# DefCor-Net: Physics-Aware Ultrasound Deformation Correction

Zhongliang Jiang*[a],*, Yue Zhou*[a], Dongliang Cao[b], Nassir Navab[a,c]

[a]Computer Aided Medical Procedures, Technical University of Munich, Munich, Germany
[b]University of Bonn, North Rhine-Westphalia, Germany
[c]Computer Aided Medical Procedures, Johns Hopkins University, Baltimore, MD, USA

## ARTICLE INFO

## ABSTRACT

The recovery of morphologically accurate anatomical images from deformed ones is challenging in ultrasound (US) image acquisition, but crucial to accurate and consistent diagnosis, particularly in the emerging field of computer-assisted diagnosis. This article presents a novel anatomy-aware deformation correction approach based on a coarse-to-fine, multi-scale deep neural network (DefCor-Net). To achieve pixel-wise performance, DefCor-Net incorporates biomedical knowledge by estimating pixel-wise stiffness online using a U-shaped feature extractor. The deformation field is then computed using polynomial regression by integrating the measured force applied by the US probe. Based on real-time estimation of pixel-by-pixel tissue properties, the learning-based approach enables the potential for anatomy-aware deformation correction. To demonstrate the effectiveness of the proposed DefCor-Net, images recorded at multiple locations on forearms and upper arms of six volunteers are used to train and validate DefCor-Net. The results demonstrate that DefCor-Net can significantly improve the accuracy of deformation correction to recover the original geometry (Dice Coefficient: from $14.3 \pm 20.9$ to $82.6 \pm 12.1$ when the force is $6N$).
**Code:** https://github.com/KarolineZhy/DefCorNet.

## 1. Introduction

The recovery of accurate anatomical images from distorted ones is a fundamental problem in medical imaging, which contributes to the achievement of reliable diagnoses and biometric measurements in clinical practices. Regarding medical ultrasound (US), it has been widely used for examining internal organs due to its merits of low cost, high accessibility, and lack of radiation. To optimize the acoustic coupling performance, a
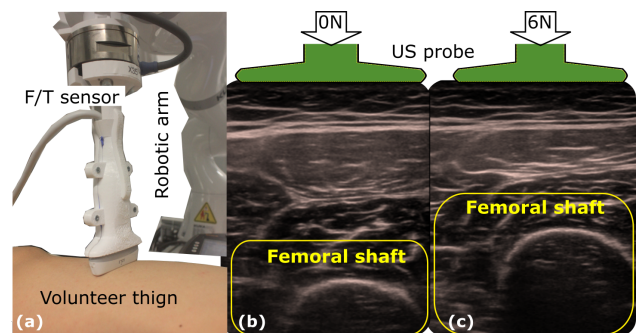
Fig. 1: Illustration of the effects of pressure-induced deformation. (a) robotic arm with a force/torque (F/T) sensor. (b) and (c) are the resulting B-mode images of the femoral shaft recorded when the contact force is $0$ $N$ and $6$ $N$, respectively.

certain pressure needs to be applied due to the inherited characteristic of US modality. Inconsistent pressure will result in distorted images with varying anatomy geometry for soft tissue, such as superficial veins (Jiang et al., 2021b), preventing the achievement of consistent diagnosis results among different clinicians. In contrast to experienced sonographers who are trained and accustomed to making diagnoses based on distorted images, image distortion will have a significant impact on the development of autonomous diagnosis systems that aim to produce standardized diagnoses (Wang et al., 2020). In addition to providing the precise geometries of soft tissues, the recovered images also reveal the precise locations of high-stiffness anatomies, e.g., bone (see Fig. 1). This is crucial for modern image-guided orthopedic surgery, which requires accurate registration between intra-operative US scans (bone surface) and pre-operative MRI or CT scans to transfer the planned trajectory (Jiang et al., 2023a,b; Wein et al., 2015; Salehi et al., 2017).

## 1.1. Robotic US Imaging

Due to their stability and accuracy, robotic US systems (RUSS) are viewed as a promising solution for obtaining high-quality images. In contrast to free-hand US imaging, the US acquisition parameters (contact force and probe orientation) can be accurately maneuvered in RUSS (Gilbertson and Anthony, 2015). To precisely control the probe orientation, Jiang et al. estimated the normal direction of unknown constraint surfaces using both force cues and US confidence map (Jiang et al., 2020a). The US confidence map provides pixel-wise quality measurements (signal strength) based on the physics of sound propagation model (Karamalis et al., 2012). To enable the autonomous adjustments of probe orientation for both linear and convex probes, Jiang et al. developed a mechanical model-based algorithm relying on accurate measurements of contact force (Jiang et al., 2020b). In addition, Ma et al. employed a depth camera to estimate the normal direction of the chest surface for screening of the lung (Ma et al., 2021). Although the camera-based approach is less accurate than the force-based method, the computation is performed more quickly. In addition to optimizing probe orientation, a number of force control schemes have been proposed for RUSS (Gilbertson and Anthony, 2015; Zettinig et al., 2017; Pierrot et al., 1999). Such approaches can maintain a given force between the probe and the contact surface, which can alleviate the inter-variation of pressure-induced deformation between the images acquired along a scanning path. Nevertheless, due to variations in the physical properties of human tissues (e.g., bone and muscle), image deformation for various anatomies in individual B-mode images remains inconsistent. To maintain the focus of this article, we refer readers to two recent survey articles for more detailed developments in the field of RUSS (Jiang et al., 2023c; Li et al., 2021).

## 1.2. US Imaging Deformation Correction

To recover uncompressed images from the deformed ones, Treece et al. combined non-rigid image-based registration and external position sensing (Treece et al., 2002). This pioneering study considered both global accuracy and local accuracy

based on the position sensor measurements and correction of the pressure-induced errors, respectively. But this approach was developed for axial deformation. To further consider the deformation in the lateral direction, Burcher et al. built an elastic model using the finite element method (FEM), where the tissue parameters were manually assigned (Burcher et al., 2001). Similarly, Flach et al. developed a correction approach in FEM based on the assumption that the tissues are homogeneous (Flack et al., 2016). To address this limitation, Sun et al. proposed a novel trajectory-based method to recover compression-free images using an empirical regressive model with respect to the contact force (Sun et al., 2010). To obtain the displacement trajectories of individual pixels under varying contact forces, a template-based image-flow technique was applied to radio frequency data. Similarly, Virga et al. used 4th-order polynomial models to depict specific pixel displacement by manually annotating 5 feature points in all frames obtained with different forces; and then, a novel graph-based inpainting approach was employed to compute 2D deformation fields for unsampled pixels (Virga et al., 2018). Although some promising results were reported on in-vivo human tissues, the lengthy computational time (186 $s$ per inpainting procedure) hinders further clinical translation.

In contrast to the aforementioned deformation correction approaches that directly link the deformation to the contact force, another folder of methods also takes tissue stiffness into consideration. Dahmani et al. first estimated the relative mechanical parameters of involved tissues based on the images themselves with a reference baseline; then, a model-based correction approach was developed to recover deformed images (Dahmani et al., 2017). Regarding the deformation recovery task, their method demonstrated better performance than the image-based method (Sun et al., 2010). But the validations were only carried out on simplified synthetic US images. In addition, Jiang et al. presented a stiffness-based deformation correction method, incorporating image pixel displacements, contact forces, and nonlinear tissue stiffness (Jiang et al., 2021b). In their study, the tissue stiffness was characterized using the probe displacements and contact forces recorded during robotic palpation, while the pixel displacements were computed using the LukasKanade algorithm (Lucas et al., 1981). To demonstrate the propagation capability of the proposed approach for various patients, the optimized deformation regression model computed at sampled positions on a stiff phantom can be successfully propagated to unsampled positions on the same stiff phantom and on a commercial soft phantom, respectively, by only substituting the estimated tissue stiffness. Nevertheless, limited by the number of coefficients of the polynomial model, the accuracy for individual pixels is sacrificed to ensure overall optimal performance for all pixels. Consequently, the performance of this approach is spatially location-dependent. In addition, the sparse optical flow technique hinders the achievement of dense deformation correction.

## 1.3. Dense Flow Estimation

The deformation correction process is often considered an inverse problem, where the deformation field is computed first

and further used for recovery. The sparse optical flow technique is often used to extract pixel displacements for computing the deformation field (Jiang et al., 2021b; Sun et al., 2010; Virga et al., 2018). Nevertheless, the sparse optical flow cannot generate a pixel-by-pixel deformation field. Thanks to the booming of deep learning, the state-of-the-art dense flow estimation approaches have achieved phenomenal success, e.g., FlowNet (Dosovitskiy et al., 2015; Ilg et al., 2017), PWC-Net (Sun et al., 2018) and Recurrent All-Pairs Field Transforms (RAFT) (Teed and Deng, 2020). Dosovitskiy *et al.* first cast optical flow estimation as a learning problem and solve it based on a standard Convolutional neural networks (CNN) architecture (Dosovitskiy et al., 2015). It was reported that FlowNet could accurately estimate the flow at a frame rate of 5 to 10 *fps*. To further improve time efficiency and flow quality, Ilg *et al.* developed FlowNet 2.0 by introducing a stacked architecture by warping the second image with intermediate optical flow (Ilg et al., 2017). To simultaneously increase the accuracy and decrease the size of a CNN model for optical flow, Sun *et al.* presented PWC-Net, a coarse-to-fine network architecture with pyramid, warping, and a cost volume (Sun et al., 2018). The results demonstrated that PWC-Net is approximately 17 times smaller in size, and two times faster than FlowNet 2.0 for inference. Instead of the coarse-to-fine architecture, Teed *et al.* proposed RAFT, which maintains and updates a single fixed flow field at high resolution (Teed and Deng, 2020). To reliably predict the flow field even when the displacement is large, a recurrent GRU-based operator is employed. In addition to its superior estimation accuracy relative to the aforementioned methods, RAFT has also been reported to possess strong cross-dataset generalization (Teed and Deng, 2020).

### 1.4. Ultrasound Elastography

US elastography (Sigrist et al., 2017) is an evolutionary method that uses US waves to assess the mechanical properties of tissues, such as stiffness. Particularly for strain elastography (also called static elastography), displacement estimation between two frames is one of the most crucial steps for high-quality elastography (Rivaz et al., 2014; Hashemi et al., 2018; Rivaz et al., 2008). From this point of view, the deformation correction procedure is an inverse problem of strain elastography. To compute the pixel-wise displacement, Rivaz *et al.* minimized cost functions incorporating the similarity of radio–frequency (RF) data intensity and displacement continuity of two frames (Rivaz et al., 2010). To improve computational effectiveness, Ashikuzzaman *et al.* considered the RF data from three consecutive frames to ensure that both spatial and temporal estimations are data-adaptive (Ashikuzzaman et al., 2019). To improve the quality of elastography, Guo *et al.* employed a partial differential equation (PDE)–based regularization algorithm to iteratively reduce noise contained in the displacement estimations (Guo et al., 2015). Considering the establishment of elastography is often time-consuming, Boctor *et al.* proposed a rapid approach to segment stiff lesions for monitoring ablative therapy based on the tissue deformation map (Boctor et al., 2006).

