading predictors of adverse cardiovascular outcomes, is the condition where heart's mass abnormally increases secondary to anatomical changes in the Left Ventricle (LV) [10]. These anatomical changes include an increase in the septal and LV wall thickness, and

^{*} Co-Corresponding Authors



Fig. 1: (a) IVS, LVID and LVPW measurements visualized on a PLAX echo frame. (b) If the wall landmark labels are smoothed by an isotropic Gaussian distribution, points along the visualized wall and ones perpendicular are penalized equally. Ideally, points along the walls must be penalized less.

the enlargement of the LV chamber. More specifically, Inter-Ventricular Septal (IVS), LV Posterior Wall (LVPW) and LV Internal Diameter (LVID) are assessed to investigate LVH and the risk of heart failure [21]. As shown in Figure 1 (a), four landmarks on a parasternal long axis (PLAX) echo frame can characterize IVS, LVPW and LVID, and allow cardiac function assessment. To automate this, machine learning-based (ML) landmark detection methods have gained traction.

It is difficult for such ML models to achieve high accuracy due to the sparsity of positive training signals (four or six) pertaining to the correct pixel locations. In an attempt to address this, previous works use 2D Gaussian distributions to smooth the ground truth landmarks of the LV [9,13,18]. However, as shown in Figure 1 (b), for LV landmark detection where landmarks are located at the wall boundaries (as illustrated by the dashed line), we argue that an isotropic Gaussian label smoothing approach confuses the model by being agnostic to the structural information of the echo frame and penalizing the model similarly whether the predictions are perpendicular or along the LV walls.

In this work, to address the challenge brought by sparse annotations and label smoothing, we propose a hierarchical framework based on Graph Neural Networks (GNNs) [25] to detect LV landmarks in ultrasound images. As shown in Figure 2, our framework learns useful representations on a hierarchical grid graph built from the input echo image and performs multi-level prediction tasks.

Our contributions are summarized below.

- We propose a novel GNN framework for LV landmark detection, performing message passing over hierarchical graphs constructed from an input echo;
- We introduce a hierarchical supervision that is automatically induced from sparse annotations to alleviate the issue of label smoothing;
- We evaluate our model on two LV landmark datasets and show that it not only achieves state-of-the-art mean absolute errors (MAEs) (1.46 mm and

1.86 mm across three LV measurements) but also outperforms other methods in out-of-distribution (OOD) testing (achieving 4.3 mm).



Fig. 2: Overview of our proposed model architecture. Hierarchical Feature Construction provides node features for the hierarchical graph representation of each echo frame where the nodes in the main graph correspond to pixels in the image, and nodes in the auxiliary graphs correspond to patches of different granularity in the image. Graph Neural Networks are used to process the hierarchical graph representation and produce node embeddings for the auxiliary graphs and the main graph. Multi-Layer Perceptrons (MLPs) are followed by a Sigmoid output function to map the node embeddings into land-mark heatmaps of different granularity over the input echo frame.

2 Related Work

Various convolution-based LV landmark detection works have been proposed. Sofka *et al.* [26] use Fully Convolutional Networks to generate prediction heatmaps followed by a center of mass layer to produce the coordinates of the landmark locations. Another work [18] uses a modified U-Net [24] model to produce a segmentation map followed by a focal loss to penalize pixel predictions in close proximity of the ground truth landmark locations modulated by a Gaussian distribution. Jafari *et al.* [13] use a similar U-Net model with Bayesian neural networks [8] to estimate the uncertainty in model predictions and reject samples that exhibit high uncertainties. Gilbert *et al.* [6] smooth ground truth labels by placing 2D Gaussian heatmaps around landmark locations at angles that are statistically obtained from training data. Lastly, Duffy *et al.* [4] use atrous convolutions [1] to make predictions for LVID, IVS and LVPW measurements.

Other related works focus on the detection of cephalometric landmarks from X-ray images. These works are highly transferable to the task of LV landmark detection as they must also detect a sparse number of landmarks. McCouat et

4 M. Mokhtari et al.

al. [20] is one of these works that abstains from using Gaussian label smoothing, but still relies on one-hot labels and treats landmark detection as a pixel-wise classification task. Chen *et al.* [2] is another cephalometric landmark detection work that creates a feature pyramid from the intermediate layers of a ResNet [11].

