# Force Estimation from OCT Volumes using 3D CNNs

Nils Gessert[1] · Jens Beringhoff[1] ·
Christoph Otte[1] · Alexander Schlaefer[1]

**Abstract** *Purpose* Estimating the interaction forces of instruments and tissue is of interest, particularly to provide haptic feedback during robot assisted minimally invasive interventions. Different approaches based on external and integrated force sensors have been proposed. These are hampered by friction, sensor size, and sterilizability. We investigate a novel approach to estimate the force vector directly from optical coherence tomography image volumes.

*Methods* We introduce a novel Siamese 3D CNN architecture. The network takes an undeformed reference volume and a deformed sample volume as an input and outputs the three components of the force vector. We employ a deep residual architecture with bottlenecks for increased efficiency. We compare the Siamese approach to methods using difference volumes and two-dimensional projections. Data was generated using a robotic setup to obtain ground truth force vectors for silicon tissue phantoms as well as porcine tissue.

*Results* Our method achieves a mean average error of $7.7 \pm 4.3$ mN when estimating the force vector. Our novel Siamese 3D CNN architecture outperforms single-path methods that achieve a mean average error of $11.59 \pm 6.70$ mN. Moreover, the use of volume data leads to significantly higher performance compared to processing only surface information which achieves a mean average error of $24.38 \pm 22.00$ mN. Based on the tissue dataset, our methods shows good generalization in between different subjects.

*Conclusions* We propose a novel image-based force estimation method using optical coherence tomography. We illustrate that capturing the deformation of subsurface structures substantially improves force estimation. Our approach can provide accurate force estimates in surgical setups when using intraoperative optical coherence tomography.

**Keywords** Force Estimation · OCT · 3D CNN · Siamese CNN

✉ Nils Gessert, E-mail: nils.gessert@tuhh.de, Tel.: +49 (0)40 42878 3389, https://orcid.org/0000-0001-6325-5092

[1] Hamburg University of Technology, Schwarzenbergstraße 95, 21073 Hamburg

# 1 Introduction

Robot-assisted minimally invasive surgery has become increasingly popular as it addresses various shortcomings of conventional minimally invasive surgery (MIS) [1]. Robotic systems allow for motion scaling, tremor compensation and more degrees of freedom for tool movement which improves precision and reduces physical trauma [2]. However, these systems often lack force feedback [3], which would be helpful to control the instrument-tissue interaction during surgery. Typically, haptic feedback is generated on the patient side with haptic sensors, such as force sensors [4]. The information is fed back to a haptic interface that delivers the information to the human operator, e.g., as vibrotactile or kinesthetic feedback [5]. One of the key challenges of generating reliable haptic feedback is accurate sensing of the forces at the patient [6]. Lack of haptic feedback may lead to complications, increased completion time or severe injuries [7]. Although various approaches to realize force feedback have been proposed, the problem is still considered an open research challenge [8].

One approach is to directly incorporate force sensing devices into the robotic setup [9]. The devices can be placed inside or outside of the patient. If the device is placed outside the patient, e.g., in between tool and robot, only indirect measurement is possible. In addition to the forces at the tool tip, forces acting on the tool, e.g., due to friction, are measured which cannot be separated [10]. When placing the device closer to the tool-tissue interaction point, e.g., inside the tool tip, problems such as sterilization and biocompatibility arise [11].

Due to these shortcomings, vision-based force estimation procedures have been proposed. First methods used a deformable template matching method to derive the force acting on an elstic object [12]. Similar methods relying on mechanical deformation models have been studied for MIS scenarios [13,14,15]. Also, learning forces from image information using neural networks has been proposed [16]. More recent approaches have combined template matching and machine learning models [17,18]. Recently, recurrent neural networks (RNN) have been proposed to learn forces based on deformation tracked over time [19,20]. The tissue surface is reconstructed from stereoscopic camera images and features representing surface deformation are defined. Then, the RNN is trained in a supervised fashion using ground-truth labels from a force sensor. Moreover, force estimation using optical coherence tomography (OCT) as an imaging modality has been proposed [21]. Surfaces are extracted from OCT volumes and forces are estimated based on surface deformation.

So far, these approaches only rely on visual information capturing the surface deformation. Most proposed methods make use of stereoscopic cameras [22] which are limited to the observation of surface deformations without imaging capabilities for inner tissue structures. Moreover, the proposed systems fuse features from visual deformation with robotic position feedback which limits the trained model to a specific robotic setup. Usually, retraining for a new setup is possible. However, quick adaptations might be difficult.

We introduce a force estimation approach using volumetric OCT data and a novel Siamese 3D CNN architecture. The 3D CNN directly processes OCT volumes and outputs the three components of the related force vector.

