

Deep Distance Transform for Tubular Structure Segmentation in CT Scans

Yan Wang^{1†} Xu Wei^{2*†} Fengze Liu¹ Jieneng Chen^{3*} Yuyin Zhou¹ Wei Shen^{1‡}
Elliot K. Fishman⁴ Alan L. Yuille¹

¹Johns Hopkins University ²University of California San Diego ³Tongji University

⁴The Johns Hopkins University School of Medicine

Abstract

*Tubular structure segmentation in medical images, e.g., segmenting vessels in CT scans, serves as a vital step in the use of computers to aid in screening early stages of related diseases. But automatic tubular structure segmentation in CT scans is a challenging problem, due to issues such as poor contrast, noise and complicated background. A tubular structure usually has a cylinder-like shape which can be well represented by its skeleton and cross-sectional radii (scales). Inspired by this, we propose a geometry-aware tubular structure segmentation method, Deep Distance Transform (DDT), which combines intuitions from the classical **distance transform** for skeletonization and modern deep segmentation networks. DDT first learns a multi-task network to predict a segmentation mask for a tubular structure and a distance map. Each value in the map represents the distance from each tubular structure voxel to the tubular structure surface. Then the segmentation mask is refined by leveraging the shape prior reconstructed from the distance map. We apply our DDT on six medical image datasets. Results show that (1) DDT can boost tubular structure segmentation performance significantly (e.g., over 13% DSC improvement for pancreatic duct segmentation), and (2) DDT additionally provides a geometrical measurement for a tubular structure, which is important for clinical diagnosis (e.g., the cross-sectional scale of a pancreatic duct can be an indicator for pancreatic cancer).*

1. Introduction

Tubular structures are ubiquitous throughout the human body, with notable examples including blood vessels, pancreatic duct and urinary tract. They occur in specific environments at the boundary of liquids, solids or air and surrounding tissues, and play a prominent role in sustaining physiological functions of the human body.

*This work was done when Xu Wei and Jieneng Chen were at JHU.

†Equal Contribution.

‡Corresponding Author (shenwei1231@gmail.com).

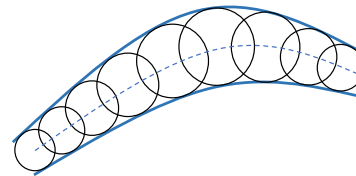


Figure 1. A tubular shape is presented as the envelope of a family of spheres with continuously changing center points and radii [9].

In this paper, we investigate automatic tubular organ/tissue segmentation from CT scans, which is important for the characterization of various diseases [18]. For example, pancreatic duct dilatation or abrupt pancreatic duct caliber change signifies high risk for pancreatic ductal adenocarcinoma (PDAC), which is the third most common cause of cancer death in the US [11]. Another example is that obstructed vessels lead to coronary heart disease, which is the leading cause of death in the US [27].

Segmenting tubular organs/tissues from CT scans is a popular but challenging problem. Existing methods addressing this problem can be roughly categorized into two groups: (1) Geometry-based methods, which build deformable shape models to fit tubular structures by exploiting their geometrical properties [43, 45, 3, 25], e.g., a tubular structure can be well represented by its *skeleton*, aka *symmetry axis* or *medial axis*, and it has a cylindrical surface. But, due to the lack of powerful learning models, these methods cannot deal with poor contrast, noise and complicated background. (2) Learning-based methods, which learn a per-pixel classification model to detect tubular structures. The performance of this type of methods is largely boosted by deep learning, especially fully convolutional networks (FCN) [23, 49, 48]. FCN and its variants have become out-of-the-box models for tubular organ/tissue segmentation and achieve state-of-the-art results [24, 47]. But, these networks simply try to learn a class label per voxel, which inevitably ignores the geometric arrangement of the voxels in a tubular structure, and consequently can not guarantee that the obtained segmentation has the right shape.

Since a tubular structure can be well represented by its skeleton and the cross-sectional radius of each skeleton

point, as shown in Fig. 1, these intrinsic geometric characteristics should be taken into account to serve as a valuable prior. To this end, a straightforward strategy is to first train a model, *e.g.*, a deep network, to directly predict whether each voxel is on the skeleton of the tubular structure or not as well as the cross-sectional radius of each skeleton point, and then reconstruct the segmentation of the tubular structure from its skeleton and radii [34]. However, such a strategy has severe limitations: (1) The ground-truth skeletons used for training are not easily obtained. Although they can be approximately computed from the ground-truth segmentation mask by 3D skeletonization methods, skeleton extraction from 3D mesh representation itself is a hard and unsolved problem [5]. Without reliable skeleton ground-truths, the performance of tubular structure segmentation cannot be guaranteed. (2) It is hard for the classifier to distinguish voxels on the skeleton itself from those immediately next to it, as they have similar features but different labels.

To tackle the obstacles mentioned above, we propose to perform tubular structure segmentation by training a multi-task deep network to predict not only a segmentation mask, but also a *distance map*, consisting of the distance transform value from each tubular structure voxel to the tubular structure surface, rather than a single skeleton/non-skeleton label. *Distance transform* [28] is a classical image processing operator to produce a distance map with the same size of the input image, each value in which is the distance from each foreground pixel/voxel to the foreground boundary. Distance transform is also known as the basis of one type of skeletonization algorithms [17], *i.e.*, the ridge of the distance map is the skeleton. Thus, the predicted distance map encodes the geometric characteristics of the tubular structure. This motivated us to design a geometry-aware approach to refine the output segmentation mask by leveraging the shape prior reconstructed from the distance map. Essentially, our approach performs tubular structure segmentation by an implicit skeletonization-reconstruction procedure with no requirements for skeleton ground-truths. We stress that the distance transform brings two benefits for our approach: (1) Distance transform values are defined on each voxel inside a tubular structure, which eliminates the problem of the discontinuity between the skeleton and its surrounding voxels; (2) distance transform values on the skeleton (the ridge of the distance map) are exactly the cross-sectional radii (scales) of the tubular structure, which is an important **geometrical measurement**. To make the distance transform value prediction more precise, we additionally propose a distance loss term used for network training, which indicates a penalty when predicted distance transform value is far away from its ground-truth.

We term our method *Deep Distance Transform* (DDT), as it naturally combines intuitions from the classical distance transform for skeletonization and modern deep seg-

mentation networks. We emphasize that DDT has two advantages over vanilla segmentation networks: (1) It guides tubular structure segmentation by taking the geometric property of tubular structures into account. This reduces the difficulty to segment tubular structures from complex surrounding structures and ensures that the segmentation results have a proper shape prototype; (2) It predicts the cross-sectional scales of a tubular structure as by-products, which are important for the further study of the tubular structure, such as clinical diagnosis and virtual endoscopy [7].

