

# PV-LVNet Direct left ventricle multitype indices estimation from 2D echocardiograms of paired apical views with deep neural networks

Rongjun Ge, Guanyu Yang, Yang Chen, Limin Luo, Cheng Feng, Heye Zhang,

Shuo Li

# ► To cite this version:

Rongjun Ge, Guanyu Yang, Yang Chen, Limin Luo, Cheng Feng, et al.. PV-LVNet Direct left ventricle multitype indices estimation from 2D echocardiograms of paired apical views with deep neural networks. Medical Image Analysis, 2019, 58, pp.101554. 10.1016/j.media.2019.101554. hal-02304385

# HAL Id: hal-02304385 https://univ-rennes.hal.science/hal-02304385

Submitted on 28 Nov 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Highlights

- An effective method for quantifying LV from multiple dimensions and views.
- A brand-new recurrent net for embedding subject and temporal information.
- An efficient location loss function for robust location and cropping.
- A creative regularization item for enhancing sequential data evolution fitting.

Johnskerkor



# PV-LVNet: Direct Left Ventricle Multitype Indices Estimation from 2D Echocardiograms of Paired Apical Views with Deep Neural Networks

Rongjun Ge<sup>a,c,d</sup>, Guanyu Yang<sup>a,c,d</sup>, Yang Chen<sup>a,b,c,d,\*</sup>, Limin Luo<sup>a,c,d</sup>, Cheng Feng<sup>e</sup>, Heye Zhang<sup>f</sup>, Shuo Li<sup>g,h,\*</sup>

<sup>a</sup>Laboratory of Image Science and Technology, School of Computer Science and Engineering, Southeast University, Nanjing, China

<sup>b</sup>School of Cyber Science and Engineering, Southeast University, Nanjing, China <sup>c</sup>Key Laboratory of Computer Network and Information Integration (Southeast

University), Ministry of Education, Nanjing, China

<sup>d</sup>Centre de Recherche en Information Biomedicale Sino-Francais (LIA CRIBs), Rennes.France

<sup>e</sup>Department of Ultrasound, The Third People's Hospital of Shenzhen, Shenzhen, China <sup>f</sup>School of Biomedical Engineering, Sun Yat-Sen University, Guangzhou, China <sup>g</sup>Department of Medical Imaging, Western University, London, Canada <sup>h</sup>Digital Imaging Group of London, London, Canada

## Abstract

Accurate direct estimation of the left ventricle (LV) multitype indices from two-dimensional (2D) echocardiograms of paired apical views, i.e., paired apical four-chamber (A4C) and two-chamber (A2C), is of great significance to clinically evaluate cardiac function. It enables a comprehensive assessment from multiple dimensions and views. Yet it is extremely challenging and has never been attempted, due to significantly varied LV shape and appearance across subjects and along cardiac cycle, the complexity brought by the paired different views, unexploited inter-frame indices relatedness hampering working effect, and low image quality preventing segmentation. We propose a paired-views LV network (PV-LVNet) to automatically and directly estimate LV multitype indices from paired echo apical views. Based on a newly designed Res-circle Net, the PV-LVNet robustly locates LV and

Preprint submitted to Journal Name

September 10, 2019

<sup>\*</sup>Corresponding author.

*Email addresses:* chenyang.list@seu.edu.cn (Yang Chen), sli287@uwo.ca (Shuo Li)

automatically crops LV region of interest from A4C and A2C sequence with location module and image resampling, then accurately and consistently estimates 7 different indices of multiple dimensions (1D, 2D & 3D) and views (A2C, A4C, and union of A2C+A4C) with indices module.

The experiments show that our method achieves high performance with accuracy up to 2.85mm mean absolute error and internal consistency up to 0.974 Cronbach's  $\alpha$  for the cardiac indices estimation. All of these indicate that our method enables an efficient, accurate and reliable cardiac function diagnosis in clinical.

*Keywords:* multitype cardiac indices, direct estimation, 2D echo, paired apical views, Res-circle Net

#### 1 1. Introduction

Accurate estimation for left ventricle (LV) indices (i.e., dimension, area 2 & volume) in two-dimensional (2D) echocardiograms (echo) of paired apical 3 views (i.e., paired apical four-chamber and two-chamber views) is of great 4 clinical significance to cardiac function evaluation (Schiller et al., 1989; Lang 5 et al., 2006, 2015). 2D echo is the most frequently used noninvasive modality 6 for the diagnosis of cardiac disease because of its unique ability to provide 7 real-time images of the beating heart, combined with its availability and 8 portability (Lang et al., 2015; Abdi et al., 2017; Gao et al., 2017, 2018). g The multitype indices of LV from 2D echo paired apical views, covering 10 long-axis dimension (LAD), short-axis dimension (SAD), area and volume, 11 which are measured from cavity as Fig.1, are most widely used to assess LV 12 chamber size and contractile function (Schiller et al., 1989; Pascual et al., 13 2003; Lang et al., 2015). It promotes comprehensive metrics from 1D (i.e., 14 LAD, SAD), 2D (i.e., area) and 3D (i.e., volume). Such paired orthogonal 15 apical four-chamber (A4C) and two-chamber (A2C) views enable a better 16 stereoscopic reproducibility of cardiac LV motion compared to the separate 17 plane observation from single view, for further comprehensive quantitative 18 functional analysis (Schiller et al., 1989; Ciampi and Villari, 2007). 19

The existing (semi-)automated cardiac indices estimation methods never refers to multitype indices in 2D echo sequences of paired apical views. These methods are mainly classified into two groups: segmentation and direct regression. However, the segmentation methods just enable limited simple index types (i.e. area) without extra interaction, and the existing direct



Figure 1: The multitype indices from the paired apical views (A2C & A4C) are critically important for clinical diagnosis, yet extremely laborious measurement. They cover the 1D and 2D metrics of each single view, and the 3D metric of union view, for a comprehensive assessment. (a) LAD: from the apex to the middle mitral valve plane. SAD: perpendicular to the long axis, at one-third of the LAD from the mitral valve plane. (b) Area: the whole LV cavity. (c) LV volume: jointly from A4C and A2C by using the biplane method of discs (modified Simpson's rule).

methods almost all focus on a single view of cardiac magnetic resonance 25 (CMR) causing limited observation and evaluation. Strong clinical evidence 26 shows that the indices from echo that cover multiple dimensions and views 27 enable a comprehensive cardiac diagnosis, yet their automated estimation is 28 still thwarted by inherently existing challenges such as 1) LV shape and 29 appearance in apical view significantly vary among subjects, and along 30 the cardiac cycle. 2) Although the paired views provide complementary 31 information, the different image structures are introduced with increased 32 3) Ambiguous relatedness inter frames hampers learning complexity. 33 procedure of sequential indices from better convergence and generalization. 34 4) Low image quality of echo, like fuzzy border, edge dropout, acoustic 35 shadows, etc., raises great challenges for automated methods, especially 36 segmentation method. 37

## 38 1.1. Related Works

**Segmentation methods** aim to achieve automated LV segmentation for 39 improving the diagnosis efficiency, however it is still an open and challenging 40 task, due to the inherent characteristics of the 2D echo, such as low signal-to-41 noise ratio, edge dropout, shadows, indirect relation between pixel intensity 42 and the physical property of the tissue, and anisotropy of ultrasonic image 43 formation (Carneiro et al., 2012). Active contours (Debreuve et al., 2001; 44 Malladi et al., 1995; Paragios, 2003) and deformable templates (Jacob et al., 45 2002; Nascimento et al., 2008) achieve good segmentation results relying on 46

the LV shape and appearance of the prior knowledge (Georgescu et al., 2005). 47 By considering use of inaccurate prior knowledge and low-level handcrafted 48 features may bound working robustness, the supervised deep learning method 49 (Mo et al., 2018; Chen et al., 2016; Carneiro et al., 2012; Oktay et al., 2018) 50 tries to learn information from data. The deep Poincar Map (Mo et al., 2018) 51 coupled deep learning with the dynamic-based labeling scheme to reduce the 52 requirement on the huge data; iMD-FCN (Chen et al., 2016) used the transfer 53 learning from cross domains to enhance the feature representation; Carneiro 54 et al. (2012) combined the deep belief networks, the decoupling rigid and 55 nonrigid classifiers and the derivative-based search to increase the robustness 56 for imaging conditions and LV shape variations; ACNNs (Oktay et al., 2018) 57 encouraged the models to follow the global anatomical properties of the 58 underlying anatomy via the non-linear representations of the shape learnt 59 from the stacked convolutional autoencoder. All of these show great potential 60 with the development of deep learning. Nevertheless, most of the working 61 LV segmentation methods in the practical clinical diagnosis are still semi-62 automatic, which need time-consuming user interaction to handle a great 63 number of medical images (Luo et al., 2018). 64

