

LitCall: Learning Implicit Topology for CNN-based Aortic Landmark Localization

Zhangxing Bian^{a,b}, Jiayang Zhong^{a,b}, Yanglong Lu^a, Charles R. Hatt^c, and Nicholas S. Burris^a

^aDepartment of Radiology, University of Michigan, Ann Arbor, MI, USA;

^bDepartment of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI, USA;

^cImbio LLC, Minneapolis, MN, USA;

ABSTRACT

Landmark detection is a critical component of the image processing pipeline for automated aortic size measurements. Given that the thoracic aorta has a relatively conserved topology across the population and that a human annotator with minimal training can estimate the location of unseen landmarks from limited examples, we proposed an auxiliary learning task to learn the implicit topology of aortic landmarks through a CNN-based network. Specifically, we created a network to predict the location of missing landmarks from the visible ones by minimizing the *Implicit Topology* loss in an end-to-end manner. The proposed learning task can be easily adapted and combined with Unet-style backbones. To validate our method, we utilized a dataset consisting of 207 CTAs, labeling four landmarks on each aorta. Our method outperforms the state-of-the-art Unet-style architectures (ResUnet, UnetR) in terms of localization accuracy, with only a light (#params=0.4M) overhead. We also demonstrate our approach in two *clinically* meaningful applications: aortic sub-region division and automatic centerline generation.

Keywords: Landmark localization, deep learning, aorta, implicit topology

1. INTRODUCTION

Fast and accurate localization of aortic landmarks can facilitate image analysis tasks such as segmentation, registration and classification of anatomic sub-regions. In clinical practice, these tasks are performed manually, leading to significant inefficiency. Deep convolutional neural networks (CNNs) can be used to automate landmark annotation using a variety of approaches including: classifying image slices,¹ model ensembling,² landmark coordinate regression³ and heatmap regression.⁴ Classifying image slices suffers from severe class imbalance, and directly regressing coordinates requires a large number of parameters and highly non-linear mapping. Heatmap regression and ensemble learning models better handle overfitting and show robustness to image variability and artifacts. Our proposed method aligns most closely with heatmap regression which assumes the probability of landmark location is not uniformly distributed over the image.

Despite differences in size, tortuosity, and the location of adjacent structures between individuals, the thoracic aorta has a relatively regular shape. Thus, we propose that the topology formed by landmarks in the thoracic aorta can serve as a prior during learning. For example, a human rater can identify the approximate location of a missing landmark after observing only a few examples as shown in Fig.1. Inspired by this, we designed an auxiliary task for our model: learn to predict the locations of missing landmarks from a fraction of visible ones by optimizing the so-called *Implicit Topology* loss in an end-to-end manner. We show that by combining the proposed learning task with Unet-style backbone, the accuracy of localization is improved with a small parameter overhead (0.4M extra parameters).

Beyond technical implementation, we also demonstrate two clinically meaningful downstream image analysis applications: automatic centerline generation and classification of anatomic sub-regions. Generation of an aortic

Further author information: (Send correspondence to Nicholas S. Burris.)

Nicholas S. Burris.: nburris@med.umich.edu

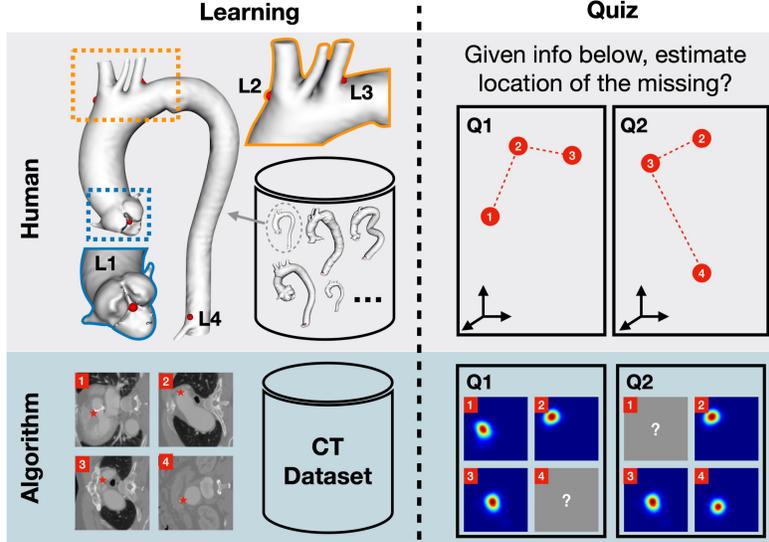


Figure 1: Humans *efficiently* learn the topology of spatial patterns. Even seeing a few examples of 4 aortic landmarks, one can give a decent answer to the quizzes above. We propose an auxiliary task for a network to learn landmark topology implicitly, which greatly improves landmark localization the accuracy prevents overfitting.