We verify DDT on six datasets, including five datasets for segmentation task, and one dataset for clinical diagnosis. For segmentation task, the performance of our DDT exceeds all backbone networks by a large margin, with even over 13% improvement in terms of Dice-Sørensen coefficient for pancreatic duct segmentation on the famous 3D-Unet [12]. The ablation study further shows the effectiveness of each proposed module in DDT. The experiment for clinical diagnosis leverages dilated pancreatic duct as cue for finding missing PDAC tumors by original deep networks, which verifies the potential of our DDT for early diagnosis of pancreatic cancer.

2. Related Work

2.1. Tubular Structure Segmentation

2.1.1 Geometry-based Methods

Various methods have been proposed to improve the performance of tubular structure segmentation by considering the geometric characteristics, and a non-exhaustive overview is given here. (1) Contour-based methods extracted the segmentation mask of a tubular structure by means of approximating its shape in the cross-sectional domain [1, 10]. (2) Minimal path approaches conducted tubular structure tracking and were usually interactive. They captured the global minimum curve (energy weighted by the image potential) between two points given by the user [9]. (3) Model-based tracking methods required to refine a tubular structure model, which most of the time adopted a 3D cylinder with elliptical or circular section. At each tracking step, they calculated the new model position by seeking for the optimal model match among all possible new model positions [8]. (4) Centerline based methods found the centerline and estimated the radius of linear structures. For example, multiscale centerline detection method proposed in [34] adopted the idea of distance transform, and reformulated centerline detection and radius estimation in terms of a regression problem in 2D. Our work fully leverages the geometric information of a tubular structure, proposing a distance transform algorithm to implicitly learn the skeleton and cross-sectional radius, and the final segmentation mask is reconstructed by adopting the shape prior of the tubular structure.

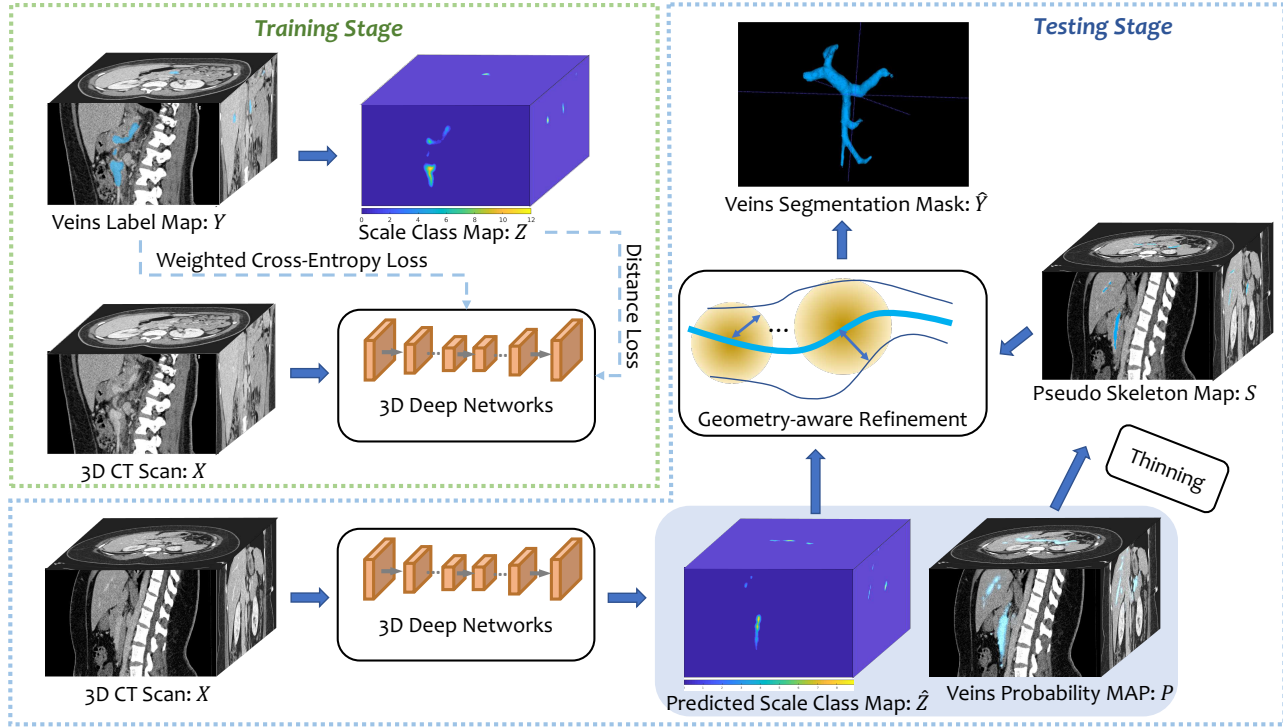


Figure 2. The training and testing stage of DDT, illustrated on an example of veins segmentation. Our DDT has two head branches: the first one is targeting on the ground-truth label map, which performs per-voxel veins/non-veins classification, and the second head branch is targeting on the scale class map, which performs scale prediction for veins voxels. Then a geometry-aware refinement approach is proposed to leverage the shape prior obtained from the scale class map and the pseudo skeleton map to refine the segmentation mask.

2.1.2 Learning-based Method

Learning-based method for tubular structure segmentation infers a rule from labeled training pairs, one for each pixel. Traditional methods such as 2-D Gabor wavelet and classifier combination [35], ridge-based segmentation [36], and random decision forest based method [2] achieved considerable progress. In the past years, various 2D and 3D deep segmentation networks have become very popular. Some multi-organ segmentation methods [41, 29] were proposed to segment multiple organs simultaneously, including tubular organs. DeepVessel [16] put a four-stage HED-like CNN and conditional random field into an integrated deep network to segment retinal vessel. Kid-Net [37], inspired from 3D-Unet [12], was a two-phase 3D network for kidney vessels segmentation. ResDSN [50, 51] and 3D-Unet [12] were used in Hyper-pairing network [47] to segment tissues in pancreas including duct by combining information from dual-phase imaging. Besides, 3D-HED and its variant were applied for vascular boundary detection [24]. Other scenarios such as using synthetic data to improve endotracheal tube segmentation [15]. Cross-modality domain adaptation framework with adversarial learning which dealt with the domain shift in segmenting biomedical images including as-

cending aorta was also proposed [13].