Our approach is different from prior works in that it aims to avoid the issue shown in Fig. 1 (b) and the sparse annotations problem by the introduction of simpler auxiliary tasks to guide the main pixel-level task, so that the ML model learns the location of the landmarks without relying on Gaussian label smoothing. It further improves the representation learning via efficient messagepassing [25,7] of GNNs among pixels and patches at different levels without having as high a computational complexity as transformers [3,19]. Lastly, while GNNs have never been applied to the task of LV landmark detection, they have been used for landmark detection in other domains. Li et al. [16] and Lin et al. [17] perform face landmark detection via modeling the landmarks with a graph and performing a cascaded regression of the locations. These methods, however, do not leverage hierarchical graphs and hierarchical supervision and instead rely on initial average landmark locations, which is not an applicable approach to echo, where the anatomy of the depicted heart can vary significantly. Additionally, Mokhtari et al. [22] use GNNs for the task of EF prediction from echo cine series. However, their work focuses on regression tasks.

3 Method

3.1 Problem Setup

We consider the following supervised setting for LV wall landmark detection. We have a dataset $D = \{X, Y\}$, where |D| = n is the number of $\{x^i, y^i\}$ pairs such that $x^i \in X$, $y^i \in Y$, and $i \in [1, n]$. Each $x^i \in \mathbb{R}^{H \times W}$ is an echo image of the heart, where H and W are height and width of the image, respectively, and each y^i is the set of four point coordinates $[(h_1^i, w_1^i), (h_2^i, w_2^i), (h_3^i, w_3^i), (h_4^i, w_4^i)]$ indicating the landmark locations in x^i . Our goal is to learn a function f: $\mathbb{R}^{H \times W} \mapsto \mathbb{R}^{4 \times 2}$ that predicts the four landmark coordinates for each input image. A figure in the supp. material further clarifies how the model generates landmark location heatmaps on different scales (Fig. 2).

3.2 Model Overview

As shown in Figure 2, each input echo frame is represented by a hierarchical grid graph where each sub-graph corresponds to the input echo frame at a different resolution. The model produces heatmaps over both the main pixel-level task as well as the coarse auxiliary tasks. While the pixel-level heatmap prediction is of main interest, we use a hierarchical multi-level loss approach where the model's prediction over auxiliary tasks is used during training to optimize the model through comparisons to coarser versions of the ground truth. The intuition behind such an approach is that the model learns nuances in the data by performing landmark detection on the easier auxiliary tasks and uses this established reasoning when performing the difficult pixel-level task.

3.3 Hierarchical Graph Construction

To learn representations that better capture the dependencies among pixels and patches, we introduce a hierarchical grid graph along with multi-level prediction tasks. As an example, the simplest task consists of a grid graph with only four nodes, where each node corresponds to four equally-sized patches in the original echo image. In the main task (the one that is at the bottom in Figure 2 and is the most difficult), the number of nodes is equal to the total number of pixels.

More formally, let us denote a graph as G = (V, E), where V is the set of nodes, and E is the set of edges in the graph such that if $v_i, v_j \in V$ and there is an edge from v_i to v_j , then $e_{i,j} \in E$. To build hierarchical task representations, for each image $x \in X$ and the ground truth $y \in Y$, K different auxiliary graphs $G_k(V_k, E_k)$ are constructed using the following steps for each $k \in [1, K]$:

- 1. $2^k \times 2^k = 4^k$ nodes are added to V_k to represent each patch in the image. Note that the larger values of k correspond to graphs of finer resolution, while the smaller values of k correspond to coarser graphs.
- 2. Grid-like, undirected edges are added such that $e_{m-1,q}, e_{m+1,q}, e_{m,q-1}, e_{m,q+1} \in E_k$ for each $m, q \in [1 \dots 2^k]$ if these neighbouring nodes exist in the graph (border nodes will not have four neighbouring nodes).
- 3. A patch feature embedding z_j^k , where $j \in [1 \dots 4^k]$ is generated and associated with that patch (node) $v_j \in V_k$. The patch feature construction technique is described in Section 3.4.
- 4. Binary node labels $\hat{y}_k \in \{0,1\}^{4^k \times 4}$ are generated such that $\hat{y}_{kj} = 1$ if at least one of the ground truth landmarks in y is contained in the patch associated with node $v_j \in V_k$. Note that for each auxiliary graph, four different one-hot labels are predicted, which correspond to each of the four landmarks required to characterize LV measurements.

