

Two-path 3D CNNs for calibration of system parameters for OCT-based motion compensation

Nils Gessert^{*a}, Martin Gromniak^{*a}, Matthias Schlüter^a, and Alexander Schlaefer^a

^aInstitute of Medical Technology, Hamburg University of Technology, Am Schwarzenberg-Campus 3, 21073 Hamburg, Germany

ABSTRACT

Automatic motion compensation and adjustment of an intraoperative imaging modality's field of view is a common problem during interventions. Optical coherence tomography (OCT) is an imaging modality which is used in interventions due to its high spatial resolution of few micrometers and its temporal resolution of potentially several hundred volumes per second. However, performing motion compensation with OCT is problematic due to its small field of view which might lead to tracked objects being lost quickly. We propose a novel deep learning-based approach that directly learns input parameters of motors that move the scan area for motion compensation from optical coherence tomography volumes. We design a two-path 3D convolutional neural network (CNN) architecture that takes two volumes with an object to be tracked as its input and predicts the necessary motor input parameters to compensate the object's movement. In this way, we learn the calibration between object movement and system parameters for motion compensation with arbitrary objects. Thus, we avoid error-prone hand-eye calibration and handcrafted feature tracking from classical approaches. We achieve an average correlation coefficient of 0.998 between predicted and ground-truth motor parameters which leads to sub-voxel accuracy. Furthermore, we show that our deep learning model is real-time capable for use with the system's high volume acquisition frequency.

Keywords: Deep Learning, 3D CNN, Motion Compensation, OCT

1. INTRODUCTION

Optical coherence tomography (OCT) is an interferometric imaging modality that allows for volumetric imaging with micrometer-level resolution. OCT has been used in intraoperative scenarios^{1,2} such as neurosurgery³ and ophthalmic surgery.⁴ Recently, systems with high-frequency acquisition have been proposed^{5,6} which allows for fast imaging and object tracking during interventions. As the field of view (FOV) of high-resolution imaging modalities is often limited, the current region-of-interest (ROI) might be lost quickly due to patient and surgical tool movement. As manual adjustment of the imaging system's FOV disrupts the surgical workflow, automatic adjustment is desirable for keeping track of the current ROI. So far, OCT-based tracking and compensation can be performed with markerless approaches using cumbersome and potentially error-prone image-based registration.⁷ Alternatively, markers can be introduced to the setup which can be invasive but promises higher accuracy. For example, detection of artificial landmarks carved into bone structures has been shown.⁸ More recently, a deep learning-based method has been proposed where a model learns to estimate the pose of a very small, arbitrary marker geometry directly from OCT volumes.⁹ For motion compensation, all these methods require a hand-eye calibration between imaging system and compensation device which is difficult for OCT.¹⁰

In this paper, we propose a calibration strategy between OCT volumes and a compensation system for marker-based tracking. We consider the setup shown in Figure 1 with an OCT system that has a mechanism for lateral and axial FOV adjustment. Two motors control mirrors for lateral beam redirection and one motor controls a reference arm for adjusting the scan distance. Thus, the motors can be used to compensate motion of the marker object and keep it within the FOV. In order to compensate motion of the object in this setup, a

Further author information: (Send correspondence to Nils Gessert)

Nils Gessert: nils.gessert@tuhh.de

Martin Gromniak: martin.gromniak@tuhh.de

* Both authors contributed equally.