In addition to model-based techniques for estimating mechanical parameters from force-displacement measurements, the development of neural networks offers an additional folder of model-free solutions. Hoerig *et al.* presented a data-driven approach to compute elasticity images using artificial neural networks (Hoerig et al., 2018). This approach does not require any prior assumptions about the underlying constitutive model, but it does require internal structure information (discrete or continuous material property distributions). Tehrani *et al.* computed the displacement by training an optical flow network first using computer vision datasets and followed by a fine-tuning process to bridge the gap between the natural images and the RF data (Tehrani et al., 2022b). To enhance the accuracy of lateral strain estimation, Tehrani *et al.* proposed a novel framework using physically inspired constraint (effective Poisson's ratio) for unsupervised regularized elastography (Tehrani et al., 2022a). In order to obtain dense elastography, RF data is often used as the input rather than the B-mode images. Yet, RF data is not accessible for most of the commercial US machines used in clinical practices. Due to the fact that the estimated stiffness of various tissues is only relative, it cannot be directly adapted to recover the force-induced deformation for a single image, i.e., individual images with force information.

### 1.5. Proposed Approach

To recover zero-compression images from deformed ones, we present a novel anatomy-aware deformation correction approach using a deep neural network (DefCor-Net) for achieving accurate anatomical US imaging without reference images. To the best of the authors' knowledge, this is the first work introducing a learning-based approach to correct the pressure-induced US deformation with pixel-wise and real-time performance. The DefCor-Net mainly consists of a stiffness estimation module (SEM) and a deformation field estimation module (DFM). The SEM first computes the local stiffness online from an input image using a U-shape network unit (Ronneberger et al., 2015). Then, the patient-specific tissue stiffness computed based on robotic palpation is used to update the computed local stiffness map, thereby enabling the capability to generalize across patients. In addition, the DFM is designed to predict pixel-wise deformation fields based on the physics of the interaction between the current force and the computed stiffness map. The DefCor-Net is implemented in a coarse-to-fine manner to capture both low-resolution and high-resolution characteristics, thereby facilitating the computation of the pixel-wise deformation field, even for relatively large deformations. The DefCor-Net is trained in an end-to-end fashion, where the ground truth of the deformation field is represented by the dense optical flow map computed using RAFT (Teed and Deng, 2020). Since RAFT can only estimate the displacement based on two frames, they cannot be used to correct individual deformed images without reference frames. Contrary to conventional polynomial model-based methods (Sun et al., 2010; Jiang et al., 2021b; Dahmani et al., 2017), DefCor-Net has a large number of trainable parameters, i.e., $7,782,936$ in total, which enables the achievement of pixel-wise correction results due to the dense stiffness map. The experiments were carried out on the in-vivo data recorded from six volunteers, including multiple sampling locations on both forearms and upper arms.

## 2. Material and Data Preparation

### 2.1. Hardware Setup

Our robotic US system consists of a lightweight robot, a US device, and an accurate F/T sensor. The F/T sensor was directly mounted on the robotic flange in order to precisely measure the contact force. The US probe was attached on the other side of the F/T sensor using a custom-designed holder.

#### 2.1.1. Robotic Manipulator

In this study, a redundant robotic arm (KUKA LBR iiwa 7 R800, KUKA Roboter GmbH, Germany) with seven joints was used. In each joint, a torque sensor is incorporated, allowing the implementation of compliant force control to address safety concerns and acoustic coupling issues. To control the robotic arm, an open-source interface[1] developed in our lab was used. This interface allows direct communication between the low-level Sunrise applications and the Robot Operating System (ROS) framework. To guarantee real-time performance, the robotic status and the control commands from user applications were updated in 500 *Hz* and 100 *Hz*, respectively.

#### 2.1.2. Tactile Sensing

To accurately measure the dynamic contact force, an external F/T sensor (Gamma, ATI IndustrialAutomation, USA) was used. The measured force can be accessed using a data acquisition device (FTD-DAQ-USB6361, National Instruments, USA). To synchronize the force information and US images, the force was published to the ROS master with precise timestamps at a rate of 100 *Hz*.

#### 2.1.3. Ultrasound System

In this study, B-mode images were acquired using an ACUSON Juniper US machine (Siemens Healthineers, Germany) with a linear probe (12L3, Siemens Healthineers, Germany). The B-mode images were updated in 58 *fps*. To access the B-mode images from the US machine, a frame grabber (Epiphan DVI2USB 3.0, Epiphan Vision, Canada) was used to take screenshots in 30 *fps*. To synchronize the force information and US images, the US images stream was published to the ROS master in real-time. Then the images and force data can be synchronized using OpenIGTLink[2], where the time tolerance was 12 *ms*. The US images were recorded using the default configuration for bone. The detailed parameters are listed as follows: imaging depth: 45 *mm*, focus: 20 *mm*, frequency: 6.7 *MHz*, and brightness: 78 *dB*.

### 2.2. Compliant Force Control Architecture

In order to estimate the patient-specific mechanical properties and record paired data (contact force, probe pose, and US images), "robotic palpation" was performed by increasing and decreasing the applied force. To this end, a compliant force controller with an input of desired force was used (Hennersperger

et al., 2016). Compared to a position controller, the compliant mode can provide a safer interaction avoiding excessive force. The Cartesian compliant controller is defined as Eq. (1).

$$\tau = \mathbf{J}^T[\mathbf{F}_d + \mathbf{K}_m e + \mathbf{D}\dot{e} + \mathbf{M}\ddot{e}] \qquad (1)$$

where $\tau \in \mathbf{R}^{7\times1}$ is the required joint torque to realize the compliant force controller, $J^T \in \mathbf{R}^{7\times6}$ is the transposed Jacobian matrix, $e \in \mathbf{R}^{6\times1} = (x_d - x_c)$ is the pose error (position and orientation) between the current pose $x_c$ and the desired pose $x_d$ in Cartesian space, $\mathbf{F}_d \in \mathbf{R}^{6\times1}$ is the desired F/T applied at the tool center point (TCP), $\mathbf{K}_m \in \mathbf{R}^{6\times6}$, $\mathbf{D} \in \mathbf{R}^{6\times6}$ and $\mathbf{M} \in \mathbf{R}^{6\times6}$ are diagonal matrices of stiffness, damping and inertia values in 6 degree of freedom (DoF), respectively. According to (Hennersperger et al., 2016), the stiffness is suggested to be in the range of $[125, 500]N/m$ for different tissues. In addition, 25 *N* was set as the maximum contact force in the low-level safety configuration in the robotic cabinet, which will be triggered once a force larger than 25 *N* occurs.

### 2.3. Robot-Assisted Data Collection

To record in-vivo US images, six healthy volunteers were employed. For each volunteer, we evenly selected four sampling points on the forearm from the elbow joint toward the wrist joint in 150 *mm* distance. In addition, two sampling points were selected on the upper arm with a distance of 40 *mm*. Each robotic palpation consists of one cycle of pressing and releasing, in which the applied force was gradually increased from zero to 15 *N* and then decreased to zero using the compliant force controller. Each palpation took around 30 *s*, which will produce approximately 900 sequential B-mode images. At each sampling location, robotic palpation were carried out twice to collect in-vivo data. In the end, we have 72 image sets on both forearms and upper arms from six volunteers. To train the DefCor-Net, 51 ( 70%) and 8 ( 10%) image sets were used for training and validation, respectively. The remaining 13 unseen image sets were used for testing.

Based on the recorded position and measured force at the probe tip, the global stiffness of human tissue can be estimated (see Section 3.1). The average stiffness values (±SD) of the forearm and upper arm over sampling points were $1.80 \pm 0.48$ *N/mm* and $0.78 \pm 0.03$ *N/mm*, respectively. The maximum compression levels that happened at the contact surface for the upper arm and forearm were $12.5 \pm 1.14$ *mm* and $7.1 \pm 0.8$ *mm*, respectively. Both results indicate that the mechanical properties of the forearm and upper arm are significantly different. During the process of data acquisition, volunteers were instructed to extend their arms naturally on a flat table.

### 2.4. Optical Flow-Based Deformation Field Generation

The estimation of the deformation field is the key to the deformation correction procedure. Benefiting from the boom of learning-based dense optical flow estimation technologies, this study explores the feasibility of using the computed dense flow map to represent the deformation field between two images. Since the deformation field between a deformed image and its corresponding uncompressed image can be estimated offline,

---

[1]https://github.com/IFL-CAMP/iiwa_stack
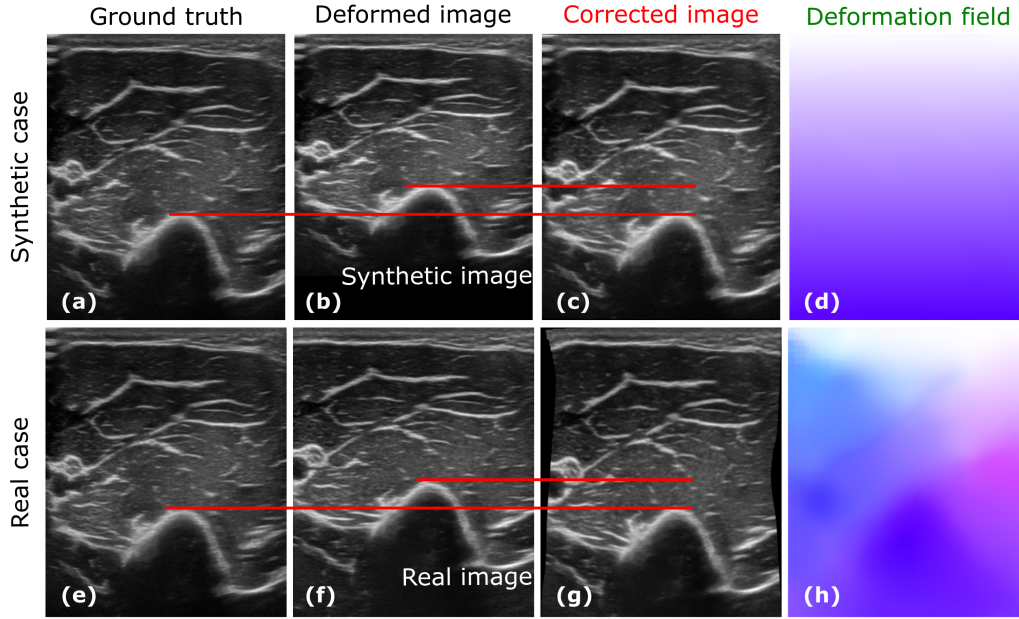[2]http://openigtlink.org/