The main graph, G_{main} , has a grid structure and contains $H \times W$ nodes regardless of the value of K, where each node corresponds to a pixel in the image. Additionally, to allow the model to propagate information across levels, we add inter-graph edges such that each node in a graph is connected to four nodes in the corresponding region in the next finer graph as depicted in Fig. 2.

3.4 Node Feature Construction

The graph representation described in Section 3.3 is not complete without proper node features, denoted by $z \in \mathbb{R}^{|V| \times d}$, characterizing patches or pixels of the image. To achieve this, the grey-scale image is initially expanded in the channel dimension using a CNN. The features are then fed into a U-Net where the decoder part is used to obtain node features such that deeper layer embeddings correspond to the node features for the finer graphs. This means that the main pixel-level graph would have the features of the last layer of the network. A figure clarifying node feature construction is provided in the supp. material (Fig. 1). 6 M. Mokhtari et al.

3.5 Hierarchical Message Passing

We now introduce how we perform message passing on our constructed hierarchical graph using GNNs to learn node representations for predicting landmarks.

The whole hierarchical graph created for each sample, *i.e.*, the main graph, auxiliary graphs, and cross-level edges, are collectively denoted as G^i , where $i \in [1, \ldots, n]$. Each G^i is fed into GNN layers followed by an MLP:

$$h_{\text{nodes}}^{l+1} = \text{ReLU}(\text{GNN}_l(G^i), h_{\text{nodes}}^l), \quad l \in [0, \dots, L]$$
(1)

$$h_{\rm out} = \sigma({\rm MLP}(h_{\rm nodes^{L+1}})), \tag{2}$$

where σ is the Sigmoid function, $h_{\text{nodes}}^l \in \mathbb{R}^{|V_{G^i}| \times d}$ is the set of d-dimensional embeddings for all nodes in the graph at layer l, and $h_{\text{out}} \in [0,1]^{|V_{G^i}| \times 4}$ is the four-channel prediction for each node with each channel corresponding to a heatmap for each of the pixel landmarks. The initial node features h_{nodes}^1 are set to the features z described in Sections 3.3 and 3.4. The coordinates $(x_{\text{out}}^p, y_{\text{out}}^p)$ for each landmark location $p \in [1, 2, 3, 4]$ are obtained by taking the expected value of individual heatmaps h_{out}^p along the x and y directions such that:

$$x_{\text{out}}^p = \sum_{s=1}^{|V_{G^i}|} \operatorname{softmax}(h_{\text{out}}^p)_s * \operatorname{loc}_x(s),$$
(3)

where similar operations are performed in the y direction for y_{out}^p . Here, we vectorize the 2D heatmap into a single vector and then feed it to the softmax. loc_x and loc_y return the x and y positions of a node in the image. It must be noted that unlike some prior works such as Duffy *et al.* [4] that use postprocessing steps such as imposing thresholds on the heatmap values, our work directly uses the output heatmaps to find the final predictions.

3.6 Training and Objective Functions

To train the network, we leverage two types of objective functions. 1) Weighted Binary Cross Entropy (BCE): Since the number of landmark locations is much smaller than non-landmark locations, we use a weighted BCE loss; 2) L2 regression of landmark coordinates: We add a regression objective which is the L2 loss between the predicted coordinates and the ground truth labels.