2.2. Learning-based Skeleton Extraction

Learning-based skeleton extraction from natural images has been widely studied in recent decades [38, 31, 34, 22, 21] and achieved promising progress with the help of deep learning [32, 20, 46, 40]. Shen *et al.* [32] showed that multi-task learning, *i.e.*, jointly learning skeleton pixel classification and skeleton scale regression, was important to obtain accurate predicted scales, and it was useful for skeleton-based object segmentation. One recent work for bronchus segmentation calculated the distance from each voxel to its nearest skeleton point [39].

However, these methods cannot be directly applied to our tasks, since they require the skeleton ground-truth, which is not easy to obtain from a tiny and highly distorted 3D mask due to the commonly existed annotation errors for medical images [42].

3. Methodology

We first define a 3D volume X of size $L \times W \times H$ as a function on the coordinate set $V = \{\mathbf{v} | \mathbf{v} \in N_L \times N_W \times N_H\}$, *i.e.*, $X : V \rightarrow R \subset \mathbb{R}$ where the value on position

\mathbf{v} is defined as $x_{\mathbf{v}} = X(\mathbf{v})$. N_L, N_W, N_H represent for the integer set ranging from 1 to L, W, H respectively, so that the Cartesian product of them can form the coordinate set. Given a 3D CT scan X , the goal of tubular structure segmentation is to predict the label \hat{Y} of all voxels in the CT scan, where $\hat{y}_{\mathbf{v}} \in \{0, 1\}$ denotes the predicted label for each voxel at position \mathbf{v} , *i.e.*, if the voxel at \mathbf{v} is predicted as a tubular structure voxel, then $\hat{y}_{\mathbf{v}} = 1$, otherwise $\hat{y}_{\mathbf{v}} = 0$. We also use \mathbf{v} to denote the voxel at position \mathbf{v} in the remaining of the paper for convenience sake. Fig. 2 illustrates our tubular structure segmentation network, *i.e.*, DDT.

3.1. Distance Transform for Tubular Structure

In this section, we discuss how to perform distance transform for tubular structure voxels. Given the ground-truth label map Y of the CT scan X in the training phase, let C_V be the set of voxels on the tubular structure surface, which can be defined by

$$C_V = \{\mathbf{v} | y_{\mathbf{v}} = 1, \exists \mathbf{u} \in \mathcal{N}(\mathbf{v}), y_{\mathbf{u}} = 0\}, \quad (1)$$

where $\mathcal{N}(\mathbf{v})$ denotes the 6-neighbour voxels of \mathbf{v} . Then, by performing distance transform on the CT scan X , the distance map D is computed by

$$d_{\mathbf{v}} = \begin{cases} \min_{\mathbf{u} \in C_V} \|\mathbf{v} - \mathbf{u}\|_2, & \text{if } y_{\mathbf{v}} = 1 \\ 0, & \text{if } y_{\mathbf{v}} = 0 \end{cases}. \quad (2)$$

Note that, for each tubular structure voxel \mathbf{v} , the distance transform assigns it a distance transform value which is the nearest distance from \mathbf{v} to the tubular structure surface C_V . Here we use Euclidean distance, as skeletons from Euclidean distance maps are robust to rotations [4].

We further quantize each $d_{\mathbf{v}}$ into one of K bins by rounding $d_{\mathbf{v}}$ to the nearest integer, which converts the continuous distance map D to a discrete quantized distance map Z , where $z_{\mathbf{v}} \in \{0, \dots, K\}$. We do this quantization, because training a deep network directly for regression is relatively unstable, since outliers, *i.e.*, the commonly existed annotation errors for medical images [42], cause a large error term, which makes it difficult for the network to converge and leads to unstable predictions [30]. Based on quantization, we rephrase the distance prediction problem as a classification problem, *i.e.*, to determine the corresponding bin for each quantized distance. We term the K bins of the quantized distances as K scale classes. We use the term scale since the distance transform values at the skeleton voxels of a tubular structure are its cross-sectional scales.

3.2. Network Training for Deep Distance Transform

Given a 3D CT scan X and its ground-truth label map Y , we can compute its *scale class map* (quantized distance map) Z according to the method given in Sec. 3.1. In this

section, we describe how to train a deep network for tubular structure segmentation by targeting on both Y and Z .

As shown in Fig. 2, our DDT model has two head branches. The first one is targeting on the ground-truth label map Y , which performs per-voxel classification for semantic segmentation with a weighted cross-entropy loss function \mathcal{L}_{cls} :

$$\mathcal{L}_{\text{cls}} = - \sum_{\mathbf{v} \in V} \left(\beta_p y_{\mathbf{v}} \log p_{\mathbf{v}}(\mathbf{W}, \mathbf{w}_{\text{cls}}) + \beta_n (1 - y_{\mathbf{v}}) \log (1 - p_{\mathbf{v}}(\mathbf{W}, \mathbf{w}_{\text{cls}})) \right), \quad (3)$$

where \mathbf{W} is the parameters of the network backbone, \mathbf{w}_{cls} is the parameters of this head branch and $p_{\mathbf{v}}(\mathbf{W}, \mathbf{w}_{\text{cls}})$ is the probability that \mathbf{v} is a tubular structure voxel as predicted by this head branch. $\beta_p = \frac{0.5}{\sum_{\mathbf{v}} y_{\mathbf{v}}}$ and $\beta_n = \frac{0.5}{\sum_{\mathbf{v}} (1 - y_{\mathbf{v}})}$ are loss weights for tubular structure and background classes respectively.

The second head branch is predicting on the scale class map Z , which performs scale prediction for tubular structure voxels (*i.e.*, $z_{\mathbf{v}} > 0$). We introduce a new distance loss function \mathcal{L}_{dis} to learn this head branch:

$$\mathcal{L}_{\text{dis}} = -\beta_p \sum_{\mathbf{v} \in V} \sum_{k=1}^K \left(\mathbf{1}(z_{\mathbf{v}} = k) \left(\log g_{\mathbf{v}}^k(\mathbf{W}, \mathbf{w}_{\text{dis}}) + \lambda \omega_{\mathbf{v}} \log (1 - \max_l g_{\mathbf{v}}^l(\mathbf{W}, \mathbf{w}_{\text{dis}})) \right) \right), \quad (4)$$