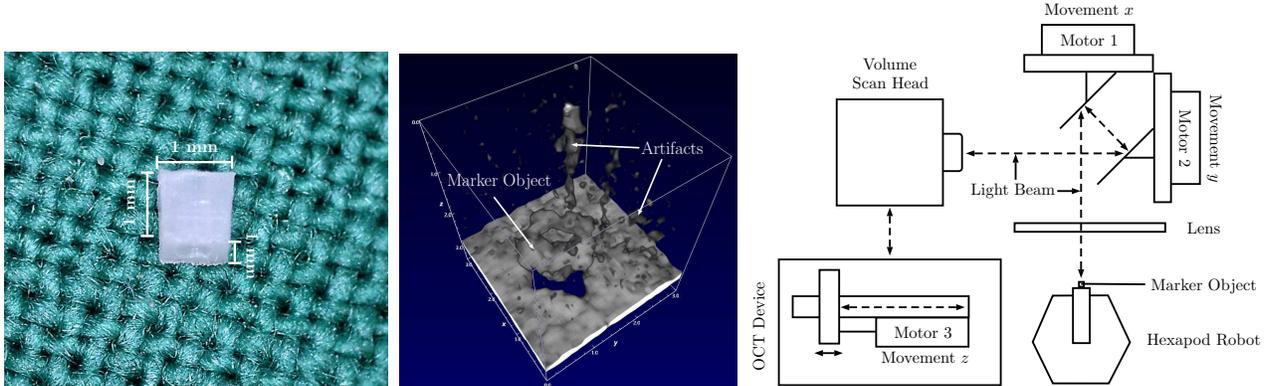


Figure 1: The object geometry to be tracked is shown under a digital microscope (left) and in a rendered OCT volume (center). The experimental setup for motion compensation and data acquisition is shown as a draft (right).

classic approach first requires an OCT-based hand-eye calibration. Second, either an image-based registration of OCT volumes¹¹ is required or a known marker geometry needs to be detected. Instead, we propose a new direct calibration approach between volumes and motors that combines both steps in a single deep learning model. For this purpose, we extend the recent idea of a 3D convolutional neural network (CNN) model for the estimation of an arbitrary marker’s pose⁹ to the calibration problem. Instead of a single volume, our model receives two volumes with an object in different areas of the FOV. The two volumes are processed with a two-path 3D CNN architecture.¹² At the output, the model predicts motor steps that need to be driven to compensate the motion between the two object states. In this way, we combine hand-eye calibration of OCT volumes with motors and marker detection in a single trainable model. For training of the two-path 3D CNN model we acquire 7 datasets of an object and we show with a separate dataset that the model learns to compensate the object’s motion. Thus, the object can be effectively used as marker to keep track of a desired region of interest. Last, we show that the model has low inference times which allows for real-time estimation, despite performing volumetric data processing.

2. METHODS AND MATERIALS

2.1 Experimental Setup

The setup for OCT-based motion compensation and the object to be tracked is shown in Figure 1. The marker object is milled from a polyoxymethylene block with a size of $1\text{ mm} \times 1\text{ mm} \times 1\text{ mm}$. We carved out an inner structure in order to have subsurface features that can be imaged by OCT and exploited by deep learning models, as recently suggested.⁹ The setup itself consists of a swept-source OCT device (OMES, Optores) with an A-Scan rate of 1.59 MHz. We use a scan head that provides volumes of size $32 \times 32 \times 460$ voxels. This leads to a potential acquisition speed of 833 volumes per second. For a uniform volume size and reduced processing time, we downsample the volumes to a size of $32 \times 32 \times 32$ voxels which covers a volume of approximately $3\text{ mm} \times 3\text{ mm} \times 3.5\text{ mm}$. The volume’s position in space can be adjusted by three stepper motors. Two motors control mirrors that can laterally move the FOV by $\approx 60\text{ mm}$. The third motor moves the mirror in the reference arm in a range of $\approx 160\text{ mm}$. For data acquisition, we also use a hexapod robot with the marker object attached to it. The robot’s purpose is to move the object into different orientations for a higher variability in object appearance.

2.2 Data Acquisition

For training the deep learning model, a large, labeled dataset is required which we acquire automatically with the setup. In each step, the hexapod moves the object into a random orientation. Then, we move the FOV to two randomly generated motor states s_1 and s_2 and acquire a volume in each state. Thus, a single labeled example consists of two volumes and the label $s_d = s_1 - s_2 = (\Delta x, \Delta y, \Delta z)$ which needs to be driven to overlay the