Direct regression methods without intermediate segmentation has 65 undergone a great development and recognition (Ravi et al., 2017; Peng 66 et al., 2016; Wu et al., 2017; Lathuilière et al., 2017; Pereira et al., 2018; 67 Zhen et al., 2014a, 2015b, 2017) for better and more efficient cardiac indices 68 estimation, but never performed on paired 2D echo apical views. By directly 69 analyzing LV biological structure, these methods provide effective tools to 70 automate the analysis of one single view from CMR, especially the short-71 axis view, and enable accurate and efficient diagnosis in clinical practice 72 (Zhen et al., 2016). With two-phase operation, LV volume (as integration 73 of cavity areas in short-axis view slices) is estimated on the handcrafted 74 cardiac image representation, including Bhattacharyya coefficient between 75 image distributions (Afshin et al., 2012, 2014), appearance features (Wang 76 et al., 2014), multiple low level image features (Zhen et al., 2014b), as 77 well as unsupervised features from multiscale convolutional deep belief 78 network (Zhen et al., 2016) and supervised descriptor learning (Zhen et al., 79 2015a). Instead of separate representation and regression, joint learning (Xue 80 et al., 2017a,c) captures task-relevant cardiac information for the indices 81 estimation. For a comprehensive assessment of cardiac function, Xue et al. 82 (2017b, 2018) achieve multitype indices estimation on short-axis view cardiac 83 CMR. However, all of these direct methods still have the limitation on 84

<sup>85</sup> 2D echo paired apical views, due to: 1) multitype indices estimation from <sup>86</sup> different views is ignored and lacked, 2) some cardiac indices in 2D echo, <sup>87</sup> like volume, are often obtained jointly from paired views, and 3) LV shape <sup>88</sup> in apical view is irregular and make it difficult to establish a standard <sup>89</sup> preprocessing method for getting LV cropping (short-axis view CMR just <sup>90</sup> need to manually find several relatively fixed landmarks).

#### 91 1.2. Contributions

In this paper, we propose a paired-views LV network (PV-LVNet) to 92 automatically achieve a high-quality estimation of LV multitype indices from 93 2D echo sequences of paired apical views. As shown in Fig.2, the network is 94 built based on our newly designed Res-circle Net, and implemented with three 95 interdependent functional parts: LV location module, image resampling and 96 LV indices module. The Res-circle Net for sequential analysis embedded with 97 subject's holistic characteristics and frame's temporal changes is used in both 98 LV location and indices modules. And functionally, the LV location module 99 with the anisotropic Euclidean distance loss shape-accordingly detects the 100 LV center in echo apical views. The image resampling further crops the LV 101 region of interest (LV-ROI) capable of efficiently reducing the interference of 102 various structure from the different views. Accepting the LV-ROI, the LV 103 indices module with the inter-frame gradient regularization and the views 104 union effectively makes the comprehensive, accurate and internally consistent 105

<sup>106</sup> indices estimation.

<sup>107</sup> The main contributions of our work include:

- For the first time, the proposed PV-LVNet enables an automatically and reliably comprehensive cardiac function clinical assessment from various dimensions and views by directly and accurately estimating LV multi-type indices on 2D echos of paired apical views.
- The newly designed Res-circle Net enables accurately and consistently
   estimating continuous changing centric positions and indices of LVs in
   echo sequence of each subject, by comprehensively combining both the
   subject-level base of cardiac cycle and the interrelated dynamic residual
   of each frame. Moreover, its residual transferring effectively reduces the
   gradient vanishing problem in recurrent net.
- The novel location loss in the form of anisotropic Euclidean distance (AED) guarantees robust and efficient location and cropping by matching the approximate bullet shape of LV in apical view echo.



Figure 2: The PV-LVNet simultaneously estimates multitype indices of various single (A4C, A2C) and union views (A4C+A2C) from paired apical 2D echo sequences, to provide a comprehensive cardiac function assessment. Based on the Res-circle Net (Sect. 2.1), it has three interdependent parts: LV location module (Sect. 2.2) for LV location, image resamping (Sect. 2.3) for LV-ROI cropping and LV indices module (Sect. 2.4) for multitype indices estimation.

• The gradient of LV indices between adjacent frames in a cardiac cycle creatively and effectively enhances sequential indices fitting, by fully exploring inter-frame relatedness to introduce frame-by-frame evolution characteristic to regularize indices estimation.

## 125 2. Methodology

121

122

123

124

As shown in Fig.2, based on the **Res-circle Net (Sect. 2.1)** to analyze echo sequences, the PV-LVNet entirely works via three interdependent parts: **LV location module (Sect. 2.2)**, **Image Resampling** (Sect. 2.3) and **LV indices module (Sect. 2.4)** for location, cropping and indices estimation.

To enable the comprehensive and efficient echo sequence analysis, the novel 130 Res-circle Net combines subject-level base for avoiding coarse sequential 131 estimation from zero level and temporal dynamic residual for developing 132 the refinement on each frame. To provide the robust LV location among 133 views for accurate indices estimation, LV location module creatively adopts 134 the loss in form of AED considering the LV shape in apical view echo. 135 To automatically crop LV-ROI with the interference of various structure in 136 paired views reduced, and build unblocked joint learning of location and 137 indices regression, image resampling, as a differentiable transformation, is 138 embedded. To achieve the various dimensional indices regression from single 139 and union views, LV indices module performs not only indices-aware feature 140 abstraction but also views union for 3D index. Moreover, to fulfill the inter-141 frame relatedness potential of indices for enhancing sequential data fitting, 142 the inter-frames gradient in the time polyline of the cardiac index is used to 143 deeply explore sequence evolution characteristics. 144

## 145 2.1. Res-circle Net for Analyzing Echo Sequence

The Res-circle Net combines both subject-level base and frame-level 146 residuals for a comprehensive analysis on echo sequence. Subject-level 147 base reflects the holistic characteristics among the different frames of 148 the same subject. It gives a whole and inherent expression on the echo 149 sequence and distinguishes different subjects. It is further extracted from 150 the representations of all frames. Frame-level residual reflects interrelated 151 temporal dynamic changes in the cardiac cycle. It enables a further 152 refinement on each frame. It is extracted by using the inter-frame relationship 153 among the whole cardiac cycle. The Res-circle Net captures interrelated 154 temporal residual of each frame, then adds the residual with subject-level 155 base together. It embeds subject and temporal information to guarantee 156 a stable and dynamic estimation for location and indices in continuous 157 moving and deforming LV. The net is implemented in the circle recurrent 158 of a novel residual learning and transferring convolutional unit named as 159 residual recurrent unit (RRU). 160

As shown in Fig.3, the Res-circle Net accepts current frame representation and links it to the integrated former residuals of the frames in the cycle, then adaptively updates the current frame-level residual and combine the residual with the subject-level base for a refined outputting. The Res-circle Net is achieved in the circle recurrent structure (Graves, 2012; Xue et al., 2017c) of RRU, which gives the memory characteristics of the cycle temporal



Figure 3: The Res-circle Net embeds both subject and interrelated temporal information together for comprehensive and reliable analysis on the echo sequence. It adaptively updates current dynamic change as residual by linking the current frame representation with the former memory in cycle, then adds such residual with the subject-level base together as the comprehensive state of the frame.

changes. Similar works to analyze data sequence can be seen in using LSTM 167 of recurrent neural networks (RNN) as: Xu et al. (2018) a lopt fully connected 168 LSTM (FC-LSTM) for the dependence crossing over a long time interval; Xue 169 et al. (2017c) further deployed circle FC-LSTM for shortening the distance 170 between the first and last frames; and convolutional LSTM (Xingjian et al., 171 2015) was developed for the spatial structure in the sequence. Specially, our 172 Res-circle Net of circle recurrent convolutional residual net is designed for 173 temporally-spatially modeling the residuals among frames and the entirety 174 of sequence, to the echo of dynamic and consecutive data. 175

The RRU has both functions of current frame state prediction and 176 residuals memory integration, as shown in Fig.4. In the output path, the 177 RRU provides the current state  $(state_i)$  for the followed regression, by adding 178 current frame-level residual  $(res_i)$  on the subject-level base (base). In the 179 hidden path, it transfers residual information  $(res_i)$  together with the formers 180  $(res\_mem_{i-1})$  to integrate residuals memory  $(res\_mem_i)$  for the next frame. 181 Instead of the frame-wise coarse estimation from the zero level, the net 182 provides such a more refined way as the subject-level base reflects the stable 183 base level of sequence and residual focuses on interrelated dynamic change of 184



Figure 4: Residual recurrent unit (RRU) has both functions of current frame state prediction and residual transfer. In output path, the current frame-level residual is added to the subject-level base for followed regression. In hidden path, the residual information is transferred together with the formers for the next frame.

each frame. Benefited from the residual connection with subject-level base
and former residuals, the net has powerful sequence analysis and temporal
modeling, and meanwhile effectively reduces the gradient vanishing problem
with the shortcut connection (Szegedy et al., 2017; He et al., 2016a,b).