4 Experiments

4.1 Datasets

Internal Dataset: Our private dataset contains 29,867 PLAX echo frames, split in a patient-exclusive manner with 23824, 3004, and 3039 frames for training, validation, and testing, respectively. External Dataset: The public Unity Imaging Collaborative (UIC) [12] LV landmark dataset consists of a combination of 3822 end-systolic and end-diastolic PLAX echo frames acquired from

seven British echocardiography labs. The provided splits contain 1613, 298, and 1911 training, validation, and testing samples, respectively. For both datasets, we down-sample the frames to a fixed size of 224×224 .

4.2 Implementation Details

Our model creates K=7 auxiliary graphs. For the node features, the initial singlelayer CNN uses a kernel size of 3 and zero-padding to output features with a dimension of $224 \times 224 \times 4$ (C=4). The U-Net's encoder contains 7 layers with $128 \times 128, 64 \times 64, 32 \times 32, 16 \times 16, 8 \times 8, 4 \times 4$, and 2×2 spatial dimensions, and 8, 16, 32, 64, 128, 256, and 512 number of channels, respectively. Three Graph Convolutional Network (GCN)[15] layers (L=3) with a hidden node dimension of 128 are used. To optimize the model, we use the Adam optimizer [14] with an initial learning rate of 0.001, β of (0.9, 0.999) and a weight decay of 0.0001, and for the weighted BCE loss, we use a weight of 9000. The model is implemented using PyTorch [23] and Pytorch Geometric [5] and is trained on two 32-GB Nvidia Titan GPUs. Our code-base is publicly available at https://github. com/MasoudMo/echoglad.

4.3 Results

We evaluate models using Mean Absolute Error (MAE) in mm, and Mean Percent Error (MPE) in percents, which is formulated as MPE = $100 \times \frac{|L_{\text{pred}} - L_{\text{true}}|}{L_{\text{true}}}$, where L_{pred} and L_{true} are the prediction and ground truth values for every measurement. We also report the Success Detection Rate (SDR) for LVID for 2 and 6 mm thresholds. This rate shows the percentage of samples where the absolute error between ground truth and LVID predictions is below the specific threshold. These thresholds are chosen based on the healthy ranges for IVS (0.6-1.1cm), LVID (2.0-5.6cm), and LVPW (0.6-0.1cm). Hence, the 2 mm threshold provides a stringent evaluation of the models, while the 6 mm threshold facilitates the assessment of out-of-distribution performance.

In-Distribution (ID) Quantitative Results. In Tab. 1, we compare the performance of our model with previous works in the ID setting where the training and test sets come from the same distribution (*e.g.*, the same clinical setting), we separately train and test the models on the private and the public dataset. The results for the public dataset are provided in the supp. material (Table 1).

Out-of-Distribution (OOD) Quantitative Results. To investigate the generalization ability of our model compared to previous works, we train all models on the private dataset (which consists of a larger number of samples compared to UIC), and test the trained models on the public UIC dataset as shown in Tab. 2. Based on our visual assessment, the UIC dataset looks very different compared to the private dataset, thus serving as an OOD test-bed.

Qualitative Results. Failure cases are shown in supp. material (Fig. 3).

Ablation Studies. In Table 3, we show the benefits of a hierarchical graph representation with a multi-scale objective for the task of LV landmark detection. We provide a qualitative view of the ablation study in supp. material (Fig. 4).

Table 1: **Quantitative results** on the private test set for models trained on the private training set. We see that our model has the best average performance over the three measurements, which shows the superiority of our model in the in-distribution setting for high-data regime.

Model	MAE [mm] \downarrow			MPE [%] \downarrow			$ SDR[\%] \text{ of LVID} < \uparrow$	
	LVID	IVS	LVPW	LVID	IVS	LVPW	$\left 2.0 \text{ mm}\right $	6.0 mm
Gilbert et al. [6]	2.9	1.4	1.4	6.5	14.5	15.2	48.1	88.9
Lin et al. $[18]$	9.4	11.2	9.0	21.2	116.5	92.9	26.0	49.1
McCouat et al. [20]	2.2	1.3	1.4	4.8	13.5	15.1	58.3	93.9
Chen $et al. [2]$	2.3	1.2	1.2	5.2	12.6	13.8	60.4	92.6
Duffy $et al. [4]$	2.5	1.2	1.2	5.4	13.2	13.5	52.1	93.0
Ours	2.2	1.1	1.1	4.8	11.2	12.2	62.4	94.4

Table 2: **Quantitative results** on the public UIC test set for models trained on the private training set. This table shows the out-of-distribution performance of the models when trained on a larger dataset and tested on a smaller external dataset. We can see that in this case, our model outperforms previous works by a large margin, which attests to the generalizability of our framework.