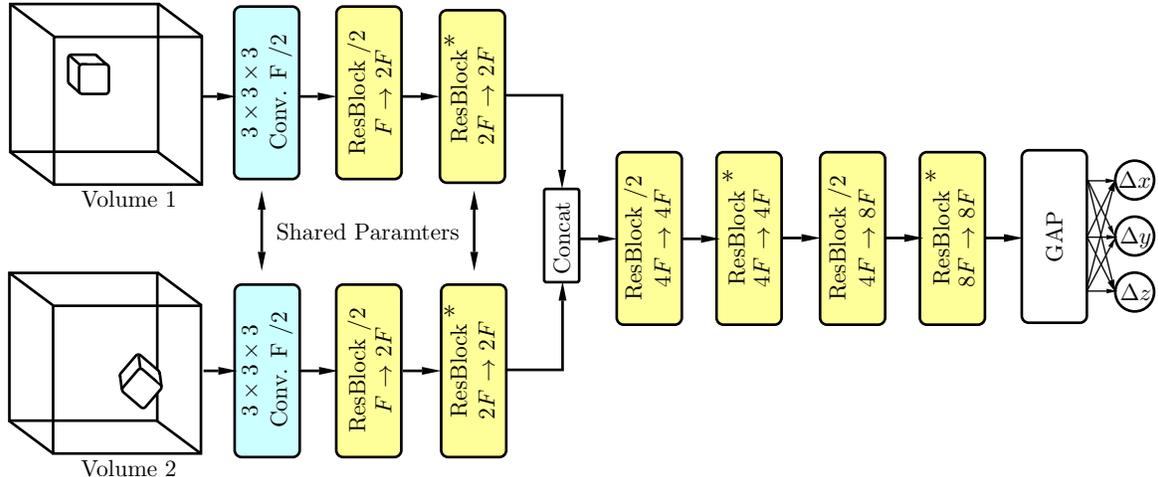


Figure 2: The proposed two-path 3D CNN architecture. In each block, the change in the number of feature maps is denoted. F refers to the base number of feature maps. ResBlock refers to residual blocks as introduced by He et al.¹³ /2 denotes spatial downsampling with a stride of 2. Note, that in the initial two path the model parameters are shared. Concat denotes tensor concatenation along the feature map dimension. GAP denotes global average pooling. GAP is followed by a single linear fully-connected layer. The ResBlocks marked with an asterisk are omitted in the reduced architecture.

Table 1: Performance results on the test set and inference times for the different models. MAE refers to the mean absolute error in motor steps. 2.5 motor steps in x and y direction roughly correspond to a shift of one voxel. For the z direction, roughly 190 motor steps correspond to a shift of one voxel. ACC denotes the average correlation coefficient between predictions and targets. F denotes the base number of feature maps, see Figure 2. Red. denotes a reduced architecture with less ResBlocks.

	MAE Δx	MAE Δy	MAE Δz	ACC	Inf. Time
Resnet $F = 60$	1.628 ± 1.326	1.426 ± 1.166	42.41 ± 34.34	0.9983	7.51 ± 0.12 ms
Resnet $F = 45$	1.990 ± 1.606	1.634 ± 1.290	48.06 ± 36.45	0.9984	5.40 ± 0.12 ms
Resnet $F = 30$	1.633 ± 1.270	1.285 ± 1.078	43.83 ± 34.84	0.9983	3.73 ± 0.16 ms
Resnet $F = 15$	1.792 ± 1.393	1.564 ± 1.243	53.12 ± 42.13	0.9978	2.51 ± 0.17 ms
Resnet $F = 15$ Red.	2.008 ± 1.619	1.728 ± 1.405	55.91 ± 45.84	0.9966	1.79 ± 0.16 ms

volumes on top of each other. In total, we acquire 7 datasets with approximately 5000 examples each. Between each dataset acquisition we rearrange the marker in order to avoid overfitting to a particular initial marker pose.

2.3 Model

The two-path 3D CNN architecture we employ is shown in Figure 2. Each path receives a volume which is processed independently up to a concatenation point. Afterwards, the features are processed jointly and finally, the state difference s_d that would be required for compensation is predicted at the output. We rely on the ResNet principle¹⁴ with identity connections for improved gradient propagation. We also share parameters between the two paths as they receive similar volumes and therefore are likely to require similar features. As the OCT system’s volume acquisition rate is very high, we investigate the performance-inference time trade-off by considering downscaled variants of the model shown in Figure 2. For this purpose we vary the base feature map size F which controls the overall capacity of the model. Also, we consider a reduced version of the architecture with less ResBlocks. We train the models with a mean squared error loss using stochastic gradient descent using the Adam¹⁵ optimizer, a constant learning rate of $5e-4$ and a batch size of 40. We split off an entire independent dataset for testing with 5000 examples. We implement our model using Tensorflow.¹⁶ Training and inference time tests are performed on an NVIDIA GeForce GTX 1080 Ti graphics card.