Fig. 2: Illustration of the computed deformation field and the corrected results. The top and bottom rows represent the representative results on synthetic data and real data acquired when the contact force is 7.5 *N*. For the synthetic case, the deformed images is generated by evenly moving the individual pixels with different displacements, where the top row is moved zero, and the bottom row is moved 70 pixels. The deformation fields are computed based on the deformed images and ground truth images. The corrected images in both cases are close to the ground truth. The Normalized Cross Correlation (NCC) similarity between two images is (a,b)=0.31 (c,b)=0.99; (e,f)=0.34, (g,f)=0.86.

the accuracy of flow estimation becomes the most crucial issue. Due to the superiority in terms of accuracy and generalization for different datasets, RAFT with a well-trained model (Sintel)[3] was employed to compute the dense optical flow field representing the deformation field. Considering RAFT is originally designed to estimate image flow based on two consecutive images, an intermediate frame $\mathbb{I}_m$ between the uncompressed image $\mathbb{I}_1$ and the current image $\mathbb{I}_2$ is used to guarantee the flow estimation accuracy by avoiding too large displacement. The final flow $f_{2,1}$ between $\mathbb{I}_1$ and $\mathbb{I}_2$ is calculated as Eq. (2).

$$f_{1,2} = f_{1,m} + warp(f_{m,2}, f_{1,m}) \qquad (2)$$

where $f_{i,j}$ represents the image flow estimated from image $\mathbb{I}_i$ to $\mathbb{I}_j$ and $warp(\mathbb{I}_j, f_{i,j}) \rightarrow \mathbb{I}_i$ represents the warping operation to recover deformed image $\mathbb{I}_j$ using known flow map $f_{i,j}$. For each pixel location $P_k$ in $\mathbb{I}_i$ and its corresponding pixel in $\mathbb{I}_j$ are satisfied the following equation.

$$c^j[P_k + f_{i,j}(P_k)] = c^i(P_k)$$
$$s.t. \quad [0,0]^T \leq P_k + f_{i,j}(P_k) \leq [W, H]^T \qquad (3)$$

where $c^i$ and $c^j$ are the intensity values of a pixel position on $\mathbb{I}_i$ and $\mathbb{I}_j$, respectively. $W$ and $H$ are the width and height of images in terms of the pixel. The bilinear interpolation is used to implement the warping operation as (Teed and Deng, 2020). To intuitively demonstrate the accuracy of the generated deformation field and the correction performance of deformed images,

two representative results from synthetic and real deformed images are depicted in Fig 2.

Due to the force-induced compression, anatomical structures on US images are moved towards the top of the image. To validate whether the RAFT can correctly compute the deformation field, a synthetic image was generated by evenly moving pixel lines (horizontal) at different distances. In our case, the top line (contact surface) was kept constant, while the bottom line was moved 70 pixels (18.2% of the image height) in the vertical direction. The computed deformation field is demonstrated in Fig. 2 (d), which conforms the expectation. No displacement happened at the top line, while maximum displacement occurred at the bottom line. In addition, the horizontal color intensity is uniform, while the vertical color intensity increases evenly from the top line to the bottom line. After applying the warping operation for the deformed image [Fig. 2 (b)] using the computed deformation field [Fig. 2 (d)], the corrected imaged is visualized in Fig. 2 (c). The position of the anatomical structures on the corrected images is very close to the ground truth image obtained when the contact force is zero (a large amount of gel was filled between the probe and contact surface) in terms of the Normalized Cross Correlation (NCC) Coefficient (improved from 0.31 to 0.99), while the pixel intensity is more bright. Besides the synthetic data, we shown another representative example of a real deformed image obtained when the force is 7.5 *N* in Fig. 2. The computed deformation field [Fig. 2 (h)] is more complex, which represents individual pixel displacement direction and magnitude. After the warping process, the corrected image [Fig. 2 (g)] also becomes very close to the ground truth (NCC similarity improved from 0.34 to 0.86).

Considering the US deformation is highly nonlinear due to

---

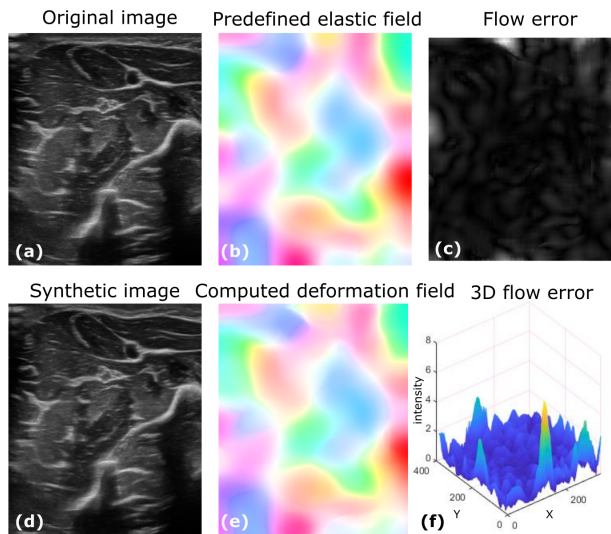[3]https://github.com/princeton-vl/RAFT

Fig. 3: The performance of the estimation of a representative nonlinear elastic distortion field. (a) the original image. (b) the predetermined elastic deformation field using elastic distortion (Simard et al., 2003). (d) the distorted image of (a) the original image by applying (b) the elastic deformation. (e) the computed deformation field using the modified RAFT. (c) and (f) are 2D and 3D visualization of the end-to-end point error (Euclidean distance) between (b) the predefined elastic field and (e) the computed deformation field.

the variations of tissue properties, we further validate the effectiveness of the modified RAFT for estimating non-linear deformation fields. To generate realistic distorted images, the elastic distortion (Simard et al., 2003) was randomly determined, and then the synthetic images can be obtained by applying the elastic distortion field to the original images. An intuitive example is demonstrated in Fig. 3. Based on the original and synthetic images, the deformation field can be calculated using RAFT [Fig. 3 (e)]. To quantitatively compare the known elastic field and computed deformation field, the end-to-end point error (Euclidean distance) is computed for individual pixel locations and visualized in the flow error map [see Fig. 3 (c)]. To make the flow error map more visible, the pixel intensity was considered as $[0, 1]$ for visualization. The locations with large error values ($> 1$ pixel) were considered as one in Fig. 3 (c). To reflect the error at each pixel location, the 3D error image is shown in Fig. 3 (f). The maximum distance error is 6.5 pixels, while the mean error is only 0.5 pixel. Furthermore, the NCC similarity value between the original and corrected images is enhanced to 0.98 from 0.71 between the original and synthetic images. Therefore, we consider the modified RAFT can provide decent deformation field ground truth for the development of learning-based deformation correction methods.

## 2.5. Data Augmentation

To improve the diversity of datasets, some classic data augmentation methods, e.g., random scaling, rotation, and translation, are widely used in the field of computer vision. However, US images have different characteristics (e.g., physical sound attenuation and tissue-based distortion) from natural camera images. Regarding the force-induced deformed images, the displacements in the vertical and lateral directions are coupled. To keep this characteristic, only two augmentation methods were applied here. To reflect the anatomy-aware ability, B-mode images were horizontally flipped at random. Besides, a crop mask ($W_c \times H$, where $W_c \leq W$) was randomly initialized to obtain the cropped images, where the image height $H$ was kept as the original images. The reason for not changing $H$ is that the deformation is physically transmitted and accumulated from the skin surface. We consider these two augmentation methods can benefit to achieving anatomy-aware correction performance by diversifying the anatomy's location.

## 3. Method

In this section, we present a novel CNN-based DefCor-Net organized in a coarse-to-fine manner for estimating the pixel-wise displacement maps of individual deformed US images acquired under various contact forces. By incorporating the concept of pixel-wise stiffness maps and patient-specific global stiffness obtained by robotic palpation, the proposed framework can predict the anatomy-aware deformation field. The proposed DefCor-Net consisting of three layers is depicted in Fig. 4. To capture potential large deformation, the original B-mode images ($W \times H$) are down-sampled twice ($\frac{W}{2} \times \frac{H}{2}$ and $\frac{W}{4} \times \frac{H}{4}$) as the inputs of the first two layers. In each layer, the networks can be grouped into SEM and DFM modules, respectively. Detailed explanations of individual parts are provided in the following subsections.

### 3.1. Pixel-Wise Stiffness Estimation Module

Considering real scenarios, human tissues are non-homogeneous and vary largely from one patient to another one. Thereby, in order to estimate the deformation field only based on the deformed B-mode image without reference, it is important to compute the pixel-wise local stiffness. This means that the system needs to be aware of the anatomies' position on the resulting B-mode images. To this end, a U-shape feature extraction unit is initially employed to extract the initial stiffness map from a single deformed image. The initial stiffness map can indicate the pixel-wise nonlinear tissue properties. To further consider the patient-specific effects, a constant global stiffness is computed for each subject via robotic palpation and fed to a stiffness update module.