centerline is an important step for most two- and three-dimensional analyses of aortic disease including deformation analysis,⁵⁻⁷ however, manually labeling seeds point to generate centerline can be cumbersome. We utilize our model to predict accurate seed locations in the ascending and distal descending aorta to automate this process. Additionally, by using the predicted landmarks as the boundary to subdivide the thoracic aorta into ascending, arch, and descending segments, this work opens the possibility of performing segment-wise assessments of aortic size and growth.

2. METHOD

Given the input CT image I , we predict a set of heatmaps \mathcal{H} through a network parametrized by θ :

$$\mathcal{H} = f_h(I; \theta), \mathcal{H} = \{h_i, \dots, h_M\}, \quad (1)$$

where M is the number of landmarks. The final landmark coordinates are obtained as the peak responses of the predicted heatmaps.

To obtain the ground-truth heatmap, we apply a Gaussian function to the ground-truth landmarks:

$$g_i(\mathbf{x}; \sigma_i) = \frac{A}{(2\pi)^{d/2} \sigma_i^d} \exp\left(-\frac{\|\mathbf{x} - \mathbf{L}_i^*\|_2^2}{2\sigma_i^2}\right) \quad (2)$$

where A is the scaling factor, \mathbf{L}_i^* is the coordinate of the i -th landmark, and σ_i controls the width of the filter response. σ is a learnable parameter that is optimized during training. To avoid a trivial solution of σ going to infinity and network predicting a constant everywhere, we minimize $\|\sigma\|_2^2$ as a extra regularization term.

For convenience, we denote the collection of g_i as set \mathcal{G} and define MSE over two sets \mathcal{H}, \mathcal{G} as:

$$\text{MSE}(\mathcal{H}, \mathcal{G}, \sigma) = \frac{1}{MN} \sum_{\mathbf{x}} \|h_i - g_i(\mathbf{x}; \sigma_i)\|_2^2, \\ h_i \in \mathcal{G}, g_i \in \mathcal{H}, |\mathcal{G}| = |\mathcal{H}|$$

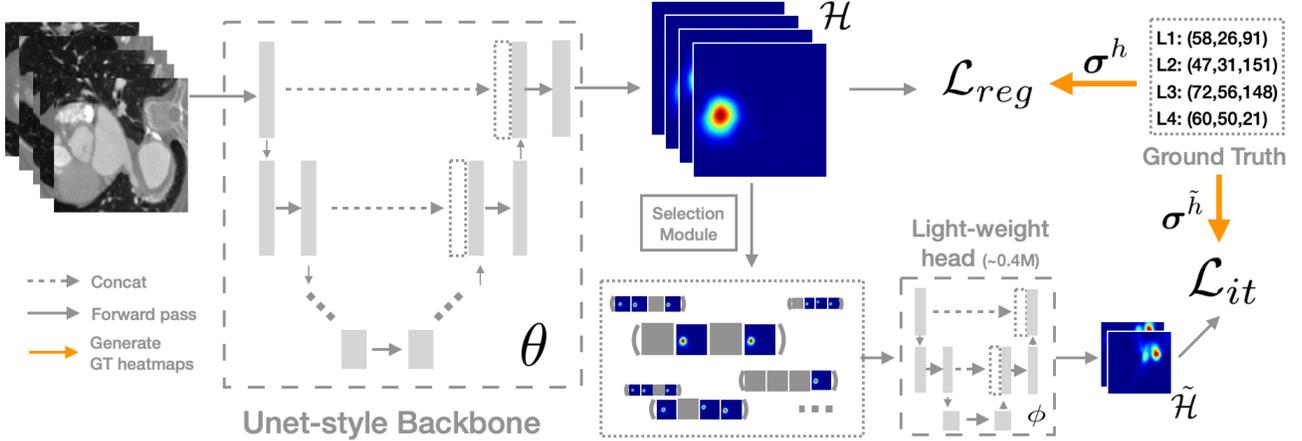


Figure 2: Illustration of our pipeline. The uncropped image is sent to a Unet-style backbone to regress heatmaps. The predicted heatmaps pass through the selection module to a light weight head, which is a three-scale resUNet, to predict the heatmaps of missing landmarks. \mathcal{L}_{reg} and \mathcal{L}_{it} indicate losses for regression and learning implicit topology.