where \mathbf{W} is the parameters of the network backbone, \mathbf{w}_{dis} is the parameters of the second head branch, $\mathbf{1}(\cdot)$ is an indication function, λ is a trade-off parameter which balances the two loss terms (we simply set $\lambda = 1$ in our implementation), $g_{\mathbf{v}}^k(\mathbf{W}, \mathbf{w}_{\text{dis}})$ is the probability that the scale of \mathbf{v} belongs to k -th scale class and $\omega_{\mathbf{v}}$ is a normalized weight defined by $\omega_{\mathbf{v}} = \frac{|\arg \max_l g_{\mathbf{v}}^l(\mathbf{W}, \mathbf{w}_{\text{dis}}) - z_{\mathbf{v}}|}{K}$. Note that, the first term of Eq. 4 is the standard softmax loss which penalizes the classification error for each scale class equally. The second term of Eq. 4 is termed as *distance loss term*, which penalizes the difference between each predicted scale class (*i.e.*, $\max_l g_{\mathbf{v}}^l(\mathbf{W}, \mathbf{w}_{\text{dis}})$) and its ground-truth scale class $z_{\mathbf{v}}$, where the penalty is controlled by $\omega_{\mathbf{v}}$. Finally, the loss function for our segmentation network is $\mathcal{L} = \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{dis}}$ and the optimal network parameters are obtained by $(\mathbf{W}^*, \mathbf{w}_{\text{cls}}^*, \mathbf{w}_{\text{dis}}^*) = \arg \min_{\mathbf{W}, \mathbf{w}_{\text{cls}}, \mathbf{w}_{\text{dis}}} \mathcal{L}$.

3.3. Geometry-aware Refinement

Shape reconstruction from skeletons is a simple and well-known operation, which often achieves better segmentation performance than pure segmentation methods [32]. Inspired by this, we propose a soft version of such reconstruction, termed as geometry-aware refinement (GAR), where skeletons are obtained by thinning probability maps

and maximal balls centered at skeleton points are softened to Gaussian kernels. GAR ensures smoothness between similar voxels, and spatial and appearance consistency of the segmentation output, especially from clutter background.

Given a 3D CT scan X in the testing phase, for each voxel \mathbf{v} , our tubular structure segmentation network, DDT, outputs two probabilities, $p_{\mathbf{v}}(\mathbf{W}^*, \mathbf{w}_{\text{cls}}^*)$, which is the probability that \mathbf{v} is a tubular structure voxel and $g_{\mathbf{v}}^k(\mathbf{W}^*, \mathbf{w}_{\text{dis}}^*)$, which is the probability that the scale of \mathbf{v} belongs to k -th scale class. For notational simplicity, we use $p_{\mathbf{v}}$ and $g_{\mathbf{v}}^k$ to denote $p_{\mathbf{v}}(\mathbf{W}^*, \mathbf{w}_{\text{cls}}^*)$ and $g_{\mathbf{v}}^k(\mathbf{W}^*, \mathbf{w}_{\text{dis}}^*)$, respectively, in the rest of the paper. $p_{\mathbf{v}}$ provides per-voxel tubular structure segmentation, and $g_{\mathbf{v}}^k$ encodes the geometric characteristics of the tubular structure. Our GAR obtains the final segmentation result by refining $p_{\mathbf{v}}$ according to $g_{\mathbf{v}}^k$. This approach is shown in Fig. 2 and is processed as follows:

- a. **Pseudo skeleton generation.** The probability map P is thinned by thresholding it to generate a binary pseudo skeleton map S for the tubular structure. If $p_{\mathbf{v}} > T^p$, $s_{\mathbf{v}} = 1$; otherwise, $s_{\mathbf{v}} = 0$, and T^p is the threshold.
- b. **Shape reconstruction.** For each voxel \mathbf{v} , its predicted scale $\hat{z}_{\mathbf{v}}$ is given by $\hat{z}_{\mathbf{v}} = \arg \max_k g_{\mathbf{v}}^k$. We fit a Gaussian kernel to soften each ball and obtain a soft reconstructed shape \tilde{Y}^s :

$$\tilde{y}_{\mathbf{v}}^s = \sum_{\mathbf{u} \in \{\mathbf{u}' | s_{\mathbf{u}'} > 0\}} c_{\mathbf{u}} \Phi(\mathbf{v}; \mathbf{u}, \Sigma_{\mathbf{u}}), \quad (5)$$

where $\Phi(\cdot)$ denotes the density function of a multivariate normal distribution, \mathbf{u} is the mean and $\Sigma_{\mathbf{u}}$ is the co-variance matrix. According to the 3-sigma rule, we set $\Sigma_{\mathbf{u}} = (\frac{\hat{z}_{\mathbf{u}}}{3})^2 \mathbf{I}$, where \mathbf{I} is an identity matrix. We notice that the peak of $\Phi(\cdot; \mathbf{u}, \Sigma_{\mathbf{u}})$ becomes smaller if $\hat{z}_{\mathbf{u}}$ is larger. To normalize the peak of each normal distribution, we introduce a normalization factor $c_{\mathbf{u}} = \sqrt{(2\pi)^3 \det(\Sigma_{\mathbf{u}})}$.

- c. **Segmentation refinement.** We use the soft reconstructed shape \tilde{Y}^s to refine the segmentation probability $p_{\mathbf{u}}$, which results in a refined segmentation map \tilde{Y}^r :

$$\tilde{y}_{\mathbf{v}}^r = \sum_{\mathbf{u} \in \{\mathbf{u}' | s_{\mathbf{u}'} > 0\}} p_{\mathbf{u}} c_{\mathbf{u}} \Phi(\mathbf{v}; \mathbf{u}, \Sigma_{\mathbf{u}}). \quad (6)$$

The final segmentation mask \hat{Y} is obtained by thresholding \tilde{Y}^r , i.e., if $\tilde{y}_{\mathbf{v}}^r > T^r$, $\hat{y}_{\mathbf{v}} = 1$, otherwise, $\hat{y}_{\mathbf{v}} = 0$, where $\tilde{y}_{\mathbf{v}}^r$ and $\hat{y}_{\mathbf{v}}$ are the value of voxel at position \mathbf{v} of \tilde{Y}^r and \hat{Y} , respectively.

As mentioned in Sec. 1, the predicted scale $\hat{z}_{\mathbf{v}}$ is a geometrical measurement for a tubular structure, which is essential for clinical diagnosis. We will show one clinical application in Sec. 4.2.

4. Experiments

In this section, we conduct the following experiments: we first evaluate our approach on five segmentation datasets, including (1) the dataset used in [47], (2) three tubular structure datasets created by radiologists in our team, and (3) hepatic vessels dataset in Medical Segmentation Decathlon (MSD) challenge [33]. Then, as we mentioned in Sec. 1, our DDT predicts cross-sectional scales as by-products, which are important for applications such as clinical diagnosis. We show that the cross-sectional scale is an important measurement for predicting the dilation degree of a pancreatic duct, which can help find the PDAC tumors missed in [51], without increasing the false positives.