OCT can provide volumetric images with a resolution of a few micrometers which allows capturing the inner structure of tissue. In contrast to surface-based
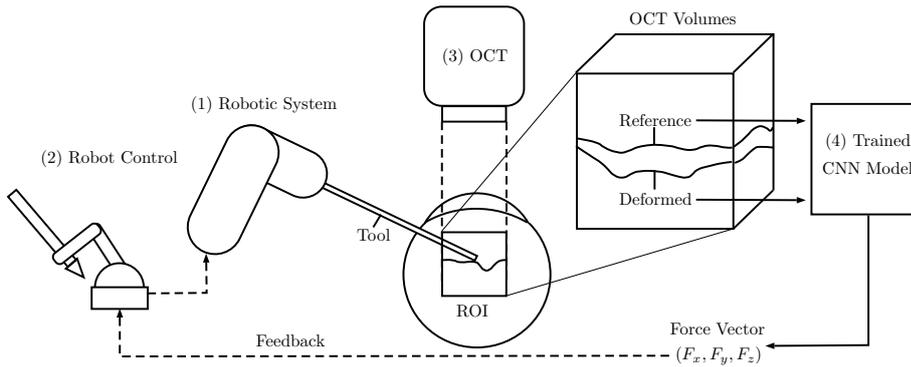
**Fig. 1:** The concept of our force estimation approach. A robotic system (1) that is controlled by a surgeon (2) performs actions in an ROI that lead to tissue-tool interactions. An intraoperative OCT device (3) repeatedly captures high resolution volumes of the ROI. The volumes are paired with a reference scan that are fed into a trained CNN (4) that performs inference in order to predict the force acting on the tissue. The force is fed back to the surgeon in order to provide visual or haptic feedback.

methods, the volumetric OCT image can also reflect tissue compression. Therefore, it is reasonable to expect OCT to provide a richer signal space with more information on subsurface deformation for accurate force estimation. We design our 3D CNN as a Siamese architecture that simultaneously processes the undeformed reference volume and the deformed sample volume to infer the force vector.

In order to evaluate our method, we acquire data for a tissue phantom. We compare the approach to methods based on difference volumes. Moreover, we compare our method to surface-based force estimation approaches by only using surfaces extracted from OCT volumes with our Siamese CNN.

Last, we validate our approach with a large dataset using porcine tissue. We use tissue from different subjects and vary the visible region for each tissue sample in order to show the robustness and generalization of our approach.

Subsequently, we describe our methodology in detail and then we provide and discuss our results. The results indicate that a precise force estimation is feasible.

## 2 Methods and Materials

### 2.1 Force Estimation with OCT

The overall concept of using OCT image volumes to estimate the force during instrument tissue interaction is shown in Figure 1. Surgery is performed with a robotic device that is remote-controlled by a surgeon. An OCT scan head is used to capture image volumes of the region of interest (ROI). First, a reference volume without tissue deformation is acquired. Then, when the tool deforms the tissue, more volumes are acquired. Both the reference and the current sample volume are fed into a trained CNN which predicts the force vector that acts on the tissue. Note, that the same vector with opposite direction acts on the tool and hence the predicted force is fed back to the surgeon and provides haptic or visual feedback.
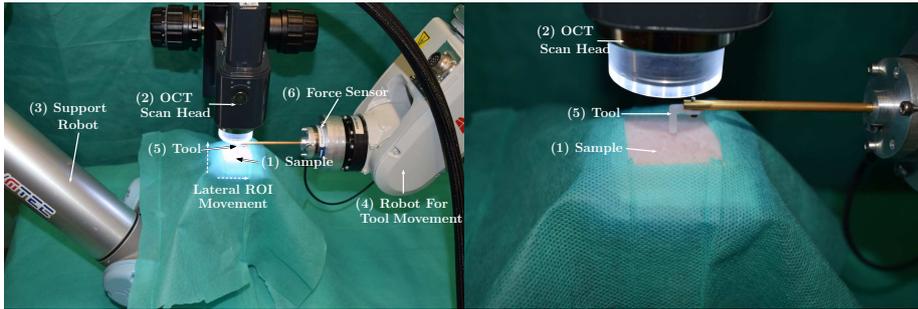
**Fig. 2:** The experimental setup we use for data acquisition. A tissue sample or phantom (1) is placed below the OCT device (2). A supporting robot (3) can move the sample in order to capture different ROIs during acquisition. The instrument robot (4) is equipped with a tool (5) and a force sensor (6) for ground-truth annotation of the OCT volumes. The main robot deforms the tissue with varying impact orientations.

Force predictions are entirely image-based and independent of the robotic system performing the motion. Therefore, only the tool and the elastic tissue parameters are relevant for the predictions made by the CNN. As a result, models can be pre-trained for specific tools and tissue types.

2.2 Experimental Setup

CNN training requires sufficiently large data sets. For automated data generation and systematic evaluation we use the setup shown in Figure 2. We place a phantom or tissue sample below the OCT scan head on a *support* robot. The robot occasionally moves the sample in order to capture different ROIs. An *instrument* robot is equipped with a tool and a force sensor at the tool base. Note, that the tool head is 3D printed and can be replaced. For each tool pose realized by the instrument robot we acquire an OCT image volume and the respective force vector. The data acquisition is performed as follows:

1. Without deformation, the OCT device acquires a reference volume and the force sensor performs a reference measurement
2. The instrument robot moves to a random orientation $\theta_x$, $\theta_y$, $\theta_z$ and deforms the sample with a random depth $d$
3. With deformation, the OCT device acquires a sample volume and the force sensor performs a measurement
4. After step (2) and (3) have been performed $L$ times, the support robot moves the sample in lateral directions by random values
5. Steps (1) to (4) are repeated $M$ times

As a result, for one iteration, we acquire $N = LM$ examples of reference and sample pairs with a force vector as a label for each pair.
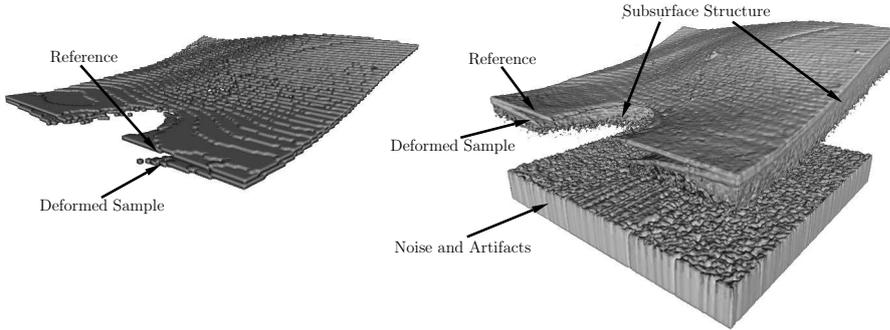
**Fig. 3:** Example data samples for training, shown as rendered volumes. On the right, full volumes are shown. On the left, the extracted surfaces are shown. In both cases, a deformed sample is overlayed with its reference measurement. A silicon phantom was used for these examples.