The RRU takes both spatial structure and temporal information into 189 account. It uses convolution process, instead of full connection in traditional 190 RNN, to extract feature for keeping spatial correlation in the cardiac image. 191 In recurrent way, it maps the current frame to the integrated residual memory 192 to get its current frame-level residual. The inherent potential spatiotemporal 193 characteristic in echo sequence is effectively mined and transmitted. Given 194 the inputting individual frame representation  $frame_i$  at each time step i, 195 the memory  $res_mem_{i-1}$  from the previous frames, and the subject-level 196 base base, RRU gets the current frame-level residual  $res_i$  for the updated 197 memory  $res\_mem_i$  and outputting state representation  $state_i$ , as: 198

$$res_{i} = LN(ELU(LN((frame_{i} \oplus res\_mem_{i-1}) * W_{1} + b_{1})) * W_{2} + b_{2})$$

$$res\_mem_{i} = ELU(res_{i} + res\_mem_{i-1})$$

$$state_{i} = ELU(base + res_{i})$$
(1)

where  $W_1$  and  $W_2$  are convolutional kernels in Conv1 and Conv2,  $b_1$  and  $b_2$ represents biases.  $\oplus$  means concatenation, \* is convolution operation, and LN, *ELU* denote the element-wise transformations of layer normalization (Ba et al., 2016) and exponential linear unit (Clevert et al., 2015).

## 203 2.2. LV Location Module for Detecting Left Ventricle Center

LV location module aims to detect continuously moving LV center in both 204 A4C and A2C sequences, as in Fig.5. It has four steps: 1) CNN-loc firstly 205 extracts cardiac subject-level base and individual frame representations of 206 the cardiac sequence and feeds them to the res-circle net; 2) **Res-circle Net** 207 then models sequential LV moving in cardiac cycle for the final location, 208 with subject's holistic position and frame's temporal changes embedded; 209 3) Fully connected (FC) layer further performs LV center coordinate 210 regression with the output of Res-circle Net fed; And 4) **AED metric** is 211 used to measure the regressed center with anisotropic scaling by considering 212 approximate bullet shape of LV in echo apical views for robust location. 213

Advantageously, LV Location Module is benefited from the special design of **CNN-loc** and **AED location metric**, besides Res-circle Net that has been proposed in Sect. 2.1.

CNN-loc. To get expressive and task-aware representation of individual
 frame and entire subject on the paired echo sequences, CNN-loc consists



Figure 5: To achieve locating continuously changing center of LV in both A4C and A2C sequences, LV location module works via: 1) CNN-loc extracts subject-level base and frame representations for both paired views. 2) Res-circle Net captures residual information of each frame by leveraging inter-frame relationship for modeling dynamic changes, and further combine subject-level base to provide the frame state for location. 3) FC layer linearly regresses LV center coordinate. 4) The metric of anisotropic Euclidean distance (AED) ensures the robust location.

of several shared layers for general expression and two shallow paths that 219 further refine on A4C and A2C respectively considering big view difference 220 and enhancing robustness, as shown in Fig.6 (a). The individual cardiac 221 distribution in each frame is extracted by the hierarchical convolutions, and 222 the global sequence base of the subject is then captured by concatenating 223 all these individual representations together with a further convolution 224 operation followed so that holistically characterizes all frames. The backbone 225 structure of CNN-lock is the stack of the successive convolutional blocks (He 226 et al., 2016a) in Fig.6 (b), which chooses identity map for the layer input 227 and output of the same size, or convolution of kernel size  $1 \times 1$  to match 228 dimensions. Such block promotes information propagation both forward and 220



Figure 6: CNN-loc gets subject-level base and frame representation of paired echo sequences. (a) CNN-loc is composed of several shared layers and two shallow paths refined on A4C and A2C. (b) The stacked block in CNN-loc. The use of short-cut connection accelerates the net convergence and improve learning performance.



(a) AED has elliptical isarithm, enabling robust metric on LV location for ROI



Figure 7: The anisotropic Euclidean distance (AED) provides an elliptical isarithm to match the approximate bullet shape of LV in apical view echo and enable a more reasonable and robust LV location metric than the isotropic Euclidean distance (IED). (a) Considering LV shape, AED gives different scaling on the horizontal and vertical direction to construct elliptical isarithm for efficient LV location in apical view echo. (b) IED causes pool metric on LV location due to its circle isarithm.

backward and hence accelerate the net convergence and improve learning performance (Szegedy et al., 2017; Yu et al., 2017). The configurations of the stacked convolutions in CNN-loc are:  $7 \times 7 \times 64$  with stride 2 for conv1, *channel* = 64, 128, 128, 256, 256 for convolutional blocks (conv{2,3,4,5}-block and conv-*path*-fram), and  $3 \times 3 \times 256$  with stride 1 for conv-*path*-seq.

**AED location metric.** To achieve a structure matching location measurement, anisotropic Euclidean distance (AED) is deployed on the regressed center with the different metric scaling on horizontal and vertical directions, as shown in Fig.7 (a). Differently and traditionally, the location metric generally uses strict isotropic Euclidean distance (IED) in Eq. (2), where the regressed result  $\hat{O} = (\hat{o}_x, \hat{o}_y)$  and the ground truth  $O = (o_x, o_y)$ .

$$distance_{IED} = \left\| O - \hat{O} \right\| \tag{2}$$

However, the shape of the LV is approximate to the bullet, so that the 241 regressed points with same IED values still cause different influences to the 242 ROL and smaller IED does not mean a more accurate location. For example, 243  $O_1$  and  $O_5$  in Fig.7 (b) fall on the same circle isarithm of the IED to the LV 244 center O, and  $\hat{O}_4$  even has smaller IED than  $\hat{O}_1$ . But only the  $\hat{O}_1$  centered 245 square contains the entire LV cavity, while  $O_4$  and  $O_5$  lead to the weak ROIs. 246 In order to overcome the shortcoming in IED, AED using anisotropic 247 scaling is a more reasonable metric that conforms to the LV shape. 248 Comparing Figs.7 (a) with (b),  $\hat{O}_4$  and  $\hat{O}_5$  that have the same IED value 249 as  $\hat{O}_1$  or smaller than  $\hat{O}_1$  are outside the elliptical isarithm of AED, which 250 means getting higher AED metric. It aligns with their poor ROI quality in 251 Fig.7 (b). Besides, the ROIs centered by the points  $O_2$  and  $O_3$  that fall on 252 the ellipse in Fig.7 (a) have the same ROI situation as  $O_1$ , that the entire 253

LV cavity is contained and close to the square border, and gains the same metric. Therefore, the AED introduces a more robust and effective location metric for LV. The AED calculation is given in Eq. (3).

$$distance_{AED} = \sqrt{\beta \cdot (\hat{o}_x - o_x)^2 + (1 - \beta) \cdot (\hat{o}_y - o_y)^2}$$
(3)

#### 257 2.3. Image Resampling for Cropping LV-ROI

Image resampling is implemented via spatial transform and bilinear 258 interpolation to automatically crop LV-ROI according to the location from 259 Sect. 2.2. Image resampling puts attention on determining the region most 260 related to the LV. It aims to reduce the disturbance from the other pathology 261 caused by various structure and extra chambers in different views, with the 262 LV-ROI being cropped. Also, the LV-ROI sequence maintains the relative 263 shapes of LVs among different frames to not destroy the inherent subject 264 characteristics and frame-by-frame LV dynamic changes along the cardiac 265 cycle for developing the sequential LV indices estimation of each subject. In 266 a similar work, Dai et al. (2016) used ROI warping layer to crop feature map 267 regions for refining further semantic segmentation. Additionally, Jaderberg 268 et al. (2015) and Vigneault et al. (2018) used STN to spatially transform 269 intermediate feature maps or inputting image for improving performance in 270 classification and medical segmentation, respectively. 271