Model	MA	E [m	ım]↓	M	PE [%	6]↓	SDR[%]	of LVID $<\uparrow$
	LVID	IVS	LVPW	LVID	IVS	LVPW	2.0 mm	6.0 mm
Gilbert et al. [6]	9.5	4.8	4.1	23.5	32.3	26.8	22.5	52.2
Lin et al. $[18]$	51.5	51.7	41.3	121.0	375.8	298.0	11.3	24.6
McCouat et al. [20]	5.9	3.6	4.4	18.5	30.5	36.4	34.6	72.3
Chen $et al. [2]$	7.4	5.3	6.9	22.5	49.4	62.4	28.9	65.3
Duffy $et al. [4]$	13.7	4.1	5.5	36.8	36.4	45.4	6.2	20.6
Ours	5.8	2.8	4.3	18.4	23.8	34.6	35.8	74.9

Table 3: Ablation results on the validation set of our private dataset. Vanilla U-Net uses a simple U-Net model, while U-Net Main Graph only uses the pixellevel graph (no aux. graphs). Main Model is our proposed approach. Lastly, Single-Scale Loss has the same framework as the Main Model but only computes the loss for the model's predictions on the main graph (no multi-scale loss).

Model	MPE [%]				
	LVID	IVS	LVPW		
Vanilla U-Net	5.31	13.17	13.47		
U-Net Main Graph	4.98	11.67	12.78		
Single-Scale Loss	5.41	12.37	12.8		
Main Model	4.91	11.45	12.36		

5 Conclusion and Future Work

In this work, we introduce a novel hierarchical GNN for LV landmark detection. The model performs better than the state-of-the-art on most measurements without relying on label smoothing. We attribute this gain in performance to two main contributions. First, our choice of representing each frame with a hierarchical graph has facilitated direct interaction between pixels at differing scales. This approach is effective in capturing the nuanced dependencies amongst the landmarks, bolstering the model's performance. Secondly, the implementation of a multi-scale objective function as a supervisory mechanism has enabled the model to construct a superior inductive bias. This approach allows the model to leverage simpler tasks to optimize its performance in the more challenging pixel-level landmark detection task.

For future work, we believe that the scalability of the framework for higherresolution images must be studied. Additionally, extension of the model to video data can be considered since the concept of intra-scale and inter-scale edges connecting nodes could be extrapolated to include temporal edges linking similar spatial locations across frames. Such an approach could greatly enhance the model's performance in unlabeled frames, mainly through the enforcement of consistency in predictions from frame to frame.

References

- Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE Transactions on Pattern Analysis and Machine Intelligence 40(4), 834–848 (2018) 3
- Chen, R., Ma, Y., Chen, N., Lee, D., Wang, W.: Cephalometric landmark detection by attentive feature pyramid fusion and regression-voting. In: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (eds.) Medical Image Computing and Computer Assisted Intervention. Springer International Publishing (2019) 4, 8
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021) 4
- Duffy, G., Cheng, P.P., Yuan, N., He, B., Kwan, A.C., Shun-Shin, M.J., Alexander, K.M., Ebinger, J., Lungren, M.P., Rader, F., Liang, D.H., Schnittger, I., Ashley, E.A., Zou, J.Y., Patel, J., Witteles, R., Cheng, S., Ouyang, D.: High-Throughput Precision Phenotyping of Left Ventricular Hypertrophy With Cardiovascular Deep Learning. JAMA Cardiology 7(4), 386–395 (04 2022) 3, 6, 8
- Fey, M., Lenssen, J.E.: Fast graph representation learning with PyTorch Geometric. In: ICLR Workshop on Representation Learning on Graphs and Manifolds (2019) 7
- Gilbert, A., Holden, M., Eikvil, L., Aase, S.A., Samset, E., McLeod, K.: Automated left ventricle dimension measurement in 2d cardiac ultrasound via an anatomically meaningful cnn approach. In: Smart Ultrasound Imaging and Perinatal, Preterm