where N is the number of voxels, thus, the heatmap regression loss can be simply defined as:

$$\mathcal{L}_{reg} = \text{MSE}(\mathcal{H}, \mathcal{G}; \sigma^h) \quad (3)$$

Then we introduce a selection module s , which randomly select k maps from set \mathcal{H} as visible heatmaps. k will gradually increase during the training. The heatmaps of a missing landmark $\tilde{\mathcal{H}}$ can be obtained by:

$$\tilde{\mathcal{H}} = f_{\tilde{h}}(s(\mathcal{H}; k); \theta, \phi), \tilde{\mathcal{H}} = \mathcal{H} - s(\mathcal{H}; k), \quad (4)$$

where $f_{\tilde{h}}$ is the function parametrized by both θ and ϕ . Note that the input for function $f_{\tilde{h}}$ only contains the predicted visible heatmap but without image features, which eliminates the trivial solution of $f_{\tilde{h}}$ imitating f_h .

Then the implicit topology loss can be simply noted as:

$$\mathcal{L}_{it} = \text{MSE}(\tilde{\mathcal{H}}, \tilde{\mathcal{G}}; \sigma^{\tilde{h}}, k) \quad (5)$$

With aforementioned, the overall loss function is:

$$\mathcal{L} = \mathcal{L}_{reg} + \alpha \mathcal{L}_{it} + \beta (\|\sigma^h\|_2^2 + \|\sigma^{\tilde{h}}\|_2^2) \quad (6)$$

Note that there are two σ 's in Eq.6, one for \mathcal{L}_{reg} and another for \mathcal{L}_{it} . This design decouples two learning tasks: predicting *accurate* landmark location by pushing σ as small as possible and estimation of an *approximate* region for the missing landmarks in a more forgiving fashion with a moderate σ .

The network is trained to minimize the loss function Eq. 6 in an end-to-end manner.

3. RESULTS

Dataset & Training Details. We collected 207 thoracic computed tomography angiography (CTA) scans. CTA scans were performed with electrocardiogram (ECG) gating during the administration of iodinated intravenous contrast, with images reconstructed at 75% of the R-R interval. All images are pre-cropped from just above the aortic arch through the upper abdomen (i.e., celiac artery). This dataset includes various aortic pathologies including aneurysm and dissection in patients with and without aortic endografts. The average volume size is $230 \times 230 \times 440$ with a voxel spacing of $0.64 \times 0.64 \times 0.75 \text{ mm}^3$. Images were resampled to have isotropic spacing of $1.8 \times 1.8 \times 1.8 \text{ mm}^3$. We center-crop/pad input images along the three dimensions to create a $96 \times 96 \times 176$