## 2.3 OCT Imaging and Image Processing

The imaging device is a spectral domain OCT system which is based on interferometry. The method captures 1D depth profiles (A-Scans) using infrared light. Repeated scanning at neighboring lateral points results in a volume scan of the ROI. We use an OCT working at 1300 nm wavelength and therefore we can capture the inner structure of an ROI in up to 1 mm depth in scattering tissue. A single raw OCT volumes has a size of $128 \times 128 \times 512$ voxels. For our 3D CNN approach, we consider downsampled versions of size $64 \times 64 \times 64$ due to time, computational and memory limitations.

We compare the volume-based method to approaches using the surface deformation only, extracted from the volumes. Maximum intensity projection (MIP) of the OCT volumes along the axial beam direction can be used for extraction [21]. The tissue phantom and tissue samples we use reflect the largest proportion of light at the surface. Therefore, the index at which the maximum intensity was observed represents a depth map of the tissue surface. Moreover, the intensity itself provides information of the surface characteristics as the intensity of the reflected light depends on the surface normal. We consider both as 2D surface representations for comparison to our volumetric approach.

Both the volumetric data and the depth images are shown in Figure 3. A deformed volume is overlayed with its corresponding reference volume. Note, that a shading of the tool is visible in the data which requires our model to be robust towards occlusion. As the occlusion is not related to the applied force, the 3D CNN can be expected to learn invariance towards the occlusion as CNNs have been shown to perform well at these kind of tasks [23].

## 2.4 Model Definition and Training
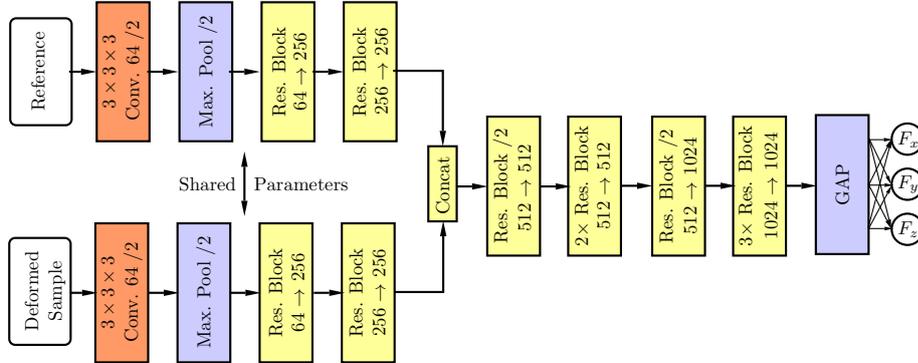
### 2.4.1 Model Architecture



**Fig. 4:** The siamese 3D CNN architecture. The model takes a deformed sample and its corresponding reference volume as its input. In the initial part, the two volumes are processed independently up to a concatenation point. At this point, the feature maps are aggregated and processed jointly. At the output, global average pooling (GAP) is applied to the remaining feature maps and a fully-connected layer leads to the force vector output. *Res. Block* refers to the residual blocks shown in Figure 5. /2 denotes a stride of two. Below each residual block, the change in the number of feature maps is denoted.

Our Siamese 3D CNN architecture is shown in Figure 4. Siamese CNN architectures take two images to be compared as their input [24]. Then, the images are initially processed independently by the same set of learnable filters. At a concatenation point, the feature maps of both images are aggregated and processed jointly by the remaining network layers [25, 26].

After the initial convolution in the network we employ residual blocks for an improved learning process [27]. The blocks are shown in Figure 5. Furthermore, our residual blocks employ the *bottleneck* concept [28]. Instead of directly applying convolutions, the input $x$ is downsampled first along its feature map dimension using a learnable $1 \times 1 \times 1$ filter. Then, the actual convolution with a larger $3 \times 3 \times 3$ filter is applied. Afterwards, the feature tensor is upsampled to its original feature map size. This method significantly reduces the number of learnable parameters and saves computational time.

The concatenation of the two paths is defined as follows. For path 1, consider a tensor $t_1$ of shape $[N_B, W, H, D, F_1]$ where $N_B$ is the batch size, $W$, $H$, $D$, are the feature maps' width, height and depth and $F_1$ is the number of feature maps. Thus, the concatenated tensor $t_c = t_1 \| t_2$ has a shape of $[N_B, W, H, D, F_1 + F_2]$. In our case, this doubles the number of feature maps which is why we keep the number of feature maps constant in the following spatial reduction ResNet block. This keeps the overall feature map sizes within the network at a reasonable level despite the concatenation.