10 M. Mokhtari et al.

and Paediatric Image Analysis: First International Workshop, SUSI 2019, and 4th International Workshop, PIPPI 2019, Held in Conjunction with MICCAI 2019 Proceedings, p. 29–37. Springer-Verlag, Berlin, Heidelberg (2019) 3, 8

- Gilmer, J., Schoenholz, S.S., Riley, P.F., Vinyals, O., Dahl, G.E.: Neural message passing for quantum chemistry. In: Proceedings of the 34th International Conference on Machine Learning - Volume 70. p. 1263–1272. JMLR.org (2017) 4
- 8. Goan, E., Fookes, C.: Bayesian Neural Networks: An Introduction and Survey, pp. 45–87. Springer International Publishing (2020) 3
- Goco, J.A.D., Jafari, M.H., Luong, C., Tsang, T., Abolmaesumi, P.: An efficient deep landmark detection network for PLAX EF estimation using sparse annotations. In: Linte, C.A., Siewerdsen, J.H. (eds.) Medical Imaging 2022: Image-Guided Procedures, Robotic Interventions, and Modeling. vol. 12034, p. 120340N. International Society for Optics and Photonics, SPIE (2022) 2
- Gradman, A.H., Alfayoumi, F.: From left ventricular hypertrophy to congestive heart failure: Management of hypertensive heart disease. Progress in Cardiovascular Diseases 48(5), 326–341 (2006), hypertension 2006 Update 1
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016) 4
- Howard, J.P., Stowell, C.C., Cole, G.D., Ananthan, K., Demetrescu, C.D., Pearce, K.A., Rajani, R., Sehmi, J.S., Vimalesvaran, K., Kanaganayagam, G.S., McPhail, E., Ghosh, A.K., Chambers, J.B., Singh, A.P., Zolgharni, M., Rana, B., Francis, D.P., Shun-shin, M.J.: Automated left ventricular dimension assessment using artificial intelligence developed and validated by a UK-wide collaborative. Circulation. Cardiovascular Imaging 14, e011951 – e011951 (2021) 6
- Jafari, M.H., Luong, C., Tsang, M., Gu, A.N., Van Woudenberg, N., Rohling, R., Tsang, T., Abolmaesumi, P.: U-land: Uncertainty-driven video landmark detection. IEEE Transactions on Medical Imaging 41(4), 793–804 (2022) 2, 3
- Kingma, D., Ba, J.: Adam: A method for stochastic optimization. International Conference on Learning Representations (2014) 7
- 15. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: International Conference on Learning Representations (2017) 7
- Li, W., Lu, Y., Zheng, K., Liao, H., Lin, C., Luo, J., Cheng, C.T., Xiao, J., Lu, L., Kuo, C.F., et al.: Structured landmark detection via topology-adapting deep graph learning. In: European Conference on Computer Vision. Springer (2020) 4
- Lin, C., Zhu, B., Wang, Q., Liao, R., Qian, C., Lu, J., Zhou, J.: Structure-coherent deep feature learning for robust face alignment. IEEE Transactions on Image Processing **30**, 5313–5326 (2021) 4
- Lin, J., Sahebzamani, G., Luong, C., Dezaki, F., Jafari, M., Abolmaesumi, P., Tsang, T.: Reciprocal landmark detection and tracking with extremely few annotations. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15165–15174. IEEE Computer Society (jun 2021) 2, 3, 8
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: 2021 IEEE/CVF International Conference on Computer Vision. pp. 9992–10002. IEEE Computer Society (oct 2021) 4
- McCouat, J., Voiculescu, I.: Contour-hugging heatmaps for landmark detection. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 20565–20573 (2022) 4, 8
- McFarland, T.M., Alam, M., Goldstein, S., Pickard, S.D., Stein, P.D.: Echocardiographic diagnosis of left ventricular hypertrophy. Circulation 50 (1978) 2