4.1. Tubular Structure Segmentation

4.1.1 Implementation Details and Evaluation Metric

Our implementation is based on PyTorch. For data pre-processing, followed by [47], we truncate the raw intensity values within the range of $[-100, 240]$ HU and normalize each CT scan into zero mean and unit variance. Data augmentation (i.e., translation, rotation and flipping) is conducted in all the methods, leading to an augmentation factor of 24. During training, we randomly sample patches of a specified size (i.e., 64) due to memory issue. We use exponential learning rate decay with $\gamma = 0.99$. During testing, we employ the sliding window strategy to obtain the final predictions. The groundtruth distance map for each tubular structure is computed by finding the euclidean distance of each foreground voxel to its nearest boundary voxels. The segmentation accuracy is measured by the well-known Dice-Sørensen coefficient (DSC) in the rest of the paper, unless otherwise specified.

4.1.2 The PDAC Segmentation Dataset [47]

We first study the PDAC segmentation dataset [47] which has 239 patients with pathologically proven PDAC. All CT scans are contrast enhanced images and our experiments are conducted on only portal venous phase. We follow the same setting and the same cross-validation as reported in [47]. DSCs for three structures were reported in [47]: abnormal pancreas, PDAC mass and pancreatic duct. We only show the average and standard deviation over all cases for pancreatic duct, which is a tubular structure.

Results and Discussions. To evaluate the performance of the proposed DDT framework, we compare it with a per-voxel classification method [47], termed as **SegBaseline** in Table 1. It can be seen that our approach outperforms the baseline reported in [47] by a large margin. It is also worth mentioning that although our DDT is only tested on venous phase, the performance is comparable with the hyper-

Table 1. Performance comparison (DSC, %) on pancreatic duct segmentation (mean \pm standard deviation of all cases). SegBaseline stands for per-voxel classification. Multi-phase HPN is a hyper-paring network combining CT scans from both **venous** (V) and **arterial** (A) phases. Noted that only CT scans in **venous** phase are used for SegBaseline and DDT. **Bold** denotes the best results.

Methods	Phase	Backbone Networks	
		3D-UNet	ResDSN
SegBaseline [47]	V	40.25 \pm 27.89	49.81 \pm 26.23
Multi-phase HPN [47]	A+V	44.93 \pm 24.88	56.77 \pm 23.33
DDT (Ours)	V	58.20 \pm 23.39	55.97 \pm 24.76

paring network [47] (*i.e.*, Multi-phase HPN), which integrates multi-phase information (*i.e.*, arterial phase and venous phase). For 3D-UNet, our DDT even outperforms the multi-phase method by more than 13% in terms of DSC.

Ablation Study. We conduct ablation experiments on the PDAC segmentation dataset, using ResDSN as the backbone. These variants of our methods are considered:

- SegfromSkel: This is the straightforward strategy mentioned in Sec. 1 for skeleton-based tubular structure segmentation, *i.e.*, segmenting by reconstructing from the predicted skeleton. The ground-truth skeleton is obtained by the mesh contraction algorithm [5], and the scale of each skeleton point is defined as its shortest distance to the duct surface. We use the same method in Sec. 3 to instantiate this strategy, but the learning target is the skeleton instead of the duct mask.
- DDT $\lambda = 0$, w/o GAR: DDT **without** distance loss term ($\lambda = 0$ in Eq. 4), and **without** geometry-aware refinement.
- DDT $\lambda = 0$, w/ GAR: DDT **without** distance loss term, and **with** geometry-aware refinement.
- DDT $\lambda = 1$, w/o GAR: DDT **with** distance loss term, and **without** geometry-aware refinement.
- DDT $\lambda = 1$, w/ GAR: DDT **with** distance loss term, and **with** geometry-aware refinement.

The results of the ablation experiments are summarized in Table 2. We also show examples of the predicted duct for better understanding how each component (*i.e.*, distance loss term and geometry-aware refinement) learns the geometry information in the supplementary material. Then, we aim at discussing parameters in the geometry-aware refinement component. In our implementation, we set $T^p = 0.98$ and $T^r = 0.5$ in Sec. 3.3. Now we vary each of them and fix the other one to the default value to see how the performance changes. As shown in Fig. 3(a), setting a larger T^p leads to better performance. This phenomenon further verifies the advantage of leveraging scale class map to refine the per-voxel segmentation results, *i.e.*, a thinner pseudo skeleton combined with a scale class map can better represent a

Table 2. Ablation study of pancreatic duct segmentation using ResDSN as backbone network. GAR indicates the proposed geometry-aware refinement.

Method	Average DSC (%)
SegBaseline [47]	49.81
SegfromSkel	51.88
DDT $\lambda = 0$, w/o GAR	52.73
DDT $\lambda = 0$, w/ GAR	54.70
DDT $\lambda = 1$, w/o GAR	53.69
DDT $\lambda = 1$, w/ GAR	55.97

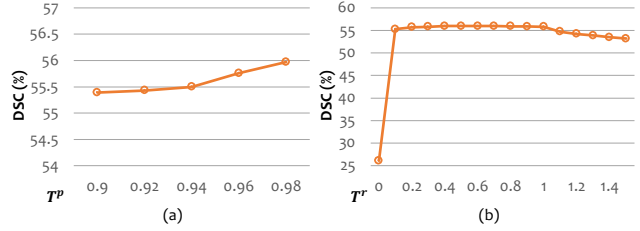


Figure 3. Performance changes by varying (a) pseudo skeleton generation parameter T^p and (b) segmentation refinement parameter T^r .

tubular structure. Fig. 3(b) shows that the performance is not sensitive within the range of $T^r \in [0.1, 1]$.

4.1.3 Tubular Structure Datasets

We then evaluate our algorithm on multiple tubular structures datasets. Radiologists in our team collected 229 abdominal CT scans of normal cases with aorta annotation, 204 normal cases with veins annotation, and 494 abdominal CT scans of biopsy-proven PDAC cases with pancreatic duct annotation. All these three datasets are under IRB approved protocol.

We conduct experiments by comparing our DDT with SegBaseline on three backbone networks: 3D-HED [24], 3D-UNet [12] and ResDSN [50, 51]. The results in terms of DSC and mean surface distance (in mm) are reported in Table 3. SegBaseline methods on all backbone networks are significantly lower than our approach. In particular, for 3D-HED, DDT outperforms SegBaseline by 8%, making its strong ability in segmenting small tubular structures like pancreatic duct in medical images. The results are obtained by cross-validation. We also illustrate segmentation results of aorta and veins in Fig. 4 for qualitative comparison. We can see that compared with SegBaseline, DDT captures geometry information, which is more robust to the noise and complicated background.