#### 3.1.1. U-Shape Feature Extraction

In order to effectively and robustly extract feature representations from images, Simonyan *et al.* proposed VGG-Net using very deep CNN structures (Simonyan and Zisserman, 2014). Specific to B-mode images, Baumgartner *et al.* proposed a CNN-based SonoNet (Baumgartner et al., 2017). The U-Net (Ronneberger et al., 2015) is another successful representative, and its variations have been widely used for segmenting detailed structures like vessel boundaries (Jiang et al., 2021a; Prevost et al., 2018). In this work, we also use a U-shape network to extract the pixel-by-pixel stiffness representations from B-mode images.
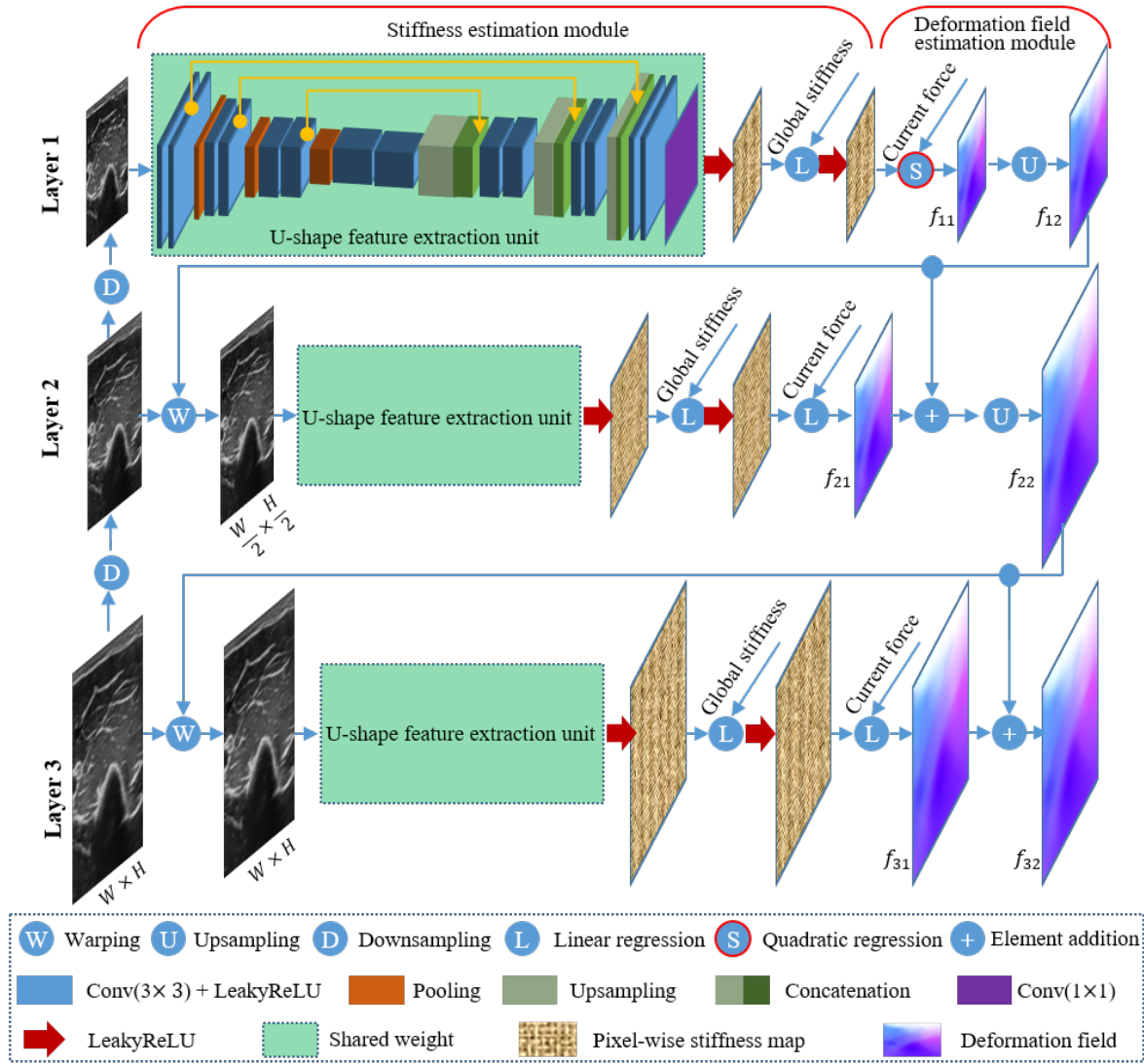
Fig. 4: Illustration of the proposed DefCor-Net in a coarse-to-fine manner.

The details of the used U-shape architecture are shown in Fig. 4, which is similar to the 3D-Unet (Çiçek et al., 2016) by replacing the 3D operations (e.g., 3D convolutions and 3D pooling) with 2D operations. Considering the intermediate stiffness representations could be negative, the LeakyReLU function is employed to guarantee the backpropagation process of the network instead of ReLU activation. The size of the outputs is kept the same as the input images. Due to the coarse-to-fine manner, the U-shape feature extractor is reused three times in all three layers to generate pixel-wise stiffness cues. Based on the experimental results, the weight-sharing strategy is implemented to reduce the number of parameters while maintaining good performance. The image-based stiffness map $\mathbb{K}$ can be expressed as Eq. (4).

$$\mathbb{K}_{us} = \text{LeakyReLU}\left(\text{UNET}(\mathbb{I}_{US})\right) \qquad (4)$$

where UNET means the U-shape feature extraction unit.

### 3.1.2. Patient-Specific Stiffness Map Updating

In addition to the variation in stiffness between tissues of the same patient, the variation between patients also affects the accuracy of the stiffness representation. To consider the patient-specific factor in the proposed framework, a set of stiffness updating processes are further applied on the computed $\mathbb{K}_{us}$. To compute the global stiffness of individual patients, the robotic palpation was executed to record the real-time contact force in the direction of the probe centerline and the displacement of the probe tip in the force direction $\lambda_z$. Five representative data recorded from volunteers' forearms have been described in Fig. 5. It can be seen from the recorded data that linear regression ($F_c = c_2\lambda_z + c_1$) is able to properly describe the relationship between the force and displacement. The R-square value of the computed regression models are 0.96, 0.97, 0.99, 0.99, and 0.98, respectively. The high R-square fitness value means that the global stiffness $K_g$ for individual patients can be represented by the gradient $c_2$ of the optimized linear regression models. It is noteworthy that the computed global stiffness may not be desired for other parts of the same volunteer. It can only

represent the tissue stiffness variable in the neighboring area of the sampling position.
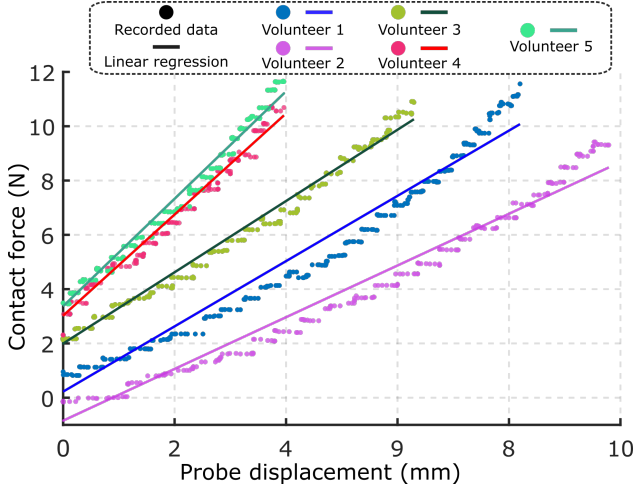


Fig. 5: Representative stiffness fitting results on the forearm of five volunteers. The R-square values of the optimized linear regressions are 0.96, 0.97, 0.99, 0.99, and 0.98, respectively.

To stabilize the deviation of individual global stiffness $K_g$ for the same objects of different patients, $K_g$ of unseen patients is normalized using the Z-score approach before being used for adjusting the computed $\mathbb{K}_{us}$ as Eq. (5).

$$k_g^n = \frac{k_g - \mu_g}{\delta_g} \qquad (5)$$

where $\mu_g = \sum_{i=1}^{N} \frac{1}{N} k_g^i$ and $\delta_g = \sqrt{\frac{\sum k_g^i - \mu^2}{N}}$ are the average and standard deviation of the recorded data, respectively. Then, the patient-specific stiffness map $\mathbb{K}_{us}^g$ is computed as Eq. (6).

$$\mathbb{K}_{us}^g = \text{LeakyReLU}\left(c_1 \mathbb{K}_{us} + c_2 k_g^n + c_3\right) \qquad (6)$$

### 3.2. Deformation Field Estimation Module

To recover a deformed image, a displacement field is required. To estimate the deformation field, polynomial regressions are often employed (Sun et al., 2010; Virga et al., 2018; Treece et al., 2002). After comparing the recovery performance using different regression models, the quadratic model was found to give a better result than linear regression, while the performance was not improved significantly for a higher order polynomial model (Sun et al., 2010). Thereby, quadratic regressions were empirically used in the first layer to characterize the displacement field with respect to the current force $F_c$ and the pixel-wise stiffness map $\mathbb{K}_{us}^g$ as Eq. (7).

$$\begin{bmatrix} \mathbf{D}_x(x,y) \\ \mathbf{D}_y(x,y) \end{bmatrix} = F_c \begin{bmatrix} \mathbf{K}_x \\ \mathbf{K}_y \end{bmatrix} \mathbf{M}_v\left(\mathbb{K}_{us}^g(x,y)\right) \qquad (7)$$

where $(x, y)$ represents a pixel position on the image and $\mathbf{D}_x$ and $\mathbf{D}_y$ represent the pixel displacements in X (lateral) and Y (axial) directions, respectively. $\mathbf{M}_v(h) = \left[h^2, h, 1\right]^T$ ($h = \mathbb{K}_{us}^g(x, y)$), and $\mathbf{K}_i \in R^{1 \times 3}$ ($i = x$ or $y$) is the unknown parameters matrix of the quadratic regressions need to be optimized. To take the

advantage of residual idea for the case with relative large deformation, the computed deformation field map $f_{11}$ is up-sampled as Eq. (8).

$$f_{i2} = UP(f_{i1}, c_s) \otimes c_s \qquad (8)$$

where $i = 1$ or $2$ is the layer identification, $UP()$ is upsampling using bilinear interpolation, $c_s = 2$ is the scaling coefficient, $\otimes$ is element-wise multiplication operation.

Compared with the methods that compute the deformation field only based on the contact force (Sun et al., 2010; Virga et al., 2018), the inclusion of tissue's mechanical properties theoretically allows the application of an optimized model for trained data on unseen objects by updating their specific stiffness. Preliminary work on this idea had been reported in (Jiang et al., 2021b), where the authors successfully applied an optimized regression model of a stiff phantom to a soft phantom with significantly different stiffness values. Although they considered the dynamic representation of the human tissue with respect to the applied force, the stiffness used there was considered homogeneous across all pixels in the image obtained under a specific force. This means anatomies' position changes on the deformed images (e.g., bone surface move from the left to the right side of the resulting images) cannot be properly addressed. To address this challenge, the pixel-wise stiffness map $\mathbb{K}_{us}^g$ is used in Eq. (7) instead of a single value to accurately estimate the deformation field, which is expected to be independent of the anatomy's position on deformed images.

### 3.3. Coarse-to-Fine Deformation Field Extraction

In consideration of the possibility of relative large deformation, DefCor-Net is organized in a coarse-to-fine manner. After introducing the stiffness and deformation field estimation modules, this subsection elaborates on the contributions of individual layers and the whole workflow among the three layers.