Our general architecture choices are as follows. We use ReLu activation functions [29]. Before every activation we use batch normalization in order to reduce
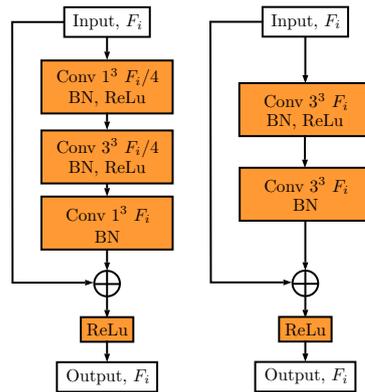
**Fig. 5:** The resdiual blocks for our architecture. $X^3$ denotes a $X \times X \times X$ filter. *BN* denotes a batch normalization layer. $F_i$ denotes the number feature maps that this layer produces. Left, the residual block we employ with a bottleneck is shown. Right, a residual block without bottleneck is shown for comparison.

internal covariate shift [30]. When we halve the spatial dimensions in our residual blocks, the $3 \times 3 \times 3$ filter uses a stride of two. We use nine residual blocks. From now on, we refer to our Siamese 3D CNN as SIAMCNN.

An alternative to a Siamese architecture is to use a single-path architecture that takes a difference or addition of volumes as its input. In this way, the deformation is captured in a single volume. We investigate performance when a single volume is passed to a 3D CNN that results from a subtraction or addition of the reference and deformed volume. The architecture is the same as the one shown in Figure 4, except that one path is removed. We refer to this architecture as DIFFCNN$_-$ for subtraction and DIFFCNN$_+$ for addition.

Lastly, we introduced surface extraction with MIPs in Section 2.3. We process the 2D depth representations with Siamese 2D CNN variants of our original architecture. The only difference to our model shown in Figure 4 is the use of 2D convolutions and 2D kernels instead of 3D convolutions and 3D kernels. We refer to this model as SURFCNN$_\text{MIP}$ for the 2D maximum intensity map as the model input and SURFCNN$_\text{DEPTH}$ for the 2D depth map as the model input.

*2.4.2 Training*

We train our models by minimizing the mean squared error (MSE) between ground-truth force labels and network predictions. We define the MSE as

$$MSE = \frac{1}{d} \sum_{i=1}^{d} \frac{1}{N_B} \sum_{j=1}^{N_B} (y_i^j - \hat{y}_i^j)^2 \qquad (1)$$

where $d$ is the number of outputs, $N_B$ the batch size, $y$ the ground-truth label and $\hat{y}$ the network's predictions. We use the Adam algorithm [31] for mini-batch gradient descent training. The initial learning rate is $l_r = 10^{-4}$. Every time the validation error plateaus, we divide the learning rate by a factor of 2 until no
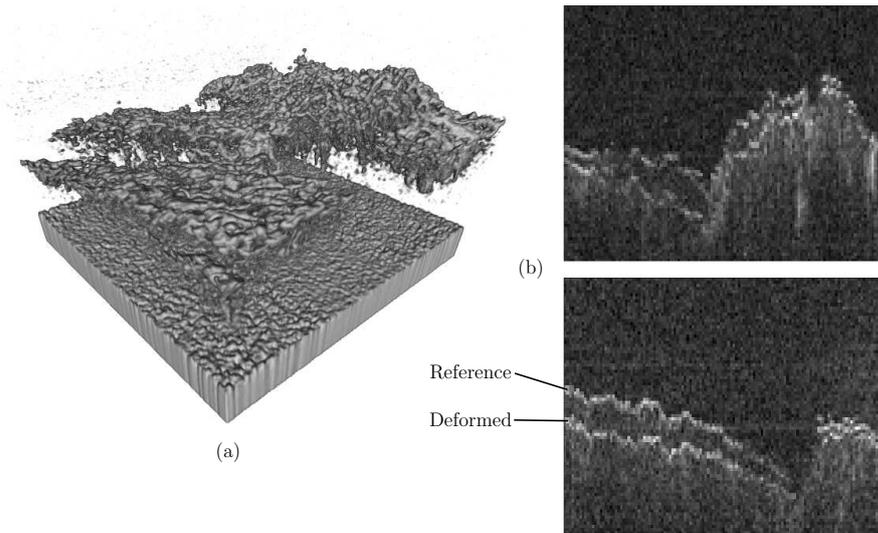
**Fig. 6:** Visualization of the tissue data. Left (a), a rendered volume of an OCT image that contains tissue is shown. Right (b), cropped, lateral slices through a volume are shown. The volume was created by overlaying a reference volume with a volume that contains deformed tissue.

further improvement can be observed. As typically done for regression, we rescale the labels to a range of $[0, 1]$ for training.

We perform hyperparameter selection on the validation set with a grid search with limited bounds for relevant hyperparameters which include the number of residual blocks, the position of the concatenation point, the total number of feature maps and the learning rate schedule. Besides, we follow standard architecture design principles for filter and feature map size per layer [32], batch normalization parameters [30] and Adam parameters [31].

## 2.5 Materials and Datasets

In our experimental setup we use a Thorlabs Telesto I SD-OCT device. Its lateral resolution is $15\,\mu\text{m}$ and its depth resolution is $7.5\,\mu\text{m}$. Its FOV covers a volume of $10\,\text{mm} \times 10\,\text{mm} \times 2.66\,\text{mm}$ resulting in image volumes with a size of $128 \times 128 \times 512$ voxels. The instrument robot performing the deformations is an ABB IRB120 6-axis manipulator. The robot performs rotations of the tool tip with ranges of $[-30°, 30°]$, $[0°, 10°]$ and $[-10°, 10°]$ for $\theta_x$, $\theta_y$ and $\theta_z$, respectively. The deformation depth $d$ is in the interval $[0.5\,\text{mm}, 1.5\,\text{mm}]$. The support robot is a UR-5 6-axis manipulator. The force sensor for ground-truth annotation is an ATI Nano43.