In our work, the image resampling transforms the images into the pattern that are centred on the predicted LV centre, and crops them to the predefined dimensions images. Given the predicted LV centre  $\hat{O} = (\hat{o}_x, \hat{o}_y)$  and the source echo image *I*, the target LV-ROI image  $I^{ROI}(\hat{O})$  is obtained by the image resampling as formulated as the differentiable linear transformation:

$$I^{ROI}(\hat{O}) = B(T(\hat{O})) \cdot I.$$
(4)

In Eq. (4),  $T(\cdot)$  is the spatial transform that firstly translates the echo 277 image I horizontally and vertically to be centred on O and then scales 278 the translated image to crop a  $153.6 \, pixel \times 153.6 \, pixel$  image (physical 279 dimensions 79.49  $\sim$  115.80mm  $\times$  79.49  $\sim$  115.80mm with pixel space 280  $0.5175mm/pixel \sim 0.7539mm/pixel$ ) centred on the predicted LV centre. 281  $B(\cdot)$  means bilinear interpolation further calculates the pixel value and 282 produces the LV-ROI in a sufficiently fine resolution which is set as same 283 as the original echo image, for the following indices estimation. 284



Figure 8: To estimate multitype indices from single/union views, LV indices module works via: 1) CNN-ind1+Feature Concatenation+CNN-ind2 gets feature representation on both entire subject and individual frame for all single and union views. 2) Res-circle Net models frame-by-frame dynamic residuals in the cardiac cycle by inter-frame relationship, then add them with the subject-level base of the holistic shape, for embedding subject and temporal information. 3) FC layer regresses indices with the outputs of the Res-circle Net. 4) Inter-frames gradient regularizes indices changes among frames to enhance sequential indices estimation.

## 285 2.4. LV Indices Module for Estimating Multitype Indices

LV indices module is designed to estimate multitype sequential cardiac 286 indices in union and single views from continuously deformed LVs, as shown 287 in Fig. 8. It consists of four components: 1) CNN-ind1 + Feature 288 **Concatenation** + **CNN-ind2** makes frame and subject feature extraction, 289 as well as union view representation. 2) Res-circle Net combines subject 290 holistic shape and temporal deformation. 3) FC layer further regresses on 291 the feature representation from Res-circle Net to estimate all indices. And 4) 292 **Inter-frames Gradient** is meanwhile introduced to regularize the indices 293 evolution among frames. 294

The superiority of LV indices module benefits from the special in 295 CNN-ind1 + Feature Concatenation + CNN-ind2 and Inter-frames 296 Gradient Regularization, besides Res-circle Net demonstrate in Sect. 2.1. 297 **CNN-ind1** + Feature Concatenation + CNN-ind2. In order to 298 get both the frame and the subject features for all union and single views, 299 it is further split and developed from CNN-loc that CNN-ind1 conducts the 300 preliminary view-specialized representation on paired fed A4C and A2C ROI 301 sequence, Feature Concatenation integrates the union view information via 302



Figure 9: Inter-frames gradient regularization promotes sequential indices regression. (a) Frame-by-frame evolution of index is reflected by the polyline of index value vs. frame. (b) Inter-frames gradient regularizes frame-by-frame evolution of estimated results to strengthen sequential indices fitting. It reveals index changes among frames, and thus characterizes index evolution. Evolution is an important metric in measuring the similarity between two sequential data.

unifying the A4C and the A2C representations along the feature channel,
and CNN-ind2 of the individual and the holistical features extraction is
further performed on all the A4C, A2C and union view. In the procedure,
the union view aims to construct the 3D spatial information from the two
orthogonal views for the volume of 3D indices estimation and meanwhile
further strengthening the contact among all views.

Inter-frames Gradient Regularization. For the accurate sequential 309 indices estimation, the gradient inter frames is used with considering the 310 evolution characteristics in the cardiac cycle. The frame-by-frame evolution 311 of index in the cardiac cycle is shown in Fig. 9(a) with the polyline of index 312 value vs. frame, it reflects the trend over time. And the gradient can be 313 explored to depict these evolution characteristics of time polyline in 9(a), 314 so that enhance the sequential indices fitting elegantly with the interrelated 315 fluctuation regularized on the sequence of the preliminary regressed index. As 316 shown in 9(b), it measures the slope of the secant passing through adjacent 317 discrete points. Given the regressed result  $\hat{y}^f$ , and normalized the frame 318 interval  $\Delta t$  as  $\pm 1$  for both adjacent frames, the inter-frames gradient  $k^f$  at 319 each frame step is defined as: 320

$$\hat{k}^{f} = \left(\hat{k}^{f-}, \hat{k}^{f+}\right), and \begin{cases} \hat{k}^{f-} = \hat{y}^{f} - \hat{y}^{f-1} \\ \hat{k}^{f+} = -(\hat{y}^{f} - \hat{y}^{f+1}) \end{cases}$$
(5)

where  $\hat{k}^{f-}$  and  $\hat{k}^{f+}$  mean left and right gradient of frame f, respectively.  $\hat{k}^{f}$ thus effectively characterizes index evolution of frame f between adjacent frames f-1 and f+1.

Therefore, the inter-frames gradient of each index among the cardiac cycle is introduced to fit the trend of polylines of regressed results and ground truth, as shown in Fig. 9(b), to enhance sequential LV indices estimation. Euclidean distance is used to calculate the gap of the change rate of each frame between regressed results and ground truth. The fitness of sequential indices evolution is measured as Eq. (6). In addition to the fitting of each index value, such evolution of the index further gives full play to the constraint between adjacent frames, and can be used as a regularization item to strengthen sequential objects estimation.

$$Reg_{grad} = \sum_{f=1}^{N} \sum_{t} \left\| k_t^f - \hat{k}_t^f \right\|_2 \tag{6}$$

where  $f \in \{1, 2, ..., N\}$  for all frames in cardiac cycle,  $t \in \{LAD_{A4C}, SAD_{A4C}, SAD_{A4C},$ 

## 335 3. Joint Loss Function for Different Tasks

The loss function in our work is designed for optimizing the two trainable modules (LV location module and indices module, while image resampling, as a powerful linear transformation, needs no training) of different tasks in the integrated PV-LVNet, so that the task-inter relevance and dependence enable the modules to mutually promote refinement of each other. The joint loss  $\mathcal{L}_{joint}$  is constructed as:

$$\mathcal{L}_{joint} = \lambda_1 \mathcal{L}_{loc} + \lambda_2 \mathcal{L}_{ind} + \lambda_3 R(\theta) \tag{7}$$

where  $\mathcal{L}_{loc}$  and  $\mathcal{L}_{ind}$  are the loss functions of location and indices estimation,  $R(\theta) = \|\theta\|_2^2$ , known as Tikhonov regularization for improving the training generality, is used as the regularization item of the network parameter vector  $\theta$  with  $l_2$ -norm.  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are set as 1000.0, 1.0 and 0.1.

The location loss function  $\mathcal{L}_{loc}$  aims to guarantee a robust location of LV for LV-ROI cropping. It is constructed with AED for taking account of the approximate bullet shape of LV. The definition of  $\mathcal{L}_{loc}$  is given by:

$$\mathcal{L}_{loc} = \frac{1}{N} \sum_{f=1}^{N} distance_{AED}^{f} \tag{8}$$

where  $distance_{AED}^{f}$  denotes  $distance_{AED}$  (defined in Eq. (3)) for the predicted centre in each frame f.

The indices loss function  $\mathcal{L}_{ind}$  aims to boost high-quality indices regression. It utilizes not only the MAE of indices value estimation error in each frame but also the trend between indices of adjacent frames for both accuracy and inter consistency of the sequential indices estimation, as:

$$\mathcal{L}_{ind} = \frac{1}{N} \sum_{t} \sum_{f=1}^{N} \left| \hat{y}_t^f - y_t^f \right| + Reg_{grad} \tag{9}$$

where the first item is the MAE loss of indices,  $Reg_{grad}$  (defined as Eq.(6)) is the inter-frames gradient regularization item for indices evolution.