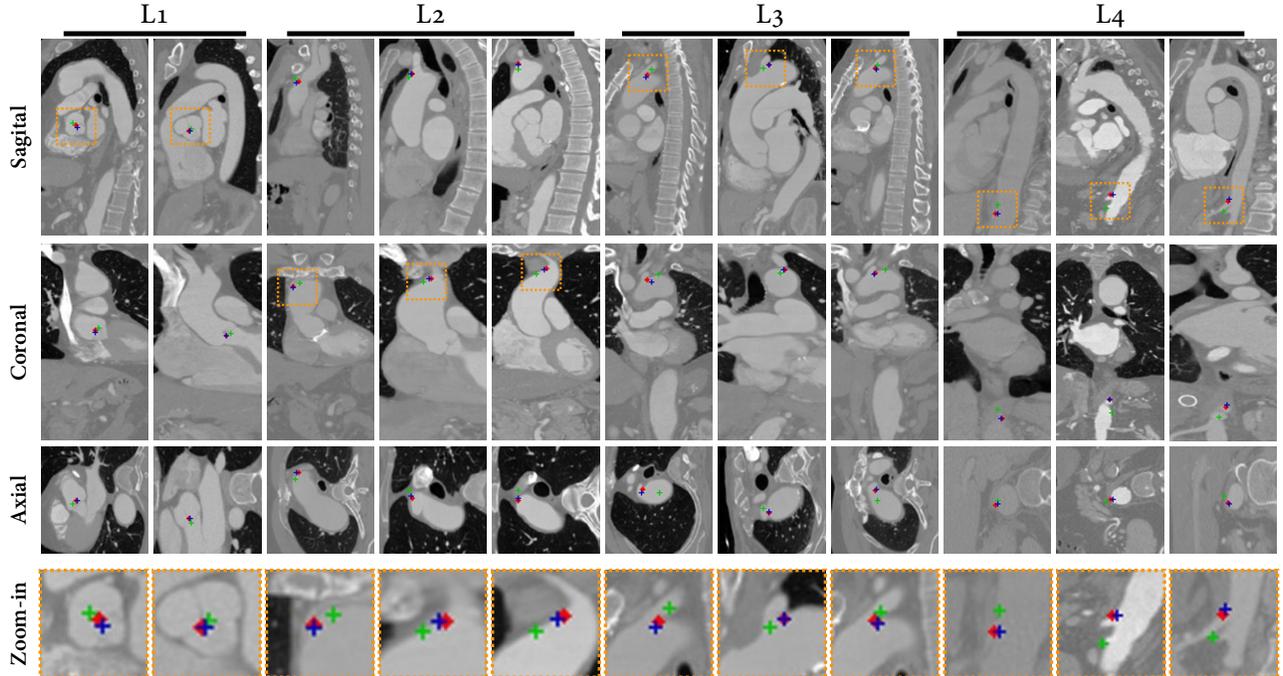


Figure 3: Qualitative result on unseen examples. For clarity, we only visualize the result of two best models, namely ResUnet+Lit and ResUnet, for comparison. $\blacksquare \leftrightarrow$ ground truth location; $\blacksquare \leftrightarrow$ ResUnet+Lit; $\blacksquare \leftrightarrow$ ResUnet.

patch, which is also the size of target heatmaps. We then clip the image intensity $[-1000,1000]$ and normalize to $[0,1]$. Image data was augmented by adding Gaussian noise and random rotations during training.

For training the networks, we used the loss function in Eq. 6 with parameters $\alpha = 0.1, \beta = 10^{-4}, A = 10^6$ that were empirically determined, with $M = 4$ landmarks in total. We initialized σ with 10 voxels. Using the selection module, heatmaps were randomly erased with a probability p , which linearly increased from 0 to 0.5 during over the entire training procedure. We used the AdamW optimizer and CyclicLR ($base_lr = 2 \times 10^{-3}, max_lr = 10^{-2}$) as a learning rate scheduler, with a mini-batch size of 2 and 20,000 total iterations. The training and inference code were implemented in `pytorch`. Using 5-fold cross-validation, 80% of the subjects were randomly chosen for training and the remaining were used for validation. A GTX 2080Ti GPU was used for training and inference.

We used three Unet-like architectures: vanilla Unet,⁸ ResUnet,⁹ and UnetR.¹⁰ Most recently, UnetR combined the original Unet with a transformer module to achieve state-of-the-art performance on several 3D medical image segmentation benchmarks. Since all of these architectures are proposed to predict segmentation, i.e., pixel-wise class label, we can simply adapt these architectures as a backbone in our *implicit topology learning* framework to perform landmark localization tasks, where the number of output channels is equal to the number of landmarks rather than number of label classes. The light-weight head is a smaller version of ResUnet, consisting of three scales with the number of channels (16,32,64). We use a strided convolution of $stride = 4$ between each scale. We also trained larger models ("Unet-L" and "ResUnet-L") with 16 more channels for each scale to investigate the effect of increasing model parameters on performance.

The primary evaluation metric was landmark localization accuracy, defined as the Euclidean distance between predicted landmarks and ground-truth landmarks.

Quantitative & Qualitative Results. The mean and SD of the landmark localization error is reported in Table 1. Cumulative landmark error distributions in Fig.4. Qualitative examples are illustrated in in Fig.3. It can be observed in Table.1 that UnetR suffered from overfitting due to the larger model capacity. Interestingly, when combining UnetR with the implicit topology learning task, the overfitting problem is alleviated and performance improves by 10% with only a 0.4% increase of parameters.