We acquired two datasets for the evaluation of our method. For the first dataset a silicon tissue phantom was used. In total, approximately 6600 pairs of image volumes were acquired. We divide the set into training, validation and test set with ratios of $80\,\%$, $10\,\%$ and $10\,\%$. We fine-tune our models on the validation set
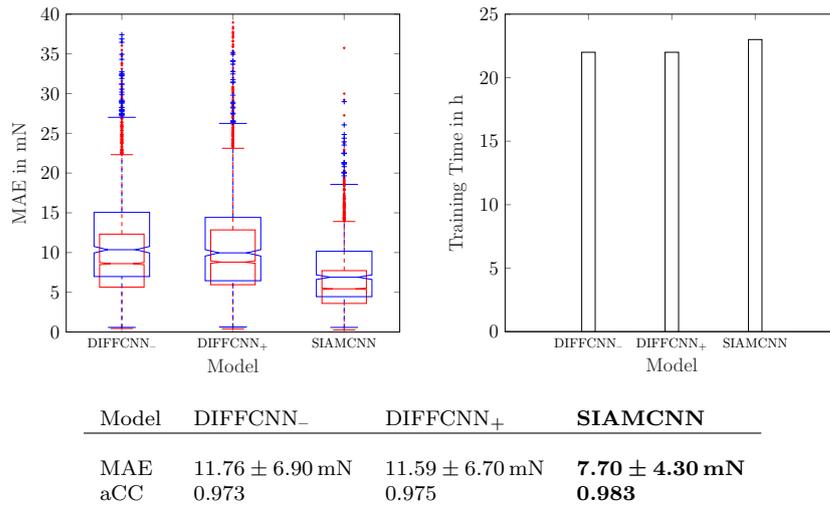
**Fig. 7 & Table 1:** Comparison for using the difference and addition of volumes compared to our Siamese approach. Top left, boxplots of the training MAE (red) and test MAE (blue) are shown. Top right, training durations until convergence are shown. Bottom, the MAE (with standard deviation) and the aCC are shown.

| Model | DIFFCNN_ | DIFFCNN_+ | **SIAMCNN** |
|---|---|---|---|
| MAE | $11.76 \pm 6.90$ mN | $11.59 \pm 6.70$ mN | **$7.70 \pm 4.30$ mN** |
| aCC | 0.973 | 0.975 | **0.983** |

and provide results for the test set. For the second dataset, porcine tissue was used. We acquired data with 17 different tissue samples from varying subjects. In total, approximately 8500 pairs of samples were acquired. We divide the set into training and test set with ratios of 80 % and 20 %. For this dataset we do not use a validation set since we use the tuned models derived from the phantom dataset. An example volume is shown in Figure 6.

We use the TensorFlow Environment [33] for implementation and train our models with an nVidia GTX 1080 Ti graphics card.

## 2.6 Evaluation Strategy

We use the mean average error (MAE) and average correlation coefficient (aCC) between network predictions and ground-truth labels for evaluation which are typical error metrics for regression [34]. Furthermore, we show the per sample MAE error distribution with boxplots, both for the training and the test set. For relevant model variations we show the training times until convergence.

## 3 Results

First, we compare SIAMCNN and DIFFCNN. The results are shown in Figure 7. Generally, the aCC shows that the SIAMCNN model accurately learned force estimation. Joining the reference and sample volume with a subtraction or addition does not make a difference in terms of performance. SIAMCNN outperforms both single-path approaches while having a similar training time. Note, that real-time
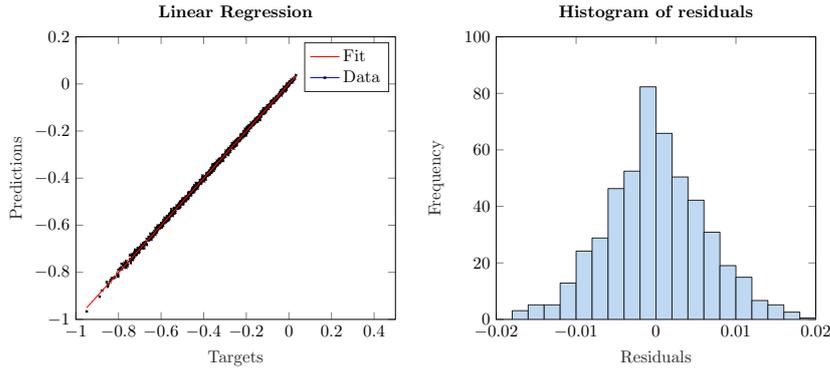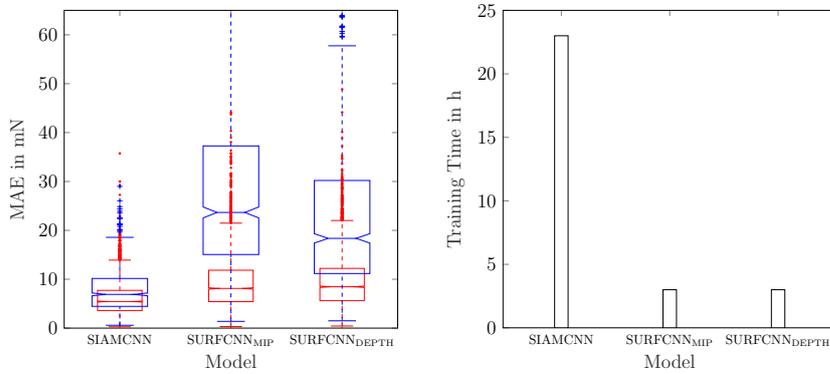
**Fig. 8:** Left, the linear regression plot between predictions and targets is shown. Right, the corresponding histogram of residuals is shown. Values are given in Newton. With an $R^2 = 0.99$ there is a strong relationship.