#### 357 4. Experiment Configurations

**Dataset.** A dataset of 2D echos with the ground truth is used to evaluate 358 our method, which includes 2000 echo images from 50 subjects collected from 359 2 hospitals. Each subject provides both paired A4C and A2C views echos, 360 with the temporal resolution of 20 frames per cardiac cycle and the resize of 361  $256 \times 256$ . All ground truth of location and indices are manually annotated 362 by two experienced cardiac radiologists with double-checking. In training, 363 location labels are normalized to  $[-1,1] \times [-1,1]$  through subtracting half of 364 the image dimension (128) and then being divided by the image dimension 365 (256). The labels of 1D (i.e.,  $LAD_{A4C}$ ,  $SAD_{A4C}$ ,  $LAD_{A2C}$  and  $SAD_{A2C}$ ), 2D 366 (i.e.,  $Area_{A4C}$  and  $Area_{A2C}$ ) and 3D (i.e., Volume) metrics are normalized by 367 LV-ROI dimension  $(\frac{256}{p})$ , where  $\frac{1}{p} = 0.6$  is set according to prior investigation on our dataset), area  $((\frac{256}{p})^2)$  and volume  $((\frac{256}{p})^3)$ , respectively. **Data Augmentation.** To avoid the over-fitting and improve the 368 369

Data Augmentation. To avoid the over-fitting and improve the generalization, we augment the dataset to 8000 images by three strategies as: 1) randomly rotating between -15° and 15°; 2) randomly zooming between 0.9 and 1.1 times; and 3) the combination of random rotation + zoom.

374 Configurations. The net is implemented by Tensorflow, and performed
 375 on NVIDIA P100 GPU. Ten-fold cross validation is employed for performance
 376 evaluation and comparison.

**Evaluation Metrics.** We evaluate the performance of the PV-LVNet in 377 terms of estimation accuracy and internal consistency for multitype indices 378 of all frames in the cardiac cycle. The evaluation is performed with two 379 metrics including: the mean absolute error (MAE) for measuring accuracy 380 and Cronbach's  $\alpha$  (Cronbach, 1951) for measuring internal consistency 381 between the estimated results and the corresponding ground truth. Denote 382 the estimated cardiac index and ground truth of the *i*th subject and 383 the *f*th frame as  $\hat{y}_{t,i}^f$  and  $y_{t,i}^f$ , where  $t \in \{LAD_{A4C}, SAD_{A4C}, Area_{A4C}, LAD_{A2C}, SAD_{A2C}, Area_{A2C}, Volume\}$  for index types. The **MAE** of each 384 385 cardiac index is given by  $MAE_t = \frac{1}{S \times N} \sum_{i=1}^{S} \sum_{f=1}^{N} \left| \hat{y}_{t,i}^f - y_{t,i}^f \right|$ , where S and F are the number of subjects and frames, respectively. **Cronbach's**  $\alpha$  of 386 387 each cardiac index is calculated as  $\alpha_t = 2 \cdot \left(1 - \frac{\sigma_{\hat{y}_t}^2 + \sigma_{\hat{y}_t}^2}{\sigma_{\mathcal{X}_t}^2}\right)$ , where  $\mathcal{X}_t$  is the sum 388 of estimated indices  $\hat{\mathcal{Y}}_t = \left\{ \hat{y}_{t,1}^1, \hat{y}_{t,1}^2, \hat{y}_{t,1}^3, ..., \hat{y}_{t,S}^N \right\}$  and corresponding ground 389

truth  $\mathcal{Y}_t = \{y_{t,1}^1, y_{t,1}^2, y_{t,1}^3, ..., y_{t,S}^N\}$ , i.e.,  $\mathcal{X}_t = \hat{\mathcal{Y}}_t + \mathcal{Y}_t$ . Moreover,  $\sigma_{\hat{\mathcal{Y}}_t}^2, \sigma_{\mathcal{Y}_t}^2$  and  $\sigma_{\mathcal{X}_t}^2$  are the corresponding variances for  $\hat{\mathcal{Y}}_t, \mathcal{Y}_t$  and  $\mathcal{X}_t$ .

#### <sup>392</sup> 5. Results and Analysis

We conduct a set of experiments to evaluate the performance of our proposed PV-LVNet, including: 1) overall performance; 2) effectiveness of res-circle net; 3) effectiveness of anisotropic Euclidean distance location loss; 4) effectiveness of inter-frames gradient regularization; 5) performance comparison with relevant methods; 6) performance of activation function and Hyper parameter selection.

#### 399 5.1. Overall Performance

As shown in the last column of Table 1, the proposed PV-LVNet achieves 400 excellent estimation accuracy and internal consistency on all the 7 different 401 indices, which are attributable to comprehensively analyzing sequential 402 echos, robustly locating and cropping LV, deeply exploiting inter-frame 403 indices relatedness. It gains extremely low MAE of 2.85mm, 3.16mm, 404 3.06mm, 2.98mm,  $287mm^2$ ,  $264mm^2$  and 10.7ml for  $LAD_{A2C}$ ,  $SAD_{A2C}$ , 405  $LAD_{A4C}$ ,  $SAD_{A4C}$ ,  $Area_{A2C}$ ,  $Area_{A4C}$  and LV volume, as well as high 406 Cronbach's  $\alpha$  all exceeding 0.9, with the manually obtained ground truth. 407

Moreover, our proposed PV-LVNet also achieves high coincide indices estimation along the cardiac cycle, indicating powerfully modeling the LV activity. As shown in Fig. 10, it reaches extremely low normalized root-mean-square error of 1.26% (NRMSE,  $NRMSE = \frac{1}{y} \sum_{f=1}^{N} \frac{(\hat{y}^f - y^f)^2}{N}$ ) with the sequential ground truth, on average. Such rare few deviations strongly validate that the network effectively captures the activity pattern of sequential LVs.

<sup>415</sup>Our method is also very efficient in running time. The training takes 16.36 <sup>416</sup>hours with one P100 GPU. The testing takes only 0.70 seconds per subject. <sup>417</sup>Clearly, our method enables a real-time solution for clinical application.

## <sup>418</sup> 5.2. Effectiveness of Res-circle Net

As shown in Table 2, the Res-circle Net decreases the MAE by 15.7% (e.g., 15.7% =  $\frac{1}{7} \left[ \frac{3.46-2.85}{3.46} + \frac{3.64-3.16}{3.64} + \frac{3.37-3.06}{3.37} + \frac{3.24-2.98}{3.24} + \frac{336-287}{336} + \frac{321-264}{321} + \frac{15.1-10.7}{15.1} \right] \right)$ and gains exceeding 0.9 Cronbach's  $\alpha$  on all indices, compared to the the situation of being replaced by CNN in the PV-LVNet for revealing its effectiveness. By combining subject-level holistic characteristics and

Table 1: The proposed method gains most advanced performance in the various dimensional metrics for LV of all views compared to the existing methods. It achieves higher accuracy and more excellent internal consistency, with lower MAE (18.9%  $\downarrow$ ) and higher Cronbach's  $\alpha$  (> 0.9) for each LV index. MAE and  $\alpha$  are shown in each cell.

	Multi-features+RF	SDL+AKRF	MCDBN+RF	Indices-Net	U-Net	PV-LVNet
One-dimensional Metric (mm)						
$LAD_{A2C}$	$3.52 \pm 3.10$	$3.29 \pm 2.48$	$3.44 \pm 3.18$	$3.19 \pm 2.43$	/	$2.85{\pm}2.46$
	0.895	0.913	0.898	0.923	/	0.941
$SAD_{A2C}$	$3.76 {\pm} 3.02$	$4.51 \pm 3.34$	$3.81 \pm 3.13$	$3.60 {\pm} 2.82$	/	$3.16{\pm}2.68$
	0.890	0.866	0.895	0.910	/	0.930
$LAD_{A4C}$	$3.86 {\pm} 3.48$	$3.73 {\pm} 3.05$	$3.93 {\pm} 3.38$	$3.29 \pm 2.42$	/	$3.06{\pm}2.73$
	0.864	0.904	0.863	0.896	/	0.932
SAD	3.23 + 2.91	$3.21 \pm 2.82$	$3.18 {\pm} 3.00$	$4.27 \pm 3.37$	,	$2.98{\pm}2.85$
$SAD_{A4C}$	0.901	0.907	0.903	0.887	/	0.917
Two-dimensional Metric $(mm^2)$						
$Area_{A2C}$	$331 \pm 259$	$321 \pm 274$	320±264	$361 \pm 431$	$393 \pm 338$	$287{\pm}284$
	0.870	0.884	0.885	0.876	0.887	0.907
$Area_{A4C}$	$323 \pm 266$	$280 \pm 236$	$312 \pm 255$	$354 \pm 338$	$392 \pm 305$	$264{\pm}228$
	0.902	0.934	0.915	0.885	0.901	0.940
Three-dimensional Metric (ml)						
Volume	$16.1 \pm 14.2$	$16.4{\pm}14.6$	16.1±14.0	$15.3 \pm 8.7$	/	$10.7{\pm}7.6$
	0.918	0.922	0.925	0.938	/	0.974

interrelated temporal changes existing in echo sequence, the Res-circle Net 424 outperforms CNN which just performs independent processing for each 425 frame, on accuracy and internal consistency. Adding subject-level base 426 and interrelated dynamic residual of each frame together, the res-circle net 427 enables and enhances refined sequential indices estimation by leveraging 428 inter-frame temporal relationship and avoiding coarse estimation on each 429 separate frame from zero level to improve accuracy. Moreover, introducing 430 subject-level and temporal characteristics, the Res-circle Net guarantees 431 excellent internal consistent estimation across subjects and among frames 432 with the ground truth. 433