4.1.4 Hepatic Vessels Dataset in MSD Challenge

We also test our DDT on a public hepatic vessels dataset in MSD challenge [33]. There are two targets in hepatic vessels dataset: vessels and tumor. As our goal is to segment

Table 3. Performance comparison (in average DSC, % and mean surface distance in mm) on three tubular structure datasets by using different backbones. “↑” and “↓” indicate the larger and the smaller the better, respectively. **Bold** denotes the best results for each tubular structure per measurement.

Backbone	Methods	Aorta		Veins		Pancreatic duct	
		Average DSC ↑	Mean Surface Distance ↓	Average DSC ↑	Mean Surface Distance ↓	Average DSC ↑	Mean surface Distance ↓
3D-HED [24]	SegBaseline	90.85	1.15	73.57	5.13	46.43	7.06
	DDT	92.94	0.82	76.20	3.78	54.43	4.91
3D-UNet [12]	SegBaseline	92.01	0.94	71.57	4.46	56.63	3.64
	DDT	93.30	0.61	75.59	4.07	62.31	3.56
ResDSN [50]	SegBaseline	89.89	1.12	71.10	6.25	55.91	4.24
	DDT	92.57	1.10	76.60	5.03	59.29	4.19

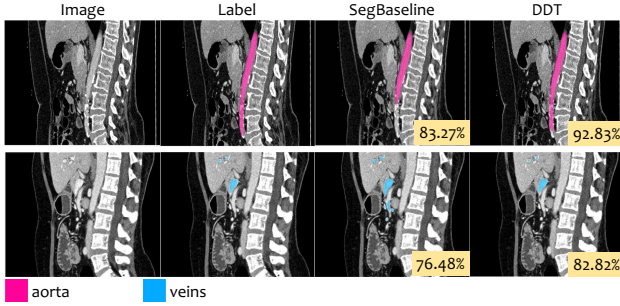


Figure 4. Illustration of aorta (upper row) and veins (lower row) segmentation results for selected example images. Numbers on the bottom right show segmentation DSCs.

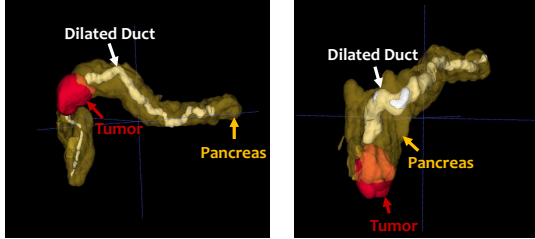


Figure 5. Examples of PDAC cases. In most PDAC cases, the tumor blocks the duct and causes it to dilate.

tubular structure, we aim at vessel segmentation. Although this challenge is over, it is still open for submissions. We train our DDT on 303 training cases, and submit vessel predictions of the testing cases to the challenge.

We simply use ResDSN [51] as our backbone network, and follow the same data augmentation as introduced in Sec. 4.1.1. We summarize some leading quantitative results reported in the leaderboard in Table 4. This comparison shows the effectiveness of our DDT.

4.2. Finding PDAC Tumor by Dilated Duct

Background. PDAC is one of the most deadly disease, whose survival is dismal as when it comes to the time of diagnosis, there are more than half of patients have evidence of metastatic disease. As mentioned in Sec. 1, dilated duct is a vital cue for the presence of a PDAC tumor. The rea-

Table 4. Comparison to competing submissions of MSD challenge: <http://medicaldecathlon.com>

Methods	Average DSC (%)
DDT (Ours)	63.43
nnU-Net [19]	63.00
UMCT [44]	63.00
K.A.V.athlon	62.00
LS Wang’s Group	55.00
MIMI	60.00
MPUNet [26]	59.00

son lies in that in most cases, the tumor blocks the duct and causes it to dilate, as shown in Fig. 5. Experienced radiologists usually trace the duct from the pancreas tail onward to see if there exists a truncated duct. If they see the predicted duct pattern as illustrated in Fig. 5, they will be alarmed and treat it as a suspicious PDAC case. For computer-aided diagnosis, given a mixture of normal and abnormal CT scans (PDAC cases), if some voxels are segmented as a tumor by a state-of-the-art deep network, we can provide radiologists with tumor locations [51]. But, as reported [51], even a state-of-the-art deep network failed to detect 8 PDACs out of 136 abnormal cases. As emphasized in [51], for clinical purposes, we shall guarantee a *high sensitivity* with a reasonable specificity. Then how can we use dilated duct as a cue to help find the PDAC tumor in an abnormal case even if it does NOT have any PDAC tumor prediction by directly applying deep networks?

Clinical Workflow. The flowchart of our strategy is illustrated in Fig. 6. We apply our DDT on the cases which do not have tumor prediction by [51]. Then the predictions of DDT are processed as follows:

1. **Find cases with predicted dilated duct.** Let’s assume a case has N predicted duct voxels. If $N = 0$, then we regard this case as negative. If $N > 0$, let’s denote the predicted associated scales (radii) are $\{\hat{z}_{v_i}\}_{i=1}^N$. If $\arg \max_i \hat{z}_{v_i} > T^s$, i.e., the largest cross-sectional scale is larger than T^s , we regard this is a dilated duct, and a tumor may present on its head location. Otherwise, we

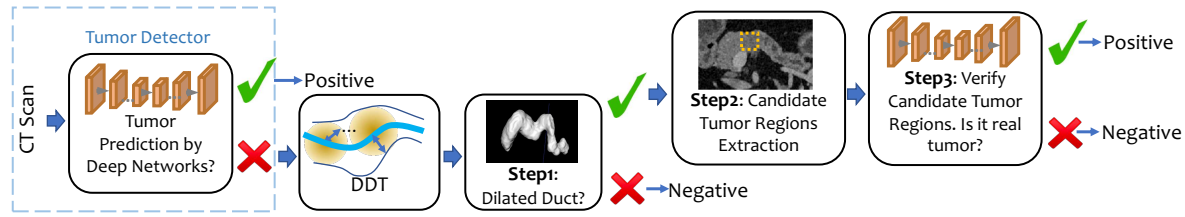


Figure 6. Flowchart of finding missing PDAC tumor by dilated duct.

treat this case as negative. We set $T^s = 3$, since the radius of a dilated duct should be larger than 1.5 mm [14], and the voxel spatial resolution of the dataset [51] is around 0.5 mm^3 .