| Model | **SIAMCNN** | SURFCNN$_{\text{MIP}}$ | SURFCNN$_{\text{DEPTH}}$ |
|-------|-------------|------------------------|---------------------------|
| MAE | **$7.70 \pm 4.30\,\text{mN}$** | $29.42 \pm 22.90\,\text{mN}$ | $24.38 \pm 22.00\,\text{mN}$ |
| aCC | **0.983** | 0.870 | 0.877 |

**Fig. 9 & Table 2:** Comparison of different 2D surface representations and volumetric inputs. Top left, boxplots of the training MAE (red) and test MAE (blue) are shown. Top right, training durations until convergence are shown. Bottom, the MAE (with standard deviation) and the aCC are shown.

force estimation is feasible with an average processing time of $16.9 \pm 1.3\,\text{ms}$ for one instance of force estimation.

Furthermore, we study the general properties of SIAMCNN. The regression plot in Figure 8 shows the linear relationship between model predictions and targets. There is a tight relationship with a high $R^2$ value of 0.99.

Furthermore, we compare our baseline model SIAMCNN to the surface-based 2D approaches SURFCNN$_{\text{MIP}}$ and SURFCNN$_{\text{DEPTH}}$. The results are shown in Figure 9. The volume-based model SIAMCNN significantly outperforms the two
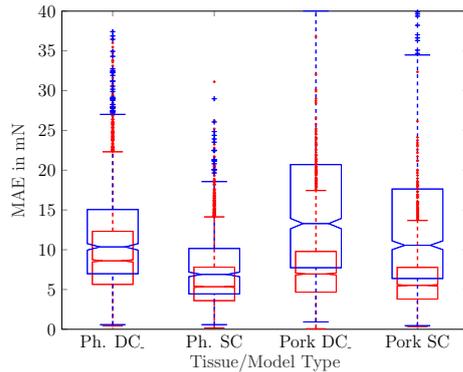
**Fig. 10:** Comparison of phantom data to tissue. Boxplots of the training MAE (red) and test MAE (blue) are shown in mN. *Ph.* refers to the phantom data. *DC* refers to DIFFCNN and *SC* refers to SIAMCNN.

|  | MAE | aCC |
|---|---|---|
| Ph. DIFFCNN_ | $11.76 \pm 6.9$ | 0.973 |
| Ph. SIAMCNN | $7.70 \pm 4.3$ | 0.983 |
| Pork DIFFCNN_ | $14.21 \pm 8.2$ | 0.969 |
| Pork SIAMCNN | $11.42 \pm 5.9$ | 0.977 |

**Table 3:** Comparison of phantom data to tissue. The MAE (with standard deviation) in mN and the aCC are shown.

surface-based approaches. In terms of training time, the 2D models require substantially less training time.

Lastly, we evaluate our method on a dataset with porcine tissue. The results are shown in Figure 10 and Table 3. The errors are slightly larger than for the phantom data. All in all, the Siamese 3D CNN is able to generalize well to a new subject that was not present during model training.

## 4 Discussion

We propose a novel method for image-based force estimation using OCT as an imaging modality. Image volumes are directly processed by a 3D CNN in order to predict a force acting between a tool and tissue. For this purpose, we introduce a novel Siamese 3D CNN architecture that processes a reference and a deformed sample simultaneously.

The results in Figure 7 and Figure 8 show that our method accurately learned force estimation with an MAE of $7.70 \pm 4.30\,$mN and an aCC of 0.983. This is achieved by only using two image volumes for one instance of force estimation. Many prior approaches depend on time series of deformations [22,35,36] which is intractable for entire volumes to be processed by CNN models. Our model can be trained within one day on a standard consumer graphics card and allows for real-time force estimation.

Still, in some application scenarios, it might be difficult to obtain a reference volume, e.g., due to rapid tissue motion. Generally, the acquisition of reference volumes can be sped up by using a faster OCT system. Recently, swept source OCT systems with A-scan rates of multiple MHz have been demonstrated and commercial systems with 1.5 MHz A-scan rate are available [37]. Moreover, full-filed OCT systems for fast parallel volume acquisition have been proposed [38]. Using such systems, small tissue patches at an instrument tip can be imaged with several hundred volumes per second, i.e., motion artifacts and latency would be minimal. Clearly, our method could be readily applied to OCT image volumes from such OCT devices.

Besides SIAMCNN, we consider a subtraction or addition of the reference and deformed volumes allowing for a single-path 3D CNN, DIFFCNN. This relates to previous approaches where differences of surfaces were considered for force estimation [21]. SIAMCNN is more accurate with an error of $7.70 \pm 4.30 \, \text{mN}$ compared to $11.76 \pm 6.90 \, \text{mN}$ for DIFFCNN$_-$. The improved performance implies that learning distinct preprocessing for both volumes within the 3D CNN is beneficial for force estimation. At the same time, there is hardly any difference in terms of training duration as the SIAMCNN and DIFFCNN models contain almost the same number of parameters. Therefore, our SIAMCNN model improves performance without demanding more resources.