Layer 1 is designed to deal with the overall large deformation using the downsampled inputs (size is ($\frac{W}{4} \times \frac{H}{4}$)). Compared with the other two layers, a quadratic regress (Eq. (7)) rather than a linear model is used in layer 1 to capture the non-linear properties of the deformation. After computing the deformation field in the first layer $f_{11}$, an up-sampling process based on the bilinear technique is implemented to obtain the field $f_{12}$ as the base deformation field of layer 2. To capture the deformation characteristics in the resolution of ($\frac{W}{2} \times \frac{H}{2}$), a warp operation is applied on the input images $\mathbb{I}_{in}$ as $warp(\mathbb{I}_{in}, f_{12})$. The warped results are further fed to the U-shape feature extraction unit with shared parameters to create an intermediate stiffness-related mask. Similar to layer 1, the patient-specific stiffness and the contact force $F_c$ are integrated consecutively using Eq. (6) and (7) to compute the intermediate deformation field $f_{21}$. Based on the residual idea, the main deformation characteristics are accounted for in layer 1, while the residual pixel displacement is refined using images with higher resolutions in layers 2 and 3. Thus, a linear function ($\mathbf{M}_v'(h) = [h, 1]^T$) is empirically used in Eq. (7) for the latter two layers. Then, the residual displacement field $f_{21}$ is added to the based flow $f_{12}$ using element-wise addition. Repeating the aforementioned steps in layer 3 and the final displacement field $f_{32}$ can be achieved.

To train the network, the $f_{32}$ will be compared with the ground truth computed using RAFT.

### 3.4. Loss Function

The proposed DefCor-Net is organized as a supervised learning problem. The ground truth of the deformation field $f_{gt}$ is computed as Section 2.4. The L1 loss is used to guide the learning process. Since the resolution of the input images are different, the sub-size $f_{gt}$ is computed as Eq. (9).

$$f_{gt}^{\ell} = Down(f_{gt}, c_s^{2-\ell}) \otimes (c_s^{-1})^{2-\ell} \tag{9}$$

where $\ell = 1$ or $2$ is the layer iteration, $Down()$ is downsampling operation based on bilinear interpolation. The L1 loss of all three layers is computed as Eq. (10).

$$\mathcal{L}_{l1} = \sum_{\ell} \left\| f_{gt}^{\ell} - f_{pred}^{\ell} \right\|_1 \tag{10}$$

where $f_{pred} = \{f_{11}, f_{21}, f_{32}\}$ is the predicted deformation field.

Besides $\mathcal{L}_{l1}$, the geometry continuity of the displayed anatomies should be further considered, which is important for clinical diagnosis, in particular for the future development of advanced computer-assisted diagnosis. To encourage such behavior, a smoothness loss $\mathcal{L}_{sm}$ proposed by (Jonschkowski et al., 2020) is modified as Eq. (11).

$$\mathcal{L}_{sm} = \sum_k \exp\left(-\lambda_x^k \left(\frac{\partial I}{\partial x}\right)^2\right) f\left(\frac{\partial^k f_{32}}{\partial x^k}\right)$$
$$+ \underbrace{\exp\left(-\lambda_y^k \left(\frac{\partial I}{\partial y}\right)^2\right)}_{geometry} \underbrace{f\left(\frac{\partial^k f_{32}}{\partial y^k}\right)}_{smoothness} \tag{11}$$

where $k = 1, 2$ represents the first- and second-order term. The smoothness term mainly consists of the partial derivatives of predicted deformation map $f_{32}$ with respect to $x$ and $y$, respectively. By minimizing the derivatives, the intensity of $f_{32}$ becomes smooth. To avoid zero gradient, $f(m) = \sqrt{m^2 + \epsilon^2}$. However, the smooth term will decrease anatomies' geometry accuracy. To address this, the geometry term is used as edge coefficient. The boundary can be identified by computing the gradient image of the original inputs with respect to $x$ and $y$, respectively. The geometry term is used to ensure the boundaries of different anatomies are not smoothed out.

Then, the final training loss is organized as Eq. (12).

$$\mathcal{L} = \lambda_1 \mathcal{L}_{l1} + \lambda_2 \mathcal{L}_{sm} \tag{12}$$

where $\lambda_1$ and $\lambda_2$ are the hyperparameters to combine the two different loss. Based on the output performance, $\lambda_1$ and $\lambda_2$ are empirically set to 1 and 10, respectively.

## 4. Results

To highlight the contributions of the proposed approach, we provide a summary of the key characteristics of existing methods (see TABLE 1) that aim to correct force-induced US deformation. In addition to the metric of axial and lateral deformation correction ability, the terms of force cue and tissue

property indicate whether the contact force and tissue mechanical property are used for US image deformation correction. The approaches without considering tissue properties will limit their performance in generalization capability among patients in theory. The metric of anatomy-aware ability indicates whether the method can recover a deformed image correctly regardless of the relative position of the anatomy of interest in US images, such as a bone surface located on the left or right side in the deformed image. The anatomy-aware ability is important in real scenarios because the position of target tissues varies significantly according to the placement of the probe. Most of the previous work cannot achieve anatomy-aware performance because they usually directly characterize the deformation with respect to the applied force using polynomial regression models. Due to the absence of tissue properties, the optimized regression model is the overall optimal interpreter for seen data, while its inference does not account for the positional variation of anatomy in US images. FEM-based approaches (Burcher et al., 2001) could be one of the potential solutions. Nevertheless, the need for precise estimation of patient-specific tissue parameters remains an open challenge in the community. In addition, FEM is time-consuming and cannot be used in real-time based on current developments.

To address this challenge, we use advanced deep learning to implicitly estimate the pixel-wised stiffness characteristics form input US images. By forcing DefCor-Net to explicitly infer nonlinear tissue stiffness maps, the presented supervised learning framework can compute the deformation field based on both applied force and pixel-wise stiffness maps. Consequently, anatomy-aware performance is implicitly achieved by being aware of the changes in stiffness value distributions in the inferred stiffness maps. Due to the fact that the stiffness value changes should be continuous, the intermediate variables (i.e., stiffness maps) are qualitatively validated whether being consistent with the underlying physics in Section-4.2. In addition, quantitative evaluations of the proposed DefCor-Net are performed on in-vivo limb US images. The validation of the anatomy-aware ability has been demonstrated in Section-4.3. The quantitative evaluation of the proposed deformation correction method based on local bone is summarized in Section 4.4. Finally, the target localization accuracy and dense flow error over the entire images are computed and discussed in Section 4.5.

### 4.1. Implementation Details

The coarse-to-fine DefCor-Net was implemented using PyTorch framework. The training and evaluation configurations were maintained based on OpenMMLab[4]. The DefCor-Net model was optimized using the ADAM optimizer (Kingma and Ba, 2014) on a workstation with a GPU of GeForce GTX 1080. The learning rate was 0.0001 for 50$k$ epochs. To guarantee the generability of the trained model for different patients, a small batch size (four) was used in this work. The sweeps used for testing are not used in the training data set. The image size

---

[4] https://github.com/open-mmlab

Table 1: Summary of the Key Characteristics of Different Approaches

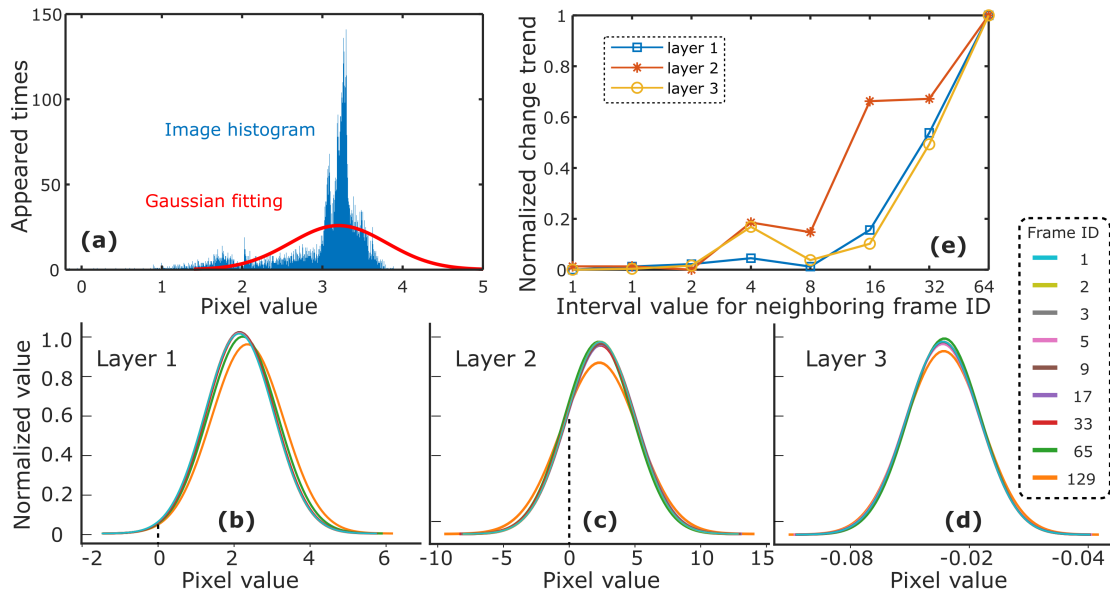| Approach | Axial def. | Lateral def. | Force cue | Tissue property | Anatomy-aware | In-vivo test |
|---|---|---|---|---|---|---|
| Pheiffer *et al.* (Pheiffer and Miga, 2015) | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Dahmani *et al.* (Dahmani et al., 2017) | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ |
| Sun *et al.* (Sun et al., 2010) | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Virga *et al.* (Virga et al., 2018) | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ |
| Burcher *et al.* (Burcher et al., 2001) | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ |
| Jiang *et al.* (Jiang et al., 2021b) | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Proposed DefCor-Net | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |



Fig. 6: Consistency of intermediate stiffness inferences with respect to the similarity levels of input images. (a) Representative images histogram and the fitted Gaussian distribution of an inferred stiffness map in layer 1. (b), (c), (d) are the fitted Gaussian distributions of the selected frames in layers 1, 2, and 3, respectively. (e) the normalized difference of the exponentially increased intervals of neighboring frames.

is $320 \times 384$. To improve the diversity of the training dataset for achieving anatomy-aware effects, the images were randomly cropped from the original size to $256 \times 384$. The training process took around eight hours. In addition, the processing time (Mean±SD) for generating the deformation map was $34 \pm 10 \, ms$ over 844 samples, which means that the approach is able to be used for the application expecting real-time performance.

### 4.2. Consistency of Intermediate Stiffness Inferences

The proposed DefCor-Net is designed by implicitly considering biomedical knowledge. It can be seen from Fig. 4 that the inferred pixel-wise stiffness map is an important intermediate variable. Since it is not possible to obtain the ground truth of pixel-wised stiffness of the input B-mode images, we cannot directly compare the inferred intermediate stiffness-related terms. However, stiffness is a physical (analog) variable, which has the characteristic that the changes in stiffness are continuous in the ideal case.