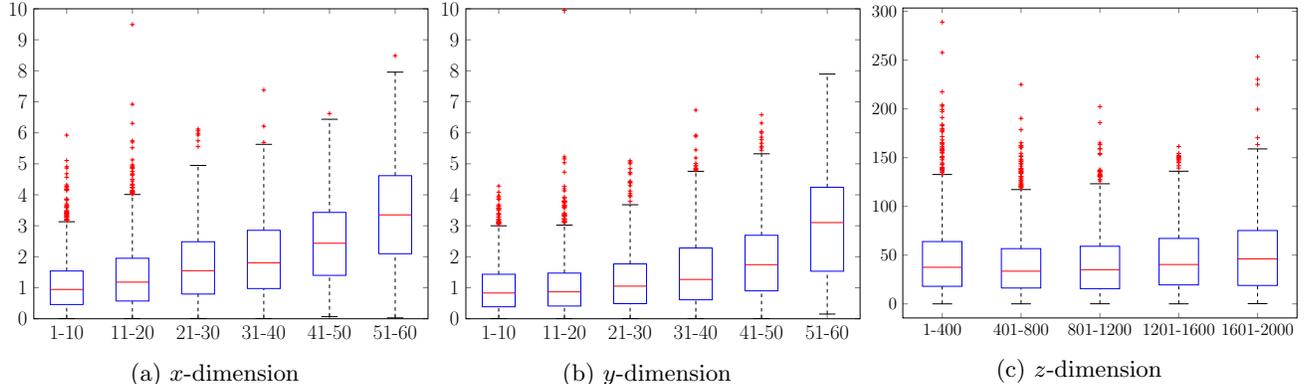


Figure 3: Absolute errors versus the magnitude of required motor steps along each dimension for $F = 30$. The horizontal axes show the absolute motor steps along each dimension while the vertical axes show the absolute errors of the predictions. The error increases only minor until ~ 40 steps in lateral x - and y -dimension and remains almost constant in the z -dimension.

3. RESULTS

The performance results on the test set and inference times for several model variants with differently sized architectures are shown in Table 1. Overall, the models’ performance is very high with an ACC larger than 0.996. When downscaling the model, performance slightly deteriorates with base feature maps size below $F = 30$. However, the inference time of the smallest model is substantially reduced to 23.8% of the largest model’s inference time while the ACC barely changes. It is notable that the inference time drops even more when removing ResBlocks for $F = 15$ on top of the feature map reduction.

Moreover, Figure 3 shows the absolute motor step errors depending on the absolute motor step distances of the labels along each dimension. With larger motor step distances that need to be compensated, the error would be expected to increase. The plot shows that this increase is only minor for a large portion of the motor step range. It can be used as an indication of the motion magnitude that can be expected to have good tracking performance.

Assuming a rough calibration factor of ~ 2.5 for lateral motor steps to voxels and factor of ~ 190 for the axial direction, our method qualitatively achieves sub-voxel accuracy. Considering the volume size of $3 \text{ mm} \times 3 \text{ mm} \times 3.5 \text{ mm}$ with a resolution of $32 \times 32 \times 32$ voxels, the absolute errors are well below $100 \mu\text{m}$.

4. DISCUSSION AND CONCLUSION

We propose a new deep learning-based method for OCT-based motion compensation. In particular, we avoid time-consuming and inaccurate volume-based registration with a subsequent hand-eye calibration by directly learning the calibration between marker object movement observed in 3D volumes and motors which move the scan area. For this purpose, we use a two-path 3D CNN architecture that predicts the required motor steps for motion compensation of an object’s movement based on two input volumes. Considering the results in Table 1, the very high average correlation coefficient shows that the learning problem is well solved. Also, the absolute errors show that we qualitatively achieve sub-voxel accuracy which translates to errors well below $100 \mu\text{m}$. With Figure 3 we can show that even large distances to be compensated only lead to a minor increase in error. Thus, our model should be robust even towards rapid and large motion. As we reattached the object several times in between dataset acquisition, the results indicate that the model is capable of learning to track the object. Thus, the object can be used as marker in a region-of-interest that should be tracked. When using a more efficient, downscaled architecture, performance is slightly reduced as a trade-off for faster inference. With 1.79 ms the model achieves a processing frequency of 559 volumes per second which is among the same magnitude as the OCT’s acquisition frequency. This shows that our model does not constitute a bottleneck in the entire compensation process despite having to perform volumetric data processing. For future work, our calibration

strategy could be extended by using different marker objects in order to achieve generalization to new marker types. Also, the more challenging motion compensation task of markerless tissue tracking could be addressed.