Prior approaches relied on deformations that were obtained from surface reconstructions [19] or surface extraction [21]. For SIAMCNN, volume processing is superior with an error of $7.70 \pm 4.30 \, \text{mN}$ compared to $24.38 \pm 22.00 \, \text{mN}$ for the best performing 2D model SURFCNN$_{\text{DEPTH}}$. This suggests that capturing subsurface tissue compression with OCT image volumes allows for learning richer feature representations. Furthermore, the training and test error distribution depicted in Figure 9 show that overfitting occurs for the 2D case, despite our early-stopping training strategy. This also indicates that the surfaces extracted from OCT volumes carry fewer generalizable features than volumes.

In a last step, we validate our results on an animal tissue dataset that contains samples from different ROIs and varying subjects. For deep learning methods, overfitting is often an issue. In our case, the 3D CNN might overfit to subject-specific features which would hinder application in practice where the model is always applied to a new subject. Therefore, it is important to ensure that the model actually learned tissue- and not subject-specific deformation features. Our test set for the tissue datasets contains samples from a subject that was not present during training. The test set results in Figure 10 show that our model was able to produce accurate results despite subject variations. Therefore, similar to previous force estimation approaches [20], our method is able to achieve invariance towards small subject variations.

As our method directly learns relevant features from the volume data, the models do not require any adjustments for other tissue types and can be directly trained on new datasets. As a drawback, this requires acquisition of new datasets for new tissue types. With intraoperative OCT systems spreading in availability [39], our method could see application in practice, given extended validation with human operators who receive the estimated forces as feedback.

## 5 Conclusions

We address force estimation for tool-tissue interaction in surgery. For sensorless measurement, we propose a novel image-based force estimation method using OCT volume data. We associate a reference volume measurement with a volume of deformed tissue for force prediction with a 3D CNN. In order to process both volumes we introduce a novel Siamese 3D CNN architecture. OCT allows to capture inner structure of tissue. We demonstrate that exploiting deep tissue structures with volumes performs significantly better than deriving deformation from the surface in OCT volumes. Lastly, we show the applicability of our method on a tissue dataset where generalization to new subjects that were not present during training is feasible.

## Conflict of Interest

The authors declare that they have no conflict of interest.

## Ethical Approval

This article does not contain any studies with human participants or animals performed by any of the authors.

## Informed Consent

Informed consent was obtained from all individual participants included in the study.

## References

1. Wilson, E.B., Bagshahi, H., Woodruff, V.D. (2014) Overview of general advantages, limitations, and strategies. In: Robotics in General Surgery, pp. 17–22. Springer
2. Kroh, M., Chalikonda, S.: (2015) Essentials of robotic surgery. Springer
3. Diana, M., Marescaux, J. (2015) Robotic surgery. British Journal of Surgery **102**(2)
4. De Lorenzo, D., De Momi, E., Dyagilev, I., Manganelli, R., Formaglio, A., Prattichizzo, D., Shoham, M., Ferrigno, G. (2011) Force feedback in a piezoelectric linear actuator for neurosurgery. The International Journal of Medical Robotics and Computer Assisted Surgery **7**(3), 268–275
5. Meli, L., Pacchierotti, C., Prattichizzo, D. (2017) Experimental evaluation of magnified haptic feedback for robot-assisted needle insertion and palpation. The International Journal of Medical Robotics and Computer Assisted Surgery **13**(4)
6. Okamura, A.M. (2009) Haptic feedback in robot-assisted minimally invasive surgery. Current opinion in urology **19**(1), 102
7. Pacchierotti, C., Meli, L., Chinello, F., Malvezzi, M., Prattichizzo, D. (2015) Cutaneous haptic feedback to ensure the stability of robotic teleoperation systems. The International Journal of Robotics Research **34**(14), 1773–1787
8. Bayle, B., Joinie-Maurin, M., Barbe, L., Gangloff, J., De Mathelin, M. (2014) Robot interaction control in medicine and surgery: Original results and open problems. In: Computational Surgery and Dual Training, pp. 169–191. Springer

9. Puangmali, P., Liu, H., Seneviratne, L.D., Dasgupta, P., Althoefer, K. (2012) Miniature 3-axis distal force sensor for minimally invasive surgical palpation. IEEE/ASME Transactions on Mechatronics **17**(4), 646–656

10. Faragasso, A., Bimbo, J., Noh, Y., Jiang, A., Sareh, S., Liu, H., Nanayakkara, T., Wurdemann, H.A., Althoefer, K. (2014) Novel uniaxial force sensor based on visual information for minimally invasive surgery. In: ICRA, 2014 IEEE International Conference on, pp. 1405–1410

11. Sokhanvar, S., Dargahi, J., Najarian, S., Arbatani, S. (2012) Clinical and regulatory challenges for medical devices tactile sensing and displays. Haptic Feedback for Minimally Invasive Surgery and Robotics

12. Greminger, M.A., Nelson, B.J. (2004) Vision-based force measurement. IEEE Transactions on Pattern Analysis and Machine Intelligence **26**(3), 290–298

13. Kim, J., Janabi-Sharifi, F., Kim, J. (2010) A haptic interaction method using visual information and physically based modeling. IEEE/ASME Transactions on Mechatronics **15**(4), 636–645