To validate the consistency among the automatically generated intermediate stiffness maps, a set of unseen images from the same US sweep were fed to the well-trained DefCor-Net. The inferred stiffness maps (the latter ones) in three layers were

stored. To assess the difference between different images, the results of nine representative frames have been shown in Fig. 6. The frame ID are exponentially increased ($1, 1 + 2^0, 1 + 2^1, \cdots, 1 + 2^7$). To quantitatively assess the consistency for neighboring sampled frames, the image histogram of all obtained stiffness maps is drawn, where the number of bins was fixed to 1000 in this work. Then, the Gaussian fitting was repeatedly carried out to fit all histograms. Representative historiography and the corresponding Gaussian fitting result of the inferred stiffness map obtained in layer 1 are shown in Fig. 6 (a). All 27 fitted Gaussian distributions of layer 1, layer 2 and layer 3 are depicted in Figs. 6 (b), (c) and (d), respectively. To intuitively display the difference between the selected nine frames, we computed the overlap area $A_{ov}$ between the fitted Gaussian distributions obtained from two neighboring sampled frames. Since the cumulative distribution function for a normal distribution is always one when $x = +\infty$, the difference between two fitted Gaussian distributions is represented by $1 - A_{ov}$. The normalized changes between neighboring frames for all three layers are depicted in Fig. 6 (e).

Due to the coarse-to-fine structure, the non-linear stiffness properties are mainly modeled in the first layer (quadratic regression). Thus, the inferred stiffness map of layer 1 is corre-
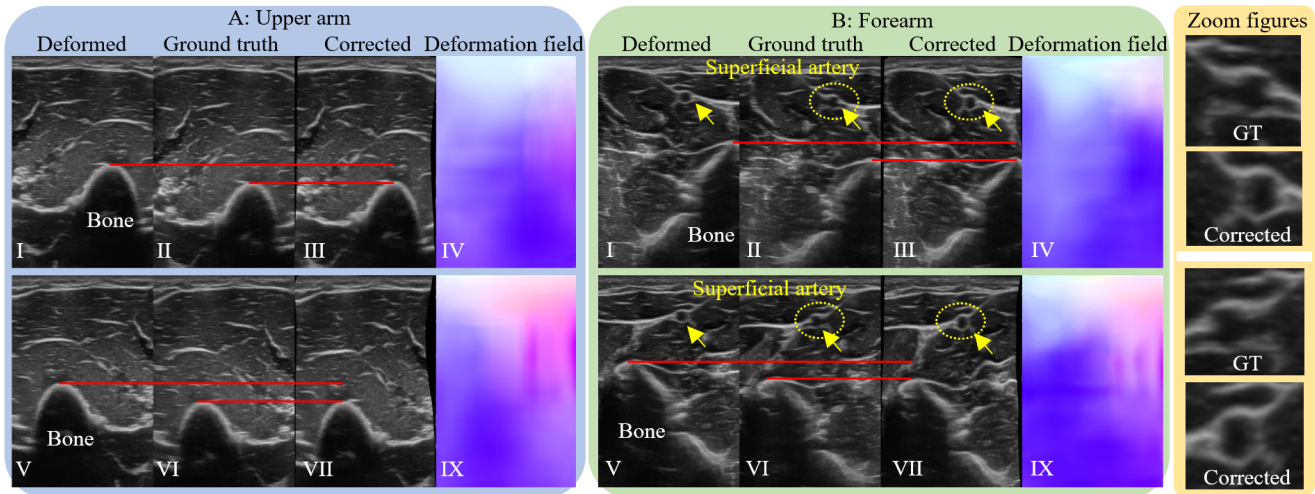
Fig. 7: Illustration of the anatomy-aware characteristic. The upper arm and forearm US images are recorded when the force is 7.0 *N* and 6.9 *N*, respectively. The superficial artery on uncompressed and corrected images is zoomed in and depicted in the right column.

sponding to the real stiffness, while the ones of the other two layers are considered residual compensations. It can be seen from Figs. 6 (a) and (b) that the pixel values of the inferred stiffness map from layer 1 are mainly distributed around two, which is consistent with the prior knowledge that the stiffness term is positive ideally. It is noteworthy that there is no such restriction on the residual compensations obtained in the other two layers. The exponentially increased intervals between neighboring frames are used to represent the different levels of similarity between the input B-mode images. Due to the characteristic of the analog variable, the stiffness changes should be positively corresponding to the similarity of the input images. This phenomenon is demonstrated in Fig. 6 (e). The changes are very small ($< 0.02$ when the frame interval is small than two) for the automatically inferred stiffness maps in all three layers. Such changes become larger when the interval increases. There is an exceptional case when the frame interval is four, where the changes are larger than the results obtained when the interval increased to eight (layer 1: 0.04 vs 0.01; layer 2: 0.18 vs 0.15; layer 3: 0.16 vs 0.04). But regarding the stiffness term (layer 1), 0.04 still can be seen as no significant difference between the neighboring frames. The experimental results demonstrate that the difference between the inferred stiffness maps is consistent with the similarity levels of the input images.

### 4.3. Validation of Anatomy-Aware Characteristic

To validate the performance of the advanced anatomy-aware ability, experiments were carried out on unseen sweeps on both the upper arm and forearm. Since the stiffness of the bone structure is far higher than the surrounding soft tissues, the bone structure will be preserved during the robotic palpation. Thereby, the pixel position of the bone on the images was used to intuitively demonstrate the changes caused by the deformation and the quality of the corrected results in Fig. 7. For each case (upper arm or forearm), two representative deformed images with significantly different anatomy (bone) positions were

used to demonstrate the anatomy-aware ability of the proposed DefCor-Net.

It can be seen from Fig. 7 that the trained DefCor-Net can properly recover the deformed images in all four representative cases for different tissues. The deformed images are significantly different from the ground truth, while the corrected ones become very close to the ground truth. As shown in Fig. 7, the position of dark blue parts on the deformed fields is consistent with the positions of the bones on their corresponding inputs. Since the bone structure is not deformed, the displacement of the pixels covered by the bone surface will be consistent in both direction and amplitude. This phenomenon has been demonstrated by the quasi-homogeneous intensity values (dark blue area) on the computed deformation field. In addition, it is noteworthy that the superficial artery geometry on the corrected images is more clear than the one displayed on the ground truth. This phenomenon is caused due to the inherited characteristic of US modality, which requires a certain force to achieve optimal acoustic coupling performance. To avoid imaging deformation, the ground truth is obtained by applying zero pressure ideally. This results in the less clear boundary of the superficial artery on the shallow layer of the ground truth than of corrected images. Two representative comparisons are demonstrated in the far right column in Fig. 7. The vascular boundaries on corrected images are more clear and more complete than the ones on the natively uncompressed images.

### 4.4. Evaluation of Deformation Correction on Bone structure

#### 4.4.1. Bone Segmentation

Due to the large stiffness against the surrounding soft tissues, the bone structure will not be deformed during scans. Taking advantage of this characteristic, we used the same mask block to partially annotate the bone structure on the deformed images, correct images, and ground truth from the whole palpation sweep (see Fig. 8). The top outline of the block needs

Table 2: Performance of the Corrected Results using Different Approaches (Dice Coefficient)

| Contact Pressure | Deformed image | Linear Scaling | DefCor-Net-C | DefCor-Net-Lin | DefCor-Net |
|---|---|---|---|---|---|
| 1 N | 84.9 ± 15.5 | 86.6 ± 15.1 | 89.0 ± 12.8 | 90.2 ± 9.4 | **95.9±3.3** |
| 2 N | 61.0 ± 28.2 | 76.2 ± 14.8 | 78.6 ± 19.8 | 88.6 ± 4.4 | **92.4±4.8** |
| 3 N | 43.2 ± 29.1 | 71.5 ± 10.0 | 55.7 ± 25.6 | 85.4 ± 6.6 | **91.1±7.5** |
| 4 N | 30.4 ± 30.2 | 74.4 ± 18.0 | 47.6 ± 28.5 | **87.8±6.8** | 87.8±7.1 |
| 5 N | 19.5 ± 25.1 | 65.3 ± 15.6 | 40.3 ± 31.8 | 84.1 ± 11.5 | **87.8±12.2** |
| 6 N | 14.3 ± 20.9 | 57.2 ± 22.8 | 32.1 ± 28.9 | 81.5 ± 7.7 | **82.6±12.1** |

to be consistent with the bone surface, while the bottom outline of the block was manually determined. Compared with the approach to annotating all images individually, a generic mask block can maximally minimize the inconsistency introduced by manual annotations. The use of a block area rather than the bone outline targets the same issue. A representative annotated result on a set of deformed, corrected images and ground truth is depicted in Fig. 8.
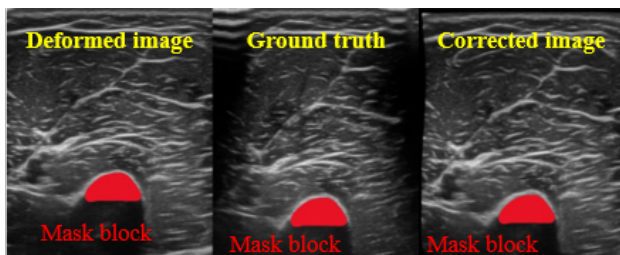


Fig. 8: Illustration of a manually annotated mask block on deformed images, ground truth, and corrected images. The deformed image was recorded on a volunteer's upper arm when the contact force was 4 *N*.

### 4.4.2. Quantitative Evaluation

To provide quantitative assessment, the popular metric of Dice Coefficient $d_c = 100 \times \frac{2 |G \cap S|}{|G| + |S|}$ was computed between the ground truth and deformed/corrected images. *G* represents the number of pixels in the mask block of the ground truth image, whereas *S* is the number of pixels in the mask block of the deformed/corrected images. To demonstrate the superiority of the proposed DefCor-Net, it was compared with other approaches, including a linear deformed correction approach and various modifications of DefCor-Net. To fully validate the correction performance on the cases with different levels of deformation, the deformed images recorded with different forces ([1, 6]) were sampled. To count the variations between different tissues (with significantly different stiffness), seven and three sweeps from the upper arm and forearm, respectively, were used for the validation. The detailed results are depicted in TABLE 2.