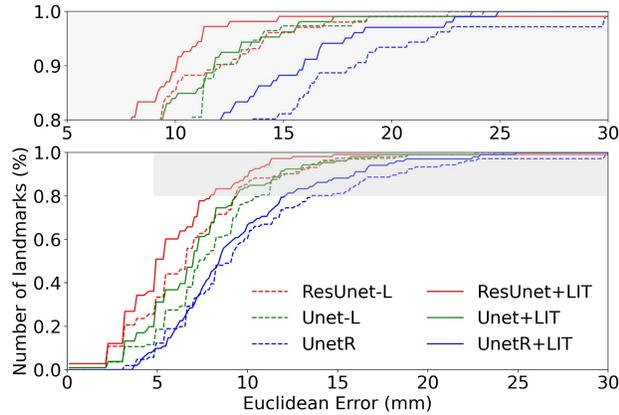


Figure 4: Cumulative landmark localization error. Range (0.8-1.0) is zoomed in for better visibility.

Table 1: Landmark localization errors. **Red** is the best over column. **Bold** is the best for certain backbone. “-L” indicates the model with a larger number of parameters. “LIT” indicates the model that learns implicit topology.

Method	Euclidean Error in (mm)		# params
	median	mean \pm std	
Unet ⁸	7.30	8.04 ± 3.57	2.0 M
Unet-L	7.28	7.98 ± 3.59	2.7 M
Unet + LIT	6.76	7.38 ± 3.58	2.4 M
ResUnet ⁹	6.63	7.01 ± 3.96	4.8 M
ResUnet-L	6.60	6.98 ± 4.10	6.3 M
ResUnet + LIT	4.92	5.93 ± 4.17	5.2 M
UnetR ¹⁰	9.07	10.40 ± 5.71	92.5 M
UnetR + LIT	8.46	9.35 ± 4.27	92.9 M

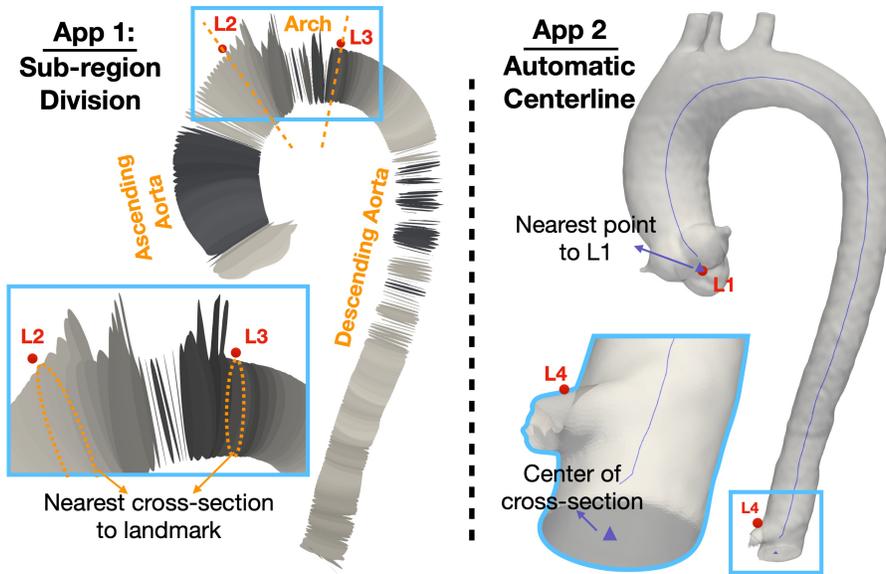


Figure 5: Two applications of aortic landmark localization: sub-region division and automatic centerline generation. When applied to divide the aorta into sub-regions, L2 and L3 are used to determine the boundary cross-sections which divides the aorta into three parts: ascending, arch, and descending aorta (App 1). When applied to automatic centerline generation, L1 and L4 are used to determine the start and end seed for generating the centerline (App 2).

Clinical Applications. Aortic centerline generation is a necessary step for most clinical and research analyses of aortic geometry. However, the vast majority of centerline algorithms require manual interaction (e.g. placement of starting and ending seed-points). Based on our method, we show this task can be fully automated without human intervention.

To automate centerline generation, we ran our previously developed aorta segmentation network¹¹ to obtain an aorta mask, and then applied the method developed in this work to obtain the L1 (ascending root) and L4 (descending celiac) landmarks, which are taken as the seeds by Vascular Modeling Toolkit (VMTK) to reconstruct the 3D mesh and centerline (Figure 5). This pipeline is fully implemented in python.