- Mokhtari, M., Tsang, T., Abolmaesumi, P., Liao, R.: Echognn: Explainable ejection fraction estimation with graph neural networks. In: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (eds.) Medical Image Computing and Computer Assisted Intervention. pp. 360–369. Springer Nature Switzerland (2022) 4
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 32. Curran Associates, Inc. (2019) 7
- Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) Medical Image Computing and Computer Assisted Intervention. pp. 234– 241. Springer International Publishing (2015) 3
- Scarselli, F., Gori, M., Tsoi, A.C., Hagenbuchner, M., Monfardini, G.: The graph neural network model. IEEE Transactions on Neural Networks 20(1), 61–80 (2009) 2, 4
- 26. Sofka, M., Milletari, F., Jia, J., Rothberg, A.: Fully convolutional regression network for accurate detection of measurement points. In: Cardoso, M.J., Arbel, T., Carneiro, G., Syeda-Mahmood, T., Tavares, J.M.R., Moradi, M., Bradley, A., Greenspan, H., Papa, J.P., Madabhushi, A., Nascimento, J.C., Cardoso, J.S., Belagiannis, V., Lu, Z. (eds.) Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support, pp. 258–266. Springer International Publishing (2017) 3

Supplementary Material

Masoud Mokhtari et al.

Electrical and Computer Engineering, University of British Columbia, Vancouver, BC, Canada



Fig. 1: Feature Generation for Graph Nodes - A CNN is initially used to expand the number of the feature maps. The intermediate features of the decoder part of a U-Net are then used as node features such that deeper representations correspond to node features of finer graphs.

Table 1: **Quantitative results** on the public UIC test set for models trained on the UIC training set. Although the number of training samples is much lower for UIC compared to our private dataset, we see that our model still outperforms previous works on average over the three measurements, which showcases the accuracy of our model in the low-data regime and in-distribution settings. Lin *et al.* is excluded since they require video inputs.

Model	MAE $[mm] \downarrow$			M	PE [%	6]↓	SDR[%]	of LVID $<\uparrow$
	LVID	IVS	LVPW	LVID	IVS	LVPW	2.0 mm	6.0 mm
Gilbert <i>et al</i> .	5.2	2.5	3.1	12.2	19.0	22.7	32.2	70.0
McCouat <i>et al</i> .	2.5	1.6	2.4	7.5	14.8	19.9	56.4	91.7
Chen <i>et al</i> .	2.3	1.5	2.3	7.1	12.5	21.4	57.3	94.6
Yao et al.	15.4	8.8	9.2	44.8	78.5	80.5	7.5	24.6
Duffy et al.	8.7	3.4	3.8	24.8	34.8	34.1	13.7	42.4
Ours	2.2	1.5	1.9	6.2	14.0	16.9	58.9	94.9



Fig. 2: **Hierarchical predictions** - An example of the model's prediction for an input echo. We show the model's prediction for the case where only three auxiliary graphs are used. We see that the model is learning the LV landmarks on different resolutions to achieve high accuracy for the pixel-level task. We show zoomed-in versions of the higher resolution task to enable comparison.



Fig. 3: Qualitative visualization of our model on two failure cases from the test set of our private dataset. The Failure 1 example is a low-quality PLAX image that also corresponds to a patient with severe LVH, a scenario that happens rarely in our dataset. The Failure 2 example belongs to a case with a low quality of PLAX with unclear boundaries for the walls and the chambers of the LV.



Fig. 4: Qualitative ablation results for the model architecture. Landmark heatmaps from top to bottom are color-coded with red, cyan, pink and green, respectively. We see that Vanilla U-Net (V. U-Net) struggles to make confident and accurate landmark predictions. While the addition of a main grid graph in U-Net Main Graph (U. M. Graph) relatively increases model's performance, it still does not produce accurate results. In contrast, the Main Model produces confident prediction heatmaps by relying on a hierarchical graph representation as well as multi-scale objectives. We also see that the removal of the multi-scale objective (Single-Scale Loss (SSL)) degrades performance.