#### 434 5.3. Effectiveness of AED location loss

As shown in Table 3, the AED location loss ensures developing accurate indices estimation. Compared with using IED in location, the AED location loss significantly decreases the MAEs by 21.3%, 11.0%, 13.8% and 30.5% on LAD, SAD, area and volume on average. These improvements are resulted from the fact that IED location loss effectively provides a robust and efficient location and cropping for indices estimation. It suits LV in apical view echo by adopting different scaled metrics on different directions to match



Figure 10: The proposed PV-LVNet effectively achieves high coincide indices estimation along the cardiac cycle to model the LV activity. The polygonal lines reflect the frame-wise value of each index for average subject. The normalized root mean square error (NRMSE) is used to measure the deviation between the polygonal lines of the estimated value and ground truth. As the results show, the network gains low NRMSE of 1.26% on average, with rarely few deviations.

the approximate bullet shape that is more strict on locations in the vertical 442 direction than the horizontal direction, while the general IED loss can only 443 provide a low-quality metric of no direction difference. Thus, LAD, area and 444 volume which are extremely sensitive to vertical direction location get the 445 highest improvements. Additionally, the SAD which is the most difficult to 446 be estimated due to its non-independent measurement and a certain degree 447 dependence on LAD still gets an obvious improvement of 11.0% with more 448 accurate LAD. 449

## 450 5.4. Effectiveness of Inter-frames Gradient Regularization

As shown in Table 4, the inter-frames gradient regularization is capable of increasing the internal consistency of the estimated results with the ground truth. It gains higher Cronbach's  $\alpha$  exceeding 0.9 on all indices and increased from 0.914 to 0.934 on average. By measuring the index change rate between adjacent frames, the inter-frames gradient is used to fit indices frame-byframe evolution in sequence. So that the estimated sequential indices are

Table 2: The Res-circle Net contributes to high estimation accuracy and excellent internal consistency. It obtains lower MAE (15.7%  $\downarrow$ ) and higher Cronbach's  $\alpha$  (> 0.9) than being replaced by CNN.

	CNN	Res-circle Net
One-	dimensional N	Metric (mm)
LAD	$3.46 \pm 2.87$	$2.85{\pm}2.46$
$LAD_{A2C}$	0.915	0.941
SAD	$3.64 \pm 2.86$	$3.16{\pm}2.18$
$SAD_{A2C}$	0.913	0.930
	$3.37 \pm 2.66$	$3.06{\pm}2.73$
$LAD_{A4C}$	0.893	0.932
SAD	$3.24{\pm}2.65$	$2.98{\pm}2.85$
$SAD_{A4C}$	0.888	0.917
Two-o	dimensional N	Ietric $(mm^2)$
Area	$336 \pm 279$	$287{\pm}284$
$Area_{A2C}$	0.885	0.907
4	$321 \pm 289$	$264{\pm}228$
$Area_{A4C}$	0.908	0.940
Three	e-dimensional	Metric $(ml)$
Volume	$15.1 \pm 11.8$	$10.7 \pm 7.6$
v orume	0.935	0.974

<sup>457</sup> regularized to get consistent variation with the ground truth.

Besides, the inter-frames gradient regularization also enhances sequential 458 data fitting to ensure stable and accurate estimation across the whole cardiac 459 cycle, as shown in Fig. 11. It not only gains consistently lower estimation 460 error, but also increases the stability by 18.7% on average. The inter-frames 461 gradient regularization mines indices inter-frame relatedness to learn the 462 fluctuation across the cardiac cycle. Such fluctuation explicitly explores the 463 constraints among indices of different frames to promote stable and accurate 464 estimation and reduce pulse estimation error for sequential indices. 465

## 466 5.5. Performance Comparison with Relevant Methods

Our PV-LVNet achieves the most advanced performance in the various 467 dimensional metrics for the LV of all views compared to the existing methods: 468 1) the two-phase direct estimation including Multi-features+RF (Zhen et al., 469 2014b), SDL+AKRF (Zhen et al., 2015a), MCDBN+RF (Zhen et al., 2016); 470 2) the end-to-end direct estimation, i.e. Indices-Net (Xue et al., 2017a); 3) the 471 indirect estimation with segmentation U-net (Ronneberger et al., 2015). As 472 shown in Table 1, our method significantly decreases the MAE by 18.9% on 473 average on all indices, compared to these methods. Besides, it simultaneously 474 maintains excellent internal consistency with the manually obtained ground 475 truth by high Cronbach's  $\alpha$  all exceeding 0.9. 476



Figure 11: The inter-frames gradient well regularizes the network to enhance sequential data fitting. The polygonal lines record the frame-wise average MAE of each index. The standard deviation (std) is used to reflect the dispersion of MAE polygonal lines across a whole cardiac cycle. As the results show, using the inter-frames gradient regularization for the sequential indices decreased std by 18.7% compared to be removed, on average. It means stable and robust estimation on each frame. Also, the polygonal lines show consistently lower estimation error with inter-frames gradient regularization.

	IED location loss	AED location loss		
	Long-axis Dimens	ion $(mm)$		
$LAD_{A2C}$	$3.89{\pm}2.89$	$\boldsymbol{2.85 {\pm} 2.46}$		
$LAD_{A4C}$	$3.62{\pm}2.38$	$3.06{\pm}2.73$		
Averge	$3.76{\pm}2.64$	$2.96{\pm}2.60$		
	Short-axis Dimens	sion (mm)		
$SAD_{A2C}$	$3.48{\pm}2.84$	$3.16{\pm}2.68$		
$SAD_{A4C}$	$3.41{\pm}2.84$	$\boldsymbol{2.98 {\pm} 2.85}$		
Average	$3.45{\pm}2.84$	$3.07 {\pm} 2.77$		
Area $(mm^2)$				
$Area_{A2C}$	$322 \pm 255$	$287{\pm}284$		
$Area_{A4C}$	$314{\pm}224$	$264{\pm}228$		
Average	$318{\pm}240$	$274{\pm}259$		
Volume (ml)				
Volume	$15.4{\pm}15.6$	$10.7{\pm}7.6$		

Table 3: The AED location loss ensures developing accurate estimation for LV indices. It brings higher estimation accuracy than the IED location loss, with lower MAE  $(17.3\% \downarrow)$  on each type of cardiac indices.

477 In detail, our method is superior to the relevant methods as:

1) The proposed PV-LVNet outperforms the two-phase direct method, 478 with the average MAE decreased by 16.2%, 12.3% and 34.0% on 1D, 2D 479 and 3D metrics, respectively. Different from these compared methods, 480 the proposed method jointly learns the deep task-aware information and 481 regresses target in an end-to-end way, instead of the split handcrafted feature 482 extraction and regression. It is obviously validated on the volume estimation. 483 The proposed method conducts a deeper learning on the concatenated feature 484 jointly with volume estimation, and gets 34.0% improvement. 485

2) The proposed PV-LVNet outperforms the existing end-to-end direct 486 method, with the average MAE decreased by 19.4% and all Cronbach's 487  $\alpha$  increased to above 0.9. All of these are own to the fact that the 488 proposed method effectively introduces the subject holistic characteristics 489 and temporal changes for developing an accurate, stable and consistent 490 estimation in a coarse-to-refine way, and deeply explores inter-frame indices 491 relatedness for enhancing sequential indices estimation. However, the 492 compared method just conducts separate estimation on each image. 493

494

3) The proposed PV-LVNet outperforms the segmentation method, with

	non- $Reg_{grad}$	$Reg_{grad}$	
One-dimensional Metric			
$LAD_{A2C}$	0.926	0.941	
$SAD_{A2C}$	0.904	0.930	
$LAD_{A4C}$	0.904	0.932	
$SAD_{A4C}$	0.902	0.917	
Two-	dimensional N	fetric	
$Area_{A2C}$	0.897	0.907	
$Area_{A4C}$	0.918	0.940	
	Volume		
Volume	0.945	0.974	

Table 4: The inter-frames gradient regularization increases internal consistency with the ground truth. It gains higher Cronbach's  $\alpha$  (> 0.9) than being removed.

estimating 5 more indices. It efficiently explores holistic characteristics and interrelated changes among the different frames in the same subject to directly analyze LV sequence and LV biological structure for adaptively learning all cardiac indices. And U-net (Ronneberger et al., 2015) just automatically provides LV area from its segmentation while the other indices need extra interaction from the expert for apex and mitral valve plane.