Recent work by our group has focused on developing techniques to quantify aortic growth/deformation over time using B-spline-based deformable registration methods^{5,6} to analyze aortic aneurysms. Currently, these techniques consider the whole aorta to compute deformation metrics and statistics. Equipped with the landmark detection method described here, one can automatically subdivide analysis of the entire thoracic aorta into three sub-regions, ascending aorta, arch, descending aorta, and descending root. However, the utility of defining aortic sub-regions is not limited to such novel registration-based growth assessment techniques, but can also be used to yield regional assessments of conventional metrics of aortic disease such as maximal diameter and volume. These unique aortic segments are affected differently by disease and thus the ability to analyze each separately may better allow regional assessments of disease.

4. CONCLUSION

We proposed a simple yet effective learning task to make the network learn the implicit topology of the thoracic aorta. The proposed method can be easily combined with Unet-style backbones and is trainable in an end-to-end manner. Localization accuracy is improved compared to baseline and the overfitting problem is alleviated. We believe this method is broadly applicable, and in future work we plan to apply our method to anatomic landmark localization tasks in other anatomies (e.g., hand joints, spine).

REFERENCES

- [1] Yang, D., Zhang, S., Yan, Z., Tan, C., Li, K., and Metaxas, D., “Automated anatomical landmark detection on distal femur surface using convolutional neural network,” in [2015 IEEE 12th international symposium on biomedical imaging (ISBI)], 17–21, IEEE (2015).
- [2] Oktay, O., Bai, W., Guerrero, R., Rajchl, M., de Marvao, A., O’Regan, D. P., Cook, S. A., Heinrich, M. P., Glocker, B., and Rueckert, D., “Stratified decision forests for accurate anatomical landmark localization in cardiac images,” *IEEE transactions on medical imaging* **36**(1), 332–342 (2016).
- [3] Zhang, J., Liu, M., and Shen, D., “Detecting anatomical landmarks from limited medical imaging data using two-stage task-oriented deep neural networks,” *IEEE Transactions on Image Processing* **26**(10), 4753–4764 (2017).
- [4] Payer, C., Štern, D., Bischof, H., and Urschler, M., “Integrating spatial configuration into heatmap regression based cnns for landmark localization,” *Medical image analysis* **54**, 207–219 (2019).
- [5] Bian, Z., Zhong, J., Hatt, C. R., and Burris, N. S., “A deformable image registration based method to assess directionality of thoracic aortic aneurysm growth,” in [Medical Imaging 2021: Image Processing], **11596**, 115962P, International Society for Optics and Photonics (2021).
- [6] Burris, N., Bian, Z., Dominic, J., Zhong, J., Van Bakel, D., Houben, I., Patel, H., Ross, B., Christensen, G., and Hatt, C., “Vascular deformation mapping as a method for 3d growth mapping during ct surveillance of thoracic aortic aneurysm,” *Radiology* (doi: <http://dx.doi.org/10.7302/246>, 2021).
- [7] Rengier, F., Weber, T. F., Giesel, F. L., Bockler, D., Kauczor, H.-U., and von Tengg-Koblighk, H., “Centerline analysis of aortic ct angiographic examinations: benefits and limitations,” *American Journal of Roentgenology* **192**(5), W255–W263 (2009).
- [8] Ronneberger, O., Fischer, P., and Brox, T., “U-net: Convolutional networks for biomedical image segmentation,” in [International Conference on Medical image computing and computer-assisted intervention], 234–241, Springer (2015).

- [9] Kerfoot, E., Clough, J., Oksuz, I., Lee, J., King, A. P., and Schnabel, J. A., “Left-ventricle quantification using residual u-net,” in [*International Workshop on Statistical Atlases and Computational Models of the Heart*], 371–380, Springer (2018).
- [10] Hatamizadeh, A., Yang, D., Roth, H., and Xu, D., “Unetr: Transformers for 3d medical image segmentation,” *arXiv preprint arXiv:2103.10504* (2021).
- [11] Zhong, J., Bian, Z., Hatt, C. R., and Burris, N. S., “Segmentation of the thoracic aorta using an attention-gated u-net,” in [*Medical Imaging 2021: Computer-Aided Diagnosis*], **11597**, 115970M, International Society for Optics and Photonics (2021).