2. **Extract candidate tumor regions by the location of dilated duct.** We use geodesic distance to find the extreme points of the duct [6]. Then we crop a set of square regions of size \mathbb{R}^3 centered on the extreme points not lying on the tail of the pancreas, since a tumor presenting on the tail of the pancreas will not block a duct. This set of square regions are candidate tumor regions.
3. **Verify candidate tumor regions.** As candidate tumor regions may come from both normal and abnormal cases. We should verify whether the candidate region is a real tumor region. From the training set, we randomly crop regions of size \mathbb{R}^3 around PDAC tumor region as positive training data, and randomly crop regions of size \mathbb{R}^3 from normal pancreas as negative training data. Then we train a ResDSN [51] to verify these candidate tumor regions. We follow the same criterion used in [51] to compute sensitivity and specificity.

Experiment Settings. We follow the same data split as used in [51]. We only test our algorithm on the 8 PDAC cases and 197 normal cases which do not have tumor prediction by [51], aiming at finding missing tumor by dilated duct, while not introducing more false positives. $\mathbb{R} = 48$.

Analysis. We compare our results with those of [51] in Table 5. In our experiment, 4 out of 8 abnormal cases and 3 out of 197 normal cases have predicted dilated duct by step 1. An example is shown in Fig. 7(a). The tubular structure residing inside the *ground-truth* pancreas, right behind the *ground-truth* tumor is our *predicted* dilated duct. This leads to overall 18 3D candidate tumor regions by step 2, shown as the yellow dashed box in Fig. 7(b) visualized in 2D. In step 3, we can successfully find all tumor regions in abnormal cases, and discard non-tumor regions in normal cases. As shown in Fig. 7(c), our algorithm can find the right tumor, which overlaps with the tumor annotation in Fig. 7(d).

It should be emphasized that dilated duct helps us to **narrow the searching space** of the tumor, so that we are able to focus on a finer region. Though we train a same network used in [51], half of the missing tumors in [51] can be found. In this way, we are imitating how radiologists detect

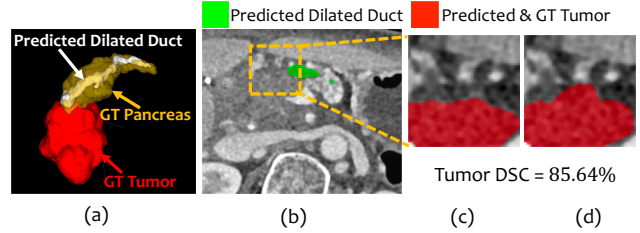


Figure 7. Examples of finding missed tumor of [51] by dilated duct. (a) The ground-truth tumor is right behind one end of the predicted dilated duct. The ground-truth pancreas is shown as a reference. (b) A cropped CT slice with predicted duct (we choose green for better visualization). The yellow dashed box is a candidate tumor region, shown in 2D. (c) and (d) are the same zoomed in image region with predicted and ground-truth tumor, respectively.

Table 5. Normal vs. abnormal classification results. Zhu *et al.* [51] + ours denotes applying our method to find the missing tumor of Zhu *et al.*. “ \uparrow ” and “ \downarrow ” indicate the larger and the smaller the better, respectively.

Methods	Misses \downarrow	Sensitivity \uparrow	Specificity \uparrow
Zhu <i>et al.</i> [51]	8/136	94.1%	98.5%
Zhu <i>et al.</i> [51] + Ours	4/136	97.1%	98.5%

PDAC, *i.e.*, they can find visible PDAC tumors easily, but for difficult ones, they will seek help from dilated ducts.

5. Conclusions

In this paper, we present Deep Distance Transform (DDT) for accurate tubular structure segmentation, which combines intuitions from the classical distance transform for skeletonization and modern deep segmentation networks. DDT guides segmentation by taking the geometric property of the tubular structure into account, which not only leads to a better segmentation result, but also provides the cross-sectional scales, *i.e.*, a geometric measure for the thickness of tubular structures. We evaluated our approach on six datasets including four tubular structures. Experiment shows the superiority of the proposed DDT for tubular structure segmentation and clinical application.

Acknowledgements. This work was supported by the Lustgarten Foundation for Pancreatic Cancer Research and also supported by NSFC No. 61672336.