In the implementation of comparison, our proposed method needs no 501 extra interaction. Our method is fed with entire echo image and does not 502 require post-processing, benefited from its robust processing ability. But 503 the other direct methods need to be performed on the the pre-handcrafted 504 region to work (Zhen et al., 2014b, 2015a, 2016; Xue et al., 2017a). The 505 segmentation method U-net is post processed as general with maximum 506 connected region extraction to improve its segmentation results for indices 507 estimation. 508

#### 509 5.6. Performance of Activation Function and Hyper Parameters Selection

Activation Function ELU vs. ReLU. As shown in Figure. 12, ELU better fits the RRU than ReLU, with the lower estimation MAE (sum of normalized multitype indices MAE). Since the activation of ELU is able to transmit not only positive value message but also negative value message among frames, which is important for stimulating the inter-frames communication. But ReLU misses the information during the negative regime because of all being forcefully pushed to zero.



Figure 12: ELU activation outperforms ReLU in the RRU with lower testing MAE and better fitting.

Hyper Parameters Setting. As shown in Figure. 13, our Hyper 517 parameters of  $\lambda_1 = 1000$ ,  $\lambda_2 = 1$  and  $\lambda_3 = 0.1$  gain the best estimation 518 accuracy compared to the other settings, with the lowest estimation MAE. 519 Defaulting  $\lambda_2$  for indices estimation as 1,  $\lambda_1$  gets the large magnitude of 1000 520 for the trade-off between the trainable tasks location and indices estimation 521 to mutually promote them;  $\lambda_3$  with the small magnitude of 0.1 balances tasks 522 training and network parameters regularization. Figure 13(a) indicates that 523 the large  $\lambda_1$  is more effective than small setting as larger ones have lower 524 rate of accuracy decay. Specifically,  $\lambda_1$  setting smaller than 1000 extremely 525 increases the estimation error. Since the unsuitable small  $\lambda_1$  decreases the 526 location supervision of LV-ROI, which leads the indices estimation in a mess. 527 The messed indices estimation then arbitrarily misleads the location through 528 the joint learning and further degrades the indices accuracy in return via the 529 chain reaction. Big is better, but not infinite. The too huge magnitude of  $\lambda_1$ 530 exceeding 1000 also has the risk of decreasing the performance. Because the 531 too huge  $\lambda_1$  weakens the effect of indices estimation in the mutual promotion, 532 so that make the indices accuracy lower. In Figure 13(b), our choice also gets 533 the best result. The big  $\lambda_3$ , as 1 and 10, have the serious problem of making 534 the learning target unclear, so that influence the learning ability. Small  $\lambda_3$ 535 keeps the learning target clear, but the too tiny  $\lambda_3$  of 0.01 and 0.001 weakens 536 the regularization on network parameters so that reduces the generalization 537 of the network and worsens the practical estimation. 538

#### 539 6. Conclusions

In this paper, we proposed the PV-LVNet for the first time achieve the direct and accurate estimation of LV multitype indices  $(LAD_{A2C}, SAD_{A2C},$  $Area_{A2C}, LAD_{A4C}, SAD_{A4C}, Area_{A4C}, Volume)$  from 2D echos of paired apical views. The PV-LVNet conducts the sufficient metrics from various dimensions (1D, 2D & 3D) and views (A2C, A4C, and union of A2C+A4C)



Figure 13: Our hyper-parameters setting gets the best estimation accuracy. (a) Influence of  $\lambda_1$  selection. (b) Influence of  $\lambda_3$  selection.

to provide a reliable comprehensive cardiac function assessment. It is 545 built based on the Res-circle Net for sequential analysis. The Res-circle 546 Net embeds both subject holistic characteristics and temporal changes 547 by combining common subject-level base among frames and interrelated 548 residuals of each frame, so that accurate and consistent location and indices 549 estimation of LVs in echo sequence are enabled. The PV-LVNet is integrated 550 of three interdependent parts for location, cropping and indices regression, as: 551 1) the LV location module utilizes AED that gives different scaled metrics on 552 different directions as the loss to suit approximate bullet shape of LV in apical 553 echos, so that robust and efficient location for indices estimation is ensured; 554 2) the Image Resampling automatically crops LV-ROI from the entire echo 555 image, so that the interference of various structures in paired views is 556 reduced; 3) by using inter-frames gradient regularization for exploring indices 557 inter-frame relatedness, the LV location module fits not only each index value 558 but also the indices evolution, so that sequential indices estimation is further 559 enhanced. The PV-LVNet reaches high accuracy on all indices estimation and 560 maintains excellent internal consistency with the ground truth, indicating its 561 great potential in clinical cardiac function evaluation. 562

#### 563 Acknowledgment

This work was supported by the Postgraduate Research & Practice Innovation Program of Jiangsu Province (No. KYCX17\_0104); the China Scholarship Council (No. 201706090248); the States Key Project of Research and Development Plan (No. 2017YFA0104302, 2017YFC0109202 and 2017YFC0107900); the National Natural Science Foundation (No. 81530060 and 61871117); and the Science and Technology Program of Guangdong (No. 2018B030333001).

### 571 References

Abdi, A.H., Luong, C., Tsang, T., Allan, G., Nouranian, S., Jue, J., Hawley,
D., et al., 2017. Automatic quality assessment of echocardiograms using

<sup>574</sup> convolutional neural networks: Feasibility on the apical four-chamber view.
<sup>575</sup> IEEE Transactions on Medical Imaging 36, 1221–1230.

Afshin, M., Ayed, I.B., Islam, A., et al., 2012. Global assessment of cardiac
 function using image statistics in mri, in: MICCAI, Springer. pp. 535–543.

Afshin, M., Ayed, I.B., Punithakumar, K., Law, M.W., et al., 2014. Regional
assessment of cardiac left ventricular myocardial function via mri statistical
features. IEEE Transactions on Medical Imaging 33, 481–494.

Ba, J., Kiros, J., Hinton, G., 2016. Layer normalization. arXiv preprint
 arXiv:1607.06450.

Carneiro, G., Nascimento, J.C., Freitas, A., 2012. The segmentation of
the left ventricle of the heart from ultrasound data using deep learning
architectures and derivative-based search methods. IEEE Transactions on
Image Processing 21, 968–982.

<sup>587</sup> Chen, H., Zheng, Y., Park, J.H., Heng, P.A., et al., 2016. Iterative multi domain regularized deep learning for anatomical structure detection and
 <sup>589</sup> segmentation from ultrasound images, in: MICCAI, Springer. pp. 487–495.

Ciampi, Q., Villari, B., 2007. Role of echocardiography in diagnosis and
 risk stratification in heart failure with left ventricular systolic dysfunction.
 Cardiovascular Ultrasound 5, 1–12.

<sup>593</sup> Clevert, D.A., et al., 2015. Fast and accurate deep network learning by <sup>594</sup> exponential linear units (elus). arXiv preprint arXiv:1511.07289.

<sup>595</sup> Cronbach, L.J., 1951. Coefficient alpha and the internal structure of tests.
 <sup>596</sup> psychometrika 16, 297–334.

<sup>597</sup> Dai, J., He, K., Sun, J., 2016. Instance-aware semantic segmentation via <sup>598</sup> multi-task network cascades, in: CVPR, pp. 3150–3158.

Debreuve, E., Barlaud, M., Aubert, G., Laurette, I., et al., 2001. Space-time
segmentation using level set active contours applied to myocardial gated
spect. IEEE Transactions on Medical Imaging 20, 643–659.

Gao, Z., Li, Y., Sun, Y., Yang, J., Xiong, H., et al., 2018. Motion tracking
of the carotid artery wall from ultrasound image sequences: a nonlinear
state-space approach. IEEE Transactions on Medical Imaging 37, 273–283.

605	Gao, Z., Xiong, H., Liu, X., Zhang, H., Ghista, D., Wu, W., Li, S., 2017.
606	Robust estimation of carotid artery wall motion using the elasticity-based
607	state-space approach. Medical Image Analysis 37, 1–21.