REFERENCES

- [1] Lankenau, E., Klinger, D., Winter, C., Malik, A., Müller, H. H., Oelckers, S., Pau, H.-W., Just, T., and Hüttmann, G., “Combining optical coherence tomography (OCT) with an operating microscope,” in [*Advances in Medical Engineering*], 343–348, Springer (2007).
- [2] Ehlers, J. P., Srivastava, S. K., Feiler, D., Noonan, A. I., Rollins, A. M., and Tao, Y. K., “Integrative advances for OCT-guided ophthalmic surgery and intraoperative OCT: microscope integration, surgical instrumentation, and heads-up display surgeon feedback,” *PloS One* **9**(8), e105224 (2014).
- [3] Finke, M., Kantelhardt, S., Schlaefer, A., Bruder, R., Lankenau, E., Giese, A., and Schweikard, A., “Automatic scanning of large tissue areas in neurosurgery using optical coherence tomography,” *The International Journal of Medical Robotics and Computer Assisted Surgery* **8**(3), 327–336 (2012).
- [4] Tao, Y. K., Srivastava, S. K., and Ehlers, J. P., “Microscope-integrated intraoperative OCT with electrically tunable focus and heads-up display for imaging of ophthalmic surgical maneuvers,” *Biomed Opt Express* **5**(6), 1877–1885 (2014).
- [5] Novais, E. A., Adhi, M., Moulton, E. M., Louzada, R. N., Cole, E. D., Husvagt, L., Lee, B., Dang, S., Regatieri, C. V., Witkin, A. J., Bauman, C. R., Hornegger, J., Jayaraman, V., Fujimoto, J. G., Duker, J. S., and Waheed, N. K., “Choroidal neovascularization analyzed on ultrahigh-speed swept-source optical coherence tomography angiography compared to spectral-domain optical coherence tomography angiography,” *American Journal of Ophthalmology* **164**, 80 – 88 (2016).
- [6] Siddiqui, M., Nam, A. S., Tozburun, S., Lippok, N., Blatter, C., and Vakoc, B. J., “High-speed optical coherence tomography by circular interferometric ranging,” *Nature Photonics* **12**(2), 111 (2018).
- [7] Laves, M.-H., Schoob, A., Kahrs, L. A., Pfeiffer, T., Huber, R., and Ortmaier, T., “Feature tracking for automated volume of interest stabilization on 4D-OCT images,” in [*SPIE Medical Imaging*], 101350W–101350W (2017).
- [8] Zhang, Y. and Wörn, H., “Optical coherence tomography as highly accurate optical tracking system,” in [*IEEE/ASME International Conference on Advanced Intelligent Mechatronics, 2014*], 1145–1150, IEEE, Piscataway, NJ (2014).
- [9] Gessert, N., Schlüter, M., and Schlaefer, A., “A deep learning approach for pose estimation from volumetric oct data,” *Medical image analysis* **46**, 162–179 (2018).
- [10] Rajput, O., Antoni, S.-T., Otte, C., Saathoff, T., Matthäus, L., and Schlaefer, A., “High accuracy 3D data acquisition using co-registered OCT and Kinect,” in [*IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems*], 32–37 (2016).
- [11] Niemeijer, M., Garvin, M. K., Lee, K., van Ginneken, B., Abramoff, M. D., and Sonka, M., “Registration of 3d spectral oct volumes using 3d sift feature point matching,” in [*Medical Imaging: Image Processing*], **7259**, 72591I (2009).
- [12] Gessert, N., Beringhoff, J., Otte, C., and Schlaefer, A., “Force estimation from OCT volumes using 3D CNNs,” *International journal of computer assisted radiology and surgery* **13**(7), 1073–1082 (2018).
- [13] He, K., Zhang, X., Ren, S., and Sun, J., “Identity mappings in deep residual networks,” in [*ECCV*], 630–645 (2016).
- [14] He, K., Zhang, X., Ren, S., and Sun, J., “Deep residual learning for image recognition,” in [*IEEE Conference on Computer Vision and Pattern Recognition*], 770–778 (2016).
- [15] Kingma, D. and Ba, J., “Adam: A method for stochastic optimization,” in [*ICLR*], (2014).
- [16] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., and Devin, M., “Tensorflow: Large-scale machine learning on heterogeneous distributed systems,” *arXiv preprint arXiv:1603.04467* (2016).