References

- [1] Luis Álvarez, Agustín Trujillo-Pino, Carmelo Cuenca, Esther González, Julio Esclarín, Luis Gómez, Luis Mazorra, Miguel Alemán-Flores, Pablo G. Tahoces, and José M. Carreira-Villamor. Tracking the aortic lumen geometry by optimizing the 3d orientation of its cross-sections. In *Proc. MICCAI*, 2017.
- [2] Roberto Annunziata, Ahmad Kheirkhah, Pedram Hamrah, and Emanuele Trucco. Scale and curvature invariant ridge detector for tortuous and fragmented structures. In *Proc. MICCAI*, pages 588–595, 2015.
- [3] Luca Antiga, Bogdan Ene-Iordache, and Andrea Remuzzi. Computational geometry for patient-specific reconstruction and meshing of blood vessels from angiography. *IEEE Trans. Med. Imaging*, 22(5):674–684, 2003.
- [4] Carlo Arcelli and Gabriella Sanniti di Baja. Ridge points in euclidean distance maps. *Pattern Recognition Letters*, 13(4):237–243, 1992.
- [5] Oscar Kin-Chung Au, Chiew-Lan Tai, Hung-Kuo Chu, Daniel Cohen-Or, and Tong-Yee Lee. Skeleton extraction by mesh contraction. *ACM Trans. Graph.*, 27(3):44:1–44:10, 2008.
- [6] Andreas Baak, Meinard Müller, Gaurav Bharaj, Hans-Peter Seidel, and Christian Theobalt. A data-driven approach for real-time full body pose reconstruction from a depth camera. In *Proc. ICCV*, 2011.
- [7] Christian Bauer and Horst Bischof. Extracting curve skeletons from gray value images for virtual endoscopy. In *International Workshop on Medical Imaging and Virtual Reality*, pages 393–402, 2008.
- [8] Christian Bauer, Thomas Pock, Erich Sorantin, Horst Bischof, and Reinhard Beichel. Segmentation of interwoven 3d tubular tree structures utilizing shape priors and graph cuts. *Medical Image Analysis*, 14(2):172–184, 2010.
- [9] Fethallah Benmansour and Laurent D. Cohen. Tubular structure segmentation based on minimal path method and anisotropic enhancement. *International Journal of Computer Vision*, 92(2):192–210, 2011.
- [10] Vicent Caselles, Ron Kimmel, and Guillermo Sapiro. Geodesic active contours. *International Journal of Computer Vision*, 22(1):61–79, 1997.
- [11] Linda C. Chu., Michael G. Goggins., and Elliot K. Fishman. Diagnosis and detection of pancreatic cancer. *The Cancer Journal*, 23(6):333–342, 2017.
- [12] Özgün Çiçek, Ahmed Abdulkadir, Soeren S. Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: Learning dense volumetric segmentation from sparse annotation. In *Proc. MICCAI*, 2016.
- [13] Qi Dou, Cheng Ouyang, Cheng Chen, Hao Chen, and Pheng-Ann Heng. Unsupervised cross-modality domain adaptation of convnets for biomedical image segmentations with adversarial loss. In *Proc. IJCAI*, 2018.
- [14] Mark D Edge, Maarouf Hoteit, Amil P Patel, Xiaoping Wang, Deborah A Baumgarten, and Qiang Cai. Clinical significance of main pancreatic duct dilation on computed tomography: Single and double duct dilation. *World J Gastroenterol.*, 13(11):1701–1705, 2007.
- [15] Maayan Frid-Adar, Rula Amer, and Hayit Greenspan. Endotracheal tube detection and segmentation in chest radiographs using synthetic data. In *Proc. MICCAI*, pages 784–792, 2019.
- [16] Huazhu Fu, Yanwu Xu, Stephen Lin, Damon Wing Kee Wong, and Jiang Liu. Deepvessel: Retinal vessel segmentation via deep learning and conditional random field. In *Proc. MICCAI*, pages 132–139, 2016.
- [17] Yaorong Ge and J. Michael Fitzpatrick. On the generation of skeletons from discrete euclidean distance maps. *IEEE Trans. Pattern Anal. Mach. Intell.*, 18(11):1055–1066, 1996.
- [18] Florent Grélard, Fabien Baldacci, Anne Vialard, and Jean-Philippe Domenger. New methods for the geometrical analysis of tubular organs. *Medical Image Analysis*, 42:89–101, 2017.
- [19] Fabian Isensee, Jens Petersen, Simon A. A. Kohl, Paul F. Jäger, and Klaus H. Maier-Hein. nnu-net: Breaking the spell on successful medical image segmentation. *CoRR*, abs/1904.08128, 2019.
- [20] Wei Ke, Jie Chen, Jianbin Jiao, Guoying Zhao, and Qixiang Ye. SRN: side-output residual network for object symmetry detection in the wild. In *Proc. CVPR*, 2017.
- [21] Tom Sie Ho Lee, Sanja Fidler, and Sven J. Dickinson. Detecting curved symmetric parts using a deformable disc model. In *Proc. ICCV*, 2013.
- [22] Alex Levinshtein, Cristian Sminchisescu, and Sven J. Dickinson. Multiscale symmetric part detection and grouping. *International Journal of Computer Vision*, 104(2):117–134, 2013.
- [23] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully Convolutional Networks for Semantic Segmentation. *Proc. CVPR*, 2015.
- [24] J. Merkow, D. Kriegman, A. Marsden, and Z Tu. Dense volume-to-volume vascular boundary detection. In *MICCAI*, 2016.
- [25] Delphine Nain, Anthony J. Yezzi, and Greg Turk. Vessel segmentation using a shape driven flow. In *Proc. MICCAI*, 2004.
- [26] Mathias Perslev, Erik Bjørnager Dam, Akshay Pai, and Christian Igel. One network to segment them all: A general, lightweight system for accurate 3d medical image segmentation. In *MICCAI*, 2019.
- [27] Wayne Rosamond, Katherine Flegal, Karen Furie, Alan Go, Kurt Greenlund, Nancy Haase, Susan M. Hailpern, Michael Ho, Virginia Howard, Brett Kissela, Steven Kittner, Donald Lloyd-Jones, Mary McDermott, James Meigs, Claudia Moy, Graham Nichol, Christopher O’Donnell, Veronique Roger, Paul Sorlie, Julia Steinberger, Thomas Thom, Matt Wilson, and Yuling Hong. Heart disease and stroke statistics-2008 update. *Circulation*, 117(4), 2008.
- [28] Azriel Rosenfeld and John L. Pfaltz. Distance functions on digital pictures. *Pattern Recognition*, 1:33–61, 1968.
- [29] Holger R. Roth, Hirohisa Oda, Yuichiro Hayashi, Masahiro Oda, Natsuki Shimizu, Michitaka Fujiwara, Kazunari Misawa, and Kensaku Mori. Hierarchical 3d fully convolutional networks for multi-organ segmentation. *CoRR*, abs/1704.06382, 2017.

- [30] Rasmus Rothe, Radu Timofte, and Luc Van Gool. Deep expectation of real and apparent age from a single image without facial landmarks. *International Journal of Computer Vision*, 126(2-4):144–157, 2018.
- [31] Wei Shen, Xiang Bai, Zihao Hu, and Zhijiang Zhang. Multiple instance subspace learning via partial random projection tree for local reflection symmetry in natural images. *Pattern Recognition*, 52:306–316, 2016.
- [32] Wei Shen, Kai Zhao, Yuan Jiang, Yan Wang, Xiang Bai, and Alan L. Yuille. Deepskeleton: Learning multi-task scale-associated deep side outputs for object skeleton extraction in natural images. *IEEE Trans. Image Processing*, 26(11):5298–5311, 2017.
- [33] Amber L. Simpson, Michela Antonelli, Spyridon Bakas, Michel Bilello, Keyvan Farahani, Bram van Ginneken, Annette Kopp-Schneider, Bennett A. Landman, Geert J. S. Litjens, Bjoern H. Menze, Olaf Ronneberger, Ronald M. Summers, Patrick Bilic, Patrick Ferdinand Christ, Richard K. G. Do, Marc Gollub, Jennifer Golia-Pernicka, Stephan Heckers, William R. Jarnagin, Maureen McHugo, Sandy Napel, Eugene Vorontsov, Lena Maier-Hein, and M. Jorge Cardoso. A large annotated medical image dataset for the development and evaluation of segmentation algorithms. *CoRR*, abs/1902.09063, 2019.
- [34] Amos Sironi, Vincent Lepetit, and Pascal Fua. Multiscale centerline detection by learning a scale-space distance transform. In *CVPR*, 2014.
- [35] João V. B. Soares, Jorge J. G. Leandro, Roberto M. Cesar, Herbert F. Jelinek, and Michael J. Cree. Retinal vessel segmentation using the 2-d gabor wavelet and supervised classification. *IEEE Trans. Med. Imaging*, 25(9):1214–1222, 2006.
- [36] Joes Staal, Michael D. Abramoff, Meindert Niemeijer, Max A. Viergever, and Bram van Ginneken. Ridge-based