- Georgescu, B., Zhou, X.S., et al., 2005. Database-guided segmentation of anatomical structures with complex appearance, in: CVPR, pp. 429–436.
- Graves, A., 2012. Supervised sequence labelling, in: Supervised sequence
  labelling with recurrent neural networks. Springer, pp. 5–13.
- He, K., Zhang, X., Ren, S., Sun, J., 2016a. Deep residual learning for image
  recognition, in: CVPR, pp. 770–778.
- He, K., Zhang, X., et al., 2016b. Identity mappings in deep residual networks,
  in: European conference on computer vision, Springer. pp. 630–645.
- Jacob, G., Noble, J.A., Behrenbruch, C., Kelion, A.D., Banning, A.P.,
  2002. A shape-space-based approach to tracking myocardial borders and
  quantifying regional left-ventricular function applied in echocardiography.
  IEEE Transactions on Medical Imaging 21, 226–238.
- Jaderberg, M., Simonyan, K., et al., 2015. Spatial transformer networks, in:
  Advances in neural information processing systems, pp. 2017–2025.
- Lang, R.M., Badano, L.P., Mor-Avi, V., Afilalo, J., Armstrong, A., Ernande,
  L., Flachskampf, F.A., et al., 2015. Recommendations for cardiac chamber
  quantification by echocardiography in adults: an update from the american
  society of echocardiography and the european association of cardiovascular
  imaging. Journal of the American Society of Echocardiography 28, 1–39.
- Lang, R.M., Bierig, M., Devereux, R.B., Flachskampf, F.A., Foster, E.,
  Pellikka, P.A., Picard, M.H., et al., 2006. Recommendations for chamber
  quantification. European Journal of Echocardiography 7, 79–108.
- Lathuilière, S., Juge, R., et al., 2017. Deep mixture of linear inverse regressions applied to head-pose estimation, in: CVPR, pp. 4817–4828.
- Luo, G., Dong, S., Wang, K., Zuo, W., Cao, S., Zhang, H., 2018. Multi-views
  fusion cnn for left ventricular volumes estimation on cardiac mr images.
  IEEE Transactions on Biomedical Engineering 65, 1924 1934.

- Malladi, R., Sethian, J.A., Vemuri, B.C., 1995. Shape modeling with front
   propagation: A level set approach. IEEE Transactions on Pattern Analysis
   and Machine Intelligence 17, 158–175.
- Mo, Y., Liu, F., McIlwraith, D., Yang, G., Zhang, J., He, T., Guo, Y., 2018.
  The deep poincaré map: A novel approach for left ventricle segmentation,
  in: MICCAI, Springer. pp. 561–568.
- Nascimento, J.C., et al., 2008. Robust shape tracking with multiple models in
   ultrasound images. IEEE Transactions on Image Processing 17, 392–406.
- Oktay, O., Ferrante, E., Kamnitsas, K., Heinrich, M., Bai, W., Caballero,
  J., Cook, S.A., de Marvao, A., et al., 2018. Anatomically constrained
  neural networks (acnns): application to cardiac image enhancement and
  segmentation. IEEE Transactions on Medical Imaging 37, 384–395.
- Paragios, N., 2003. A level set approach for shape-driven segmentation and
  tracking of the left ventricle. IEEE Transactions on Medical Imaging 22,
  773–776.
- Pascual, M., Pascual, D., Soria, F., Vicente, T., Hernandez, A., Tebar, F.,
  Valdes, M., 2003. Effects of isolated obesity on systolic and diastolic left
  ventricular function. Heart 89, 1152–1156.
- Peng, P., Lekadir, K., Gooya, A., et al., 2016. A review of heart chamber
   segmentation for structural and functional analysis using cardiac magnetic
   resonance imaging. Magnetic Resonance Materials in Physics 29, 155–195.
- Pereira, S., et al., 2018. Enhancing interpretability of automatically extracted
  machine learning features: application to a rbm-random forest system on
  brain lesion segmentation. Medical image analysis 44, 228–244.
- Ravi, D., Wong, C., Deligianni, F., et al., 2017. Deep learning for health
   informatics. IEEE Journal of Biomedical and Health Informatics 21, 4–21.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks
   for biomedical image segmentation, in: MICCAI, Springer. pp. 234–241.
- Schiller, N.B., Shah, P.M., Crawford, M., DeMaria, A., Devereux, R.,
  Feigenbaum, H., Gutgesell, H., Reichek, N., et al., 1989. Recommendations
  for quantitation of the left ventricle by two-dimensional echocardiography.
  Journal of the American Society of Echocardiography 2, 358–367.

Szegedy, C., Ioffe, S., et al., 2017. Inception-v4, inception-resnet and the 667 impact of residual connections on learning, in: AAAI, pp. 4278–4284. 668 Vigneault, D.M., Xie, W., Ho, C.Y., et al., 2018.  $\omega$ -net (omega-net): Fully 669 automatic, multi-view cardiac mr detection, orientation, and segmentation 670 with deep neural networks. Medical Image Analysis 48, 95–106. 671 Wang, Z., Salah, M.B., Gu, B., Islam, A., Goela, A., Li, S., 2014. Direct 672 estimation of cardiac biventricular volumes with an adapted bayesian 673 formulation. IEEE Transactions on Biomedical Engineering 61, 1251–1260. 674 Wu, L., Cheng, J.Z., Li, S., Lei, B., Wang, T., Ni, D., 2017. Fuiqa: Fetal 675 ultrasound image quality assessment with deep convolutional networks. 676 IEEE Transactions on Cybernetics 47, 1336–1349. 677 Xingjian, S., Chen, Z., Wang, H., Yeung, D.Y., et al., 2015. Convolutional 678 lstm network: A machine learning approach for precipitation nowcasting, 679 in: Advances in neural information processing systems, pp. 802–810. 680 Xu, C., Xu, L., Gao, Z., Zhao, S., Zhang, H., Zhang, Y., et al., 2018. Direct 681 delineation of myocardial infarction without contrast agents using a joint 682 motion feature learning architecture. Medical image analysis 50, 82–94. 683 Xue, W., Brahm, G., et al. 2018. Full left ventricle quantification via deep 684 multitask relationships learning. Medical Image Analysis 43, 54–65. 685 Xue, W., Islam, A., Bhaduri, M., Li, S., 2017a. Direct multitype cardiac 686 indices estimation via joint representation and regression learning. IEEE 687 Transactions on Medical Imaging 36, 2057–2067. 688 Xue, W. Lum, A., Mercado, A., Landis, M., et al., 2017b. Full quantification 689 of left ventricle via deep multitask learning network respecting intra-and 690 inter-task relatedness, in: MICCAI, Springer. pp. 276–284. 691

 Xue, W., Nachum, I.B., Pandey, S., Warrington, J., Leung, S., Li, S., 2017c.
 Direct estimation of regional wall thicknesses via residual recurrent neural network, in: IPMI, Springer. pp. 505–516.

Yu, L., Yang, X., Chen, H., Qin, J., Heng, P.A., 2017. Volumetric convnets
with mixed residual connections for automated prostate segmentation from
3d mr images, in: AAAI, pp. 66–72.

- Zhen, X., Islam, A., Bhaduri, M., Chan, I., Li, S., 2015a. Direct and
  simultaneous four-chamber volume estimation by multi-output regression,
  in: MICCAI, Springer. pp. 669–676.
- Zhen, X., Wang, Z., Islam, A., Bhaduri, M., Chan, I., Li, S., 2016. Multiscale deep networks and regression forests for direct bi-ventricular volume
  estimation. Medical Image Analysis 30, 120–129.
- Zhen, X., Wang, Z., Islam, A., Chan, I., Li, S., 2014a. A comparative study
  of methods for cardiac ventricular volume estimation, in: Annual MeetingRadiological Society of North America (RSNA), pp. 228–244.
- Zhen, X., Wang, Z., Yu, M., Li, S., 2015b. Supervised descriptor learning
  for multi-output regression, in: Proceedings of the IEEE conference on
  computer vision and pattern recognition, pp. 1211–1218.
- Zhen, X., Wang, Z., et al., 2014b. Direct estimation of cardiac bi-ventricular
  volumes with regression forests, in: MICCAI, Springer. pp. 586–593.
- Zhen, X., Yu, M., He, X., Li, S., 2017. Multi-target regression via robust
  low-rank learning. IEEE transactions on pattern analysis and machine
  intelligence 40, 497–504.

OUTRO

715 Conflicts of interest: none

Journal Prevention