Besides the Dice Coefficient computed between the ground truth and the deformed images, TABLE 2 also presents the Dice Coefficient computed between the ground truth and the corrected images obtained by various approaches. The linear scaling approach assumes that the compression is linearly increased as the image depth. The DefCor-Net-C represents the network replacing the U-shape feature extraction unit with a classic 2D U-Net (Ronneberger et al., 2015). The U-shape unit used in the final DefCor-Net is inspired by the 3D U-Net (Çiçek et al.,

2016), which has three layers with fewer parameters compared with the classic U-Net. The numbers of parameters in DefCor-Net-C and DefCor-Net are around 31 million and 8 million, respectively. The DefCor-Net-Lin represents the network replacing the quadratic regression with a linear regression in the first layer.

It can be seen from TABLE 2 that the Dice Coefficient of the recorded deformation images is dramatically decreased from 84 to 14 when the contact force increased from 1 *N* to 6 *N*. This means that the bone's position on the US images is significantly changed when the contact force is increased. Therefore, to provide accurate anatomical geometry and position for diagnosis, it is necessary to apply a certain deformation correction. Based on TABLE 2, the highest Dice Coefficient is achieved when the force is 1 *N* for all individual cases. This is because the deformation is tiny and relatively easy to be recovered when the force is small (1 *N*). Regarding the case with a large force, the Dice Coefficients achieved by DefCor-Net are much higher than the ones achieved by linear and DefCor-Net-C (82.6 vs 32.1 when force is 6 *N*). Besides, it is noteworthy that DefCor-Net-Lin achieved the closest results to DefCor-Net, while still worse than DefCor-Net in most cases. This is because the deformation is not linearly with respect to the applied force and tissue stiffness. This is also consistent with the results achieved by (Sun et al., 2010) that a high-order polynomial model outperforms the linear model for correcting force-induced deformation. The results demonstrate that the proposed DefCor-Net can effectively recover the deformed images.

### 4.5. Evaluation of Deformation Correction on Overall Images

Due to the bone structure often located in the lower part of images, it can only reflect the performance of the correction method over local areas. To comprehensively validate the performance of the proposed method, we further calculated the target localization accuracy (TLA) based on biomedical interfaces of different tissues (e.g., fat and muscle) on US images. Besides, a direct comparison was carried out between the computed deformation field of DefCor-Net and the ground truth calculated by RAFT. These two validations can intuitively demonstrate the performance of the deformation correction method over the whole view of B-mode images.

### 4.5.1. Target Localization Accuracy

US images often exhibit distinct curves caused by the boundaries between various tissues. Such interfaces located at different parts of the image can be used as good biomarkers to
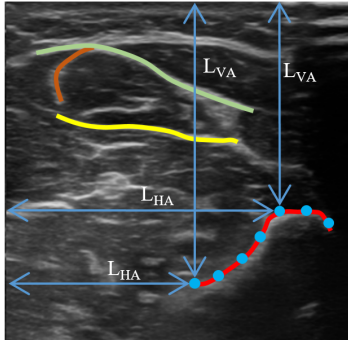
Fig. 9: Illustration of tissue interface annotations on B-mode images. The four curves with clear ending points on both sides are carefully annotated. The blue dots are the sampling points generated on one selected interface.

Table 3: Results in Terms of Target Localization Accuracy [Mean(SD)]

| Pressure | Horizontal accuracy (HA) | | Vertical accuracy (VA) | |
|---|---|---|---|---|
| | Deformed | Corrected | Deformed | Corrected |
| 1 N | 0.95±0.04 | 0.96±0.03 | 0.98±0.01 | 0.98±0.01 |
| 2 N | 0.89±0.09 | 0.96±0.03 | 0.94±0.07 | 0.97±0.03 |
| 3 N | 0.85±0.13 | 0.94±0.04 | 0.91±0.07 | 0.97±0.03 |
| 4 N | 0.80±0.15 | 0.91±0.08 | 0.88±0.09 | 0.95±0.04 |
| 5 N | 0.79±0.16 | 0.90±0.08 | 0.85±0.11 | 0.93±0.07 |
| 6 N | 0.77±0.20 | 0.91±0.08 | 0.84±0.09 | 0.95±0.03 |

calculate the TLA between the uncompressed images and deformed/corrected images. A representative illustration is depicted in Fig. 9, where we carefully annotated the selected four distinct interfaces. Each selected interface should have clear ending points at both sides to ensure its completeness on a set of corresponding uncompressed, deformed, and corrected images. To compute TLA, we first evenly (every 5 pixels in the horizontal axis) sampled the measurement points from all interfaces. Then horizontal and vertical distances ($L_{HA}$ and $L_{VA}$) at each sampling point were computed (see Fig. 9). Then TLA at each sampling point can be calculated as $HA = \frac{\|L_{HA}^i - L_{HA}^{gt}\|}{L_{HA}^{gt}}$ and $VA = \frac{\|L_{VA}^i - L_{VA}^{gt}\|}{L_{VA}^{gt}}$, where $i$ is the deformed or corrected images. The average TLA (HA and VA) computed on unseen scans from six volunteers are summarized in TABLE 3.

It can be seen from TABLE 3 that both HA and VA decrease when the contact force increases. The HA decreases from $0.98 \pm 0.04$ to $0.77 \pm 0.20$, while the VA decreases from $0.98 \pm 0.01$ to $0.84 \pm 0.09$. The TLAs are significantly enhanced after applying the proposed deformation correction process. Both HA and VA are maintained over 0.9 in different force levels. The HA and VA are increased from $0.77 \pm 0.20$ to $0.91 \pm 0.08$ and from $0.84 \pm 0.09$ to $0.95 \pm 0.03$, respectively, when the force is 6 *N*. It is also noteworthy that SD is significantly reduced after deformation correction in most cases. This implicitly demonstrates the corrected performance is stable and robust among all sampling points on the images.

### 4.5.2. Dense Error Map of Deformation Field

To further provide dense error maps, we directly computed the end-to-end point error (Euclidean distance) between the DefCor-Net-calculated deformation field and the RAFT-calculated flow maps. Fig. 10 depicts representative dense error maps obtained under varying forces. To enhance visibility, the intensity values in all error maps were scaled by the same factor (15). The average pixel-by-pixel errors are $0.88 \pm 0.21$ pixels, $1.80 \pm 0.46$ pixels, $1.87 \pm 0.71$ pixels, $1.92 \pm 1.02$ pixels, $2.47 \pm 1.18$ pixels and $2.73 \pm 1.59$ pixels when the forces are 1 *N*, 2 *N*, 3 *N*, 4 *N*, 5 *N* and 6 *N*, respectively. To ensure the deformation correction performance is robust to the pixel locations, we further counted the number of pixels (percentage $Per_{10}$, $Per_{15}$ and $Per_{20}$) with the error larger than 10 pixels (1.1 *mm*), 15 pixels (1.76 *mm*) and 20 pixels (2.34 *mm*). Regarding Figs. 10 (a), (b) and (c), all errors are less than 10 pixels. In Figs. 10 (e), (f) and (g), $Per_{10}$ increases to 2.5%, 2.0% and 2.2%, while $Per_{20}$ are only 0, 0.6% and 0.1%, respectively. This indicates that the computed deformation field using the proposed DefCor-Net is close to the ground truth flow. Therefore, we consider the proposed pixel-wise deformation correction can properly recover deformed US images caused by external pressure.
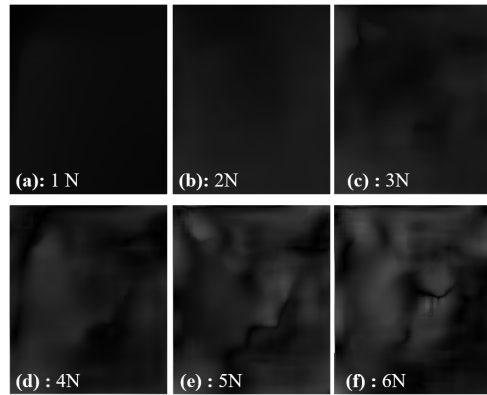


Fig. 10: Illustration of the error maps between the estimated deformation field and the ground truth flow computed using RAFT.

## 5. Discussion

The proposed novel DefCor-Net works boldly in our setup. By explicitly incorporating pixel-wise stiffness information into the loop, the estimated deformation field can be generated implicitly with respect to the location of the anatomies in real time. But, there are some limitations of the current approach that are worth noting. First, this work only aims to correct the force-induced deformation, where the contact force is the main source for the resulting deformation in images. There is still a gap in propagating the proposed method to internal organs such as the liver, in which the deformation is a result of the coupling between external pressure and physiological motions such as respiration. Second, this study only considered the data recorded from healthy volunteers. Although it works boldly in our setup, the applicability to patients with abnormal tissues, such as peripheral arterial disease, has not been validated yet. In the future, further studies should include patients' data to prove clinical transferability.

## 6. Conclusion

This work first proposes a learning-based approach DerCor-Net for achieving accurate anatomical images from deformed B-mode images caused by the inevitable pressure. The novel DerCor-Net is organized in a coarse-to-fine fashion and incorporates the biomedical knowledge between force, stiffness, and contact force, which theoretically enables the extrapolation of applying a well-trained model to different patients. Benefited from the deep neural network, over 8 million trainable parameters allow DerCor-Net to compute pixel-wise stiffness based on the input images and the deformation field ground truth obtained using an optical flow network RAFT (Teed and Deng, 2020). Besides improving the correcting accuracy, the online estimated pixel-wise stiffness map also enables the anatomy-aware characteristics for recovering the deformed image. This feature is still missed in the previous related work (Sun et al., 2010; Dahmani et al., 2017; Jiang et al., 2021b; Virga et al., 2018), however, is very important in real scenarios where the sonographer cannot always maintain the target anatomy at the same position on the US view. To validate the proposed approach, DefCor-Net was trained using the data recorded from six volunteers, both forearms and upper arms, and further validated on 13 unseen sweeps. The proposed DefCor-Net can significantly improve the accuracy of tissue geometry from deformed images (Dice Coefficient: $82.6 \pm 12.1$ vs $14.3 \pm 20.9$ when force is 6 *N*). The average Dice Coefficient of the proposed method for the images with different levels of deformation achieves 89.6.

We hope this work can reduce intra- and inter-operator variations of US acquisition by compensating for the force-induced deformations, thereby resulting in consistent and accurate diagnosis in clinical practice. This is particularly important for the further development of advanced computer-assisted diagnosis programs because the sensing and perception ability of machines are still significantly inferior to those of human operators. The future work includes integrating the proposed method to improve image-guided orthopedic surgery and exploring the way to extend this work in 3D by considering 3D continuity.

## Appendix A. Optical Flow Color Coding



Fig. A.1: Illustration of the flow color encoding used in this article.