14. Noohi, E., Parastegari, S., Žefran, M. (2014) Using monocular images to estimate interaction forces during minimally invasive surgery. In: IROS, 2014 IEEE/RSJ International Conference on, pp. 4297–4302

15. Kim, W., Seung, S., Choi, H., Park, S., Ko, S.Y., Park, J.O. (2012) Image-based force estimation of deformable tissue using depth map for single-port surgical robot. In: ICCAS, 2012 12th International Conference on, pp. 1716–1719. IEEE

16. Greminger, M.A., Nelson, B.J. (2003) Modeling elastic objects with neural networks for vision-based force measurement. In: Intelligent Robots and Systems, 2003.(IROS 2003). Proceedings. 2003 IEEE/RSJ International Conference on, vol. 2, pp. 1278–1283. IEEE

17. Karimirad, F., Chauhan, S., Shirinzadeh, B. (2014) Vision-based force measurement using neural networks for biological cell microinjection. Journal of Biomechanics **47**(5), 1157–1163

18. Mozaffari, A., Behzadipour, S., Kohani, M. (2014) Identifying the tool-tissue force in robotic laparoscopic surgery using neuro-evolutionary fuzzy systems and a synchronous self-learning hyper level supervisor. Applied Soft Computing **14**, 12–30

19. Aviles, A.I., Alsaleh, S.M., Sobrevilla, P., Casals, A. (2015) Force-feedback sensory substitution using supervised recurrent learning for robotic-assisted surgery. In: EMBC, 2015 37th Annual International Conference of the IEEE, pp. 1–4

20. Aviles, A.I., Alsaleh, S.M., Hahn, J.K., Casals, A. (2017) Towards retrieving force feedback in robotic-assisted surgery: A supervised neuro-recurrent-vision approach. IEEE Transactions on Haptics **10**(3), 431–443

21. Otte, C., Beringhoff, J., Latus, S., Antoni, S.T., Rajput, O., Schlaefer, A. (2016) Towards force sensing based on instrument-tissue interaction. In: MFI, 2016 IEEE International Conference on, pp. 180–185. IEEE

22. Rivero, A.I.A., Alsaleh, S.M., Hahn, J.K., Casals, A. (2016) Towards retrieving force feedback in robotic-assisted surgery: A supervised neuro-recurrent-vision approach. IEEE Transactions on Haptics

23. Goodfellow, I., Bengio, Y., Courville, A.: (2016) Deep Learning. MIT Press

24. Leal-Taixé, L., Canton-Ferrer, C., Schindler, K. (2016) Learning by tracking: Siamese cnn for robust target association. In: Proceedings of the IEEE CVPR Workshops, pp. 33–40

25. Zbontar, J., LeCun, Y. (2015) Computing the stereo matching cost with a convolutional neural network. In: Proceedings of the IEEE CVPR, pp. 1592–1599

26. 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 ICCV, pp. 2758–2766

27. He, K., Zhang, X., Ren, S., Sun, J. (2016) Identity mappings in deep residual networks. In: ECCV, pp. 630–645

28. He, K., Zhang, X., Ren, S., Sun, J. (2016) Deep residual learning for image recognition. In: Proceedings of the IEEE CVPR, pp. 770–778

29. Glorot, X., Bordes, A., Bengio, Y. (2011) Deep Sparse Rectifier Neural Networks. In: Aistats, vol. 15, p. 275

30. Ioffe, S., Szegedy, C. (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML, pp. 448–456

31. Kingma, D., Ba, J. (2014) Adam: A method for stochastic optimization. In: ICLR

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

33. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M. (2016) Tensorflow: Large-scale machine learning on heterogeneous distributed systems. arXiv preprint arXiv:1603.04467

34. Borchani, H., Varando, G., Bielza, C., Larrañaga, P. (2015) A survey on multi-output regression. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery **5**(5), 216–233

35. Carrasco-Zevallos, O., Keller, B., Viehland, C., Shen, L., Waterman, G., Todorich, B., Shieh, C., Hahn, P., Farsiu, S., Kuo, A., Toth, C., Izatt, J. (2016) Live volumetric (4d) visualization and guidance of in vivo human ophthalmic surgery with intraoperative optical coherence tomography. Scientific Reports **6**, 31,689

36. Aviles, A.I., Alsaleh, S., Sobrevilla, P., Casals, A. (2016) Exploring the effects of dimensionality reduction in deep networks for force estimation in robotic-assisted surgery. In: SPIE Medical Imaging, pp. 97,861X–97,861X. International Society for Optics and Photonics

37. Potsaid, B., Baumann, B., Huang, D., Barry, S., Cable, A.E., Schuman, J.S., Duker, J.S., Fujimoto, J.G. (2010) Ultrahigh speed 1050nm swept source/fourier domain oct retinal and anterior segment imaging at 100,000 to 400,000 axial scans per second. Optics express **18**(19), 20,029–20,048

38. Hillmann, D., Franke, G., Hinkel, L., Bonin, T., Koch, P., Hüttmann, G. (2013) Off-axis full-field swept-source optical coherence tomography using holographic refocusing. In: Optical Coherence Tomography and Coherence Domain Optical Methods in Biomedicine XVII, vol. 8571, p. 857104. International Society for Optics and Photonics

39. Ehlers, J.P., Srivastava, S.K., Feiler, D., Noonan, A.I., Rollins, A.M., Tao, Y.K. (2014) Integrative advances for oct-guided ophthalmic surgery and intraoperative oct: microscope integration, surgical instrumentation, and heads-up display surgeon feedback. PloS one **9**(8), e105,224