In order to properly visualize computed optical flows, Baker *et al.* applied different colors to represent pixel displacement vectors (Baker et al., 2011). The color coding is shown in Fig. A.1, in which the hue represents the flow directions and the intensity indicates the relative magnitude. The displacement of each pixel in the figure is the vector from the square center to this pixel, while the center pixel (white color) does not move.

## Declaration of Competing Interest

The authors report no conflicts of interest.

## Acknowledgments

## References

Ashikuzzaman, M., Gauthier, C.J., Rivaz, H., 2019. Global ultrasound elastography in spatial and temporal domains. IEEE transactions on ultrasonics, ferroelectrics, and frequency control 66, 876–887.

Baker, S., Scharstein, D., Lewis, J., Roth, S., Black, M.J., Szeliski, R., 2011. A database and evaluation methodology for optical flow. International journal of computer vision 92, 1–31.

Baumgartner, C.F., Kamnitsas, K., Matthew, J., Fletcher, T.P., Smith, S., Koch, L.M., Kainz, B., Rueckert, D., 2017. Sononet: real-time detection and localisation of fetal standard scan planes in freehand ultrasound. IEEE transactions on medical imaging 36, 2204–2215.

Boctor, E., Deoliveira, M., Choti, M., Ghanem, R., Taylor, R., Hager, G., Fichtinger, G., 2006. Ultrasound monitoring of tissue ablation via deformation model and shape priors, in: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2006: 9th International Conference, Copenhagen, Denmark, October 1-6, 2006. Proceedings, Part II 9, Springer. pp. 405–412.

Burcher, M.R., Han, L., Noble, J.A., 2001. Deformation correction in ultrasound images using contact force measurements, in: Proceedings IEEE Workshop on Mathematical Methods in Biomedical Image Analysis (MMBIA 2001), IEEE. pp. 63–70.

Çiçek, Ö., Abdulkadir, A., Lienkamp, S.S., Brox, T., Ronneberger, O., 2016. 3d u-net: learning dense volumetric segmentation from sparse annotation, in: International conference on medical image computing and computer-assisted intervention, Springer. pp. 424–432.

Dahmani, J., Petit, Y., Laporte, C., 2017. Model-based correction of ultrasound image deformations due to probe pressure, in: Medical Imaging 2017: Image Processing, International Society for Optics and Photonics. p. 101331D.

Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Van Der Smagt, P., Cremers, D., Brox, T., 2015. Flownet: Learning optical flow with convolutional networks, in: Proceedings of the IEEE international conference on computer vision, pp. 2758–2766.

Flack, B., Makhinya, M., Goksel, O., 2016. Model-based compensation of tissue deformation during data acquisition for interpolative ultrasound simulation, in: 2016 IEEE International Symposium on Biomedical Imaging (ISBI), IEEE. pp. 502–505.

Gilbertson, M.W., Anthony, B.W., 2015. Force and position control system for freehand ultrasound. IEEE Transactions on Robotics 31, 835–849.

Guo, L., Xu, Y., Xu, Z., Jiang, J., 2015. A pde-based regularization algorithm toward reducing speckle tracking noise: A feasibility study for ultrasound breast elastography. Ultrasonic imaging 37, 277–293.

Hashemi, H.S., Fallone, S., Boily, M., Towers, A., Kilgour, R.D., Rivaz, H., 2018. Assessment of mechanical properties of tissue in breast cancer-related lymphedema using ultrasound elastography. IEEE transactions on ultrasonics, ferroelectrics, and frequency control 66, 541–550.

Hennersperger, C., Fuerst, B., Virga, S., Zettinig, O., Frisch, B., Neff, T., Navab, N., 2016. Towards MRI-based autonomous robotic us acquisitions: a first feasibility study. IEEE Trans. Med. Imaging 36, 538–548.

Hoerig, C., Ghaboussi, J., Insana, M.F., 2018. Data-driven elasticity imaging using cartesian neural network constitutive models and the autoprogressive method. IEEE transactions on medical imaging 38, 1150–1160.

Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., Brox, T., 2017. Flownet 2.0: Evolution of optical flow estimation with deep networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2462–2470.

Jiang, Z., Grimm, M., Zhou, M., Esteban, J., et al., 2020a. Automatic normal positioning of robotic us probe based only on confidence map optimization and force measurement. IEEE Robotics and Automation Letters 5, 1342–1349.

Jiang, Z., Grimm, M., Zhou, M., Hu, Y., Esteban, J., Navab, N., 2020b. Automatic force-based probe positioning for precise robotic ultrasound acquisition. IEEE Transactions on Industrial Electronics 68, 11200–11211.

Jiang, Z., Li, C., Li, X., Navab, N., 2023a. Thoracic cartilage ultrasound-ct registration using dense skeleton graph. arXiv preprint arXiv:2307.03800 .

Jiang, Z., Li, X., Zhang, C., Bi, Y., Stechele, W., Navab, N., 2023b. Skeleton graph-based ultrasound-ct non-rigid registration. IEEE Robotics and Automation Letters .

Jiang, Z., Li, Z., Grimm, M., Mingchuan, Z., Esposito, M., Wolfgang, W., Stechele, W., Wendler, T., Navab, N., 2021a. Autonomous robotic screening of tubular structures based only on real-time ultrasound imaging feedback. IEEE Transactions on Industrial Electronics .

Jiang, Z., Salcudean, S.E., Navab, N., 2023c. Robotic ultrasound imaging: State-of-the-art and future perspectives. Medical Image Analysis , 102878.

Jiang, Z., Zhou, Y., Bi, Y., Zhou, M., Wendler, T., Navab, N., 2021b. Deformation-aware robotic 3d ultrasound. IEEE Robotics and Automation Letters 6, 7675–7682.

Jonschkowski, R., Stone, A., Barron, J.T., Gordon, A., Konolige, K., Angelova, A., 2020. What matters in unsupervised optical flow, in: European Conference on Computer Vision, Springer. pp. 557–572.

Karamalis, A., Wein, W., Klein, T., Navab, N., 2012. Ultrasound confidence maps using random walks. Med. Image Anal. 16, 1101–1112.

Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 .

Li, K., Xu, Y., Meng, M.Q.H., 2021. An overview of systems and techniques for autonomous robotic ultrasound acquisitions. IEEE Transactions on Medical Robotics and Bionics 3, 510–524.

Lucas, B.D., Kanade, T., et al., 1981. An iterative image registration technique with an application to stereo vision. volume 81. Vancouver.

Ma, X., Zhang, Z., Zhang, H.K., 2021. Autonomous scanning target localization for robotic lung ultrasound imaging, in: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE. pp. 9467–9474.

Pheiffer, T.S., Miga, M.I., 2015. Toward a generic real-time compression correction framework for tracked ultrasound. Int. J. Comput. Assist. Radiol. Surg. 10, 1777–1792.

Pierrot, F., Dombre, E., Dégoulange, E., Urbain, L., Caron, P., Boudet, S., Gariépy, J., Mégnien, J.L., 1999. Hippocrate: A safe robot arm for medical applications with force feedback. Medical Image Analysis 3, 285–300.

Prevost, R., Salehi, M., Jagoda, S., Kumar, N., Sprung, J., Ladikos, A., Bauer, R., Zettinig, O., Wein, W., 2018. 3D freehand ultrasound without external tracking using deep learning. Med. Image Anal. 48, 187–202.

Rivaz, H., Boctor, E., Foroughi, P., Zellars, R., Fichtinger, G., Hager, G., 2008. Ultrasound elastography: a dynamic programming approach. IEEE transactions on medical imaging 27, 1373–1377.

Rivaz, H., Boctor, E.M., Choti, M.A., Hager, G.D., 2010. Real-time regularized ultrasound elastography. IEEE transactions on medical imaging 30, 928–945.

Rivaz, H., Boctor, E.M., Choti, M.A., Hager, G.D., 2014. Ultrasound elastography using multiple images. Medical image analysis 18, 314–329.

Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer. pp. 234–241.

Salehi, M., Prevost, R., Moctezuma, J.L., Navab, N., Wein, W., 2017. Precise ultrasound bone registration with learning-based segmentation and speed of sound calibration, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 682–690.

Sigrist, R.M., Liau, J., El Kaffas, A., Chammas, M.C., Willmann, J.K., 2017. Ultrasound elastography: review of techniques and clinical applications. Theranostics 7, 1303.

Simard, P., Steinkraus, D., Platt, J., 2003. Best practices for convolutional neural networks applied to visual document analysis, in: Seventh International Conference on Document Analysis and Recognition, 2003. Proceedings., IEEE. pp. 958–963.

Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 .

Sun, D., Yang, X., Liu, M.Y., Kautz, J., 2018. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8934–8943.

Sun, S.Y., Anthony, B.W., Gilbertson, M.W., 2010. Trajectory-based deformation correction in ultrasound images, in: Medical Imaging 2010: Ultrasonic Imaging, Tomography, and Therapy, International Society for Optics and Photonics. p. 76290A.

Teed, Z., Deng, J., 2020. Raft: Recurrent all-pairs field transforms for optical flow, in: European conference on computer vision, Springer. pp. 402–419.

Tehrani, A.K., Ashikuzzaman, M., Rivaz, H., 2022a. Lateral strain imaging using self-supervised and physically inspired constraints in unsupervised regularized elastography. IEEE Transactions on Medical Imaging .

Tehrani, A.K., Sharifzadeh, M., Boctor, E., Rivaz, H., 2022b. Bi-directional semi-supervised training of convolutional neural networks for ultrasound elastography displacement estimation. IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control 69, 1181–1190.

Treece, G.M., Prager, R.W., Gee, A.H., Berman, L., 2002. Correction of probe pressure artifacts in freehand 3d ultrasound. Med. Image Anal. 6, 199–214.

Virga, S., Göbl, R., Baust, M., Navab, N., Hennersperger, C., 2018. Use the force: deformation correction in robotic 3D ultrasound. Int. J. Comput. Assist. Radiol. Surg. 13, 619–627.

Wang, L., Zhang, L., Zhu, M., Qi, X., Yi, Z., 2020. Automatic diagnosis for thyroid nodules in ultrasound images by deep neural networks. Medical image analysis 61, 101665.

Wein, W., Karamalis, A., Baumgartner, A., Navab, N., 2015. Automatic bone detection and soft tissue aware ultrasound–ct registration for computer-aided orthopedic surgery. Int. J. Comput. Assist. Radiol. Surg. 10, 971–979.

Zettinig, O., Frisch, B., Virga, S., Esposito, M., Rienmüller, A., Meyer, B., Hennersperger, C., Ryang, Y.M., Navab, N., 2017. 3D ultrasound registration-based visual servoing for neurosurgical navigation. Int. J. Comput. Assist. Radiol. Surg. 12, 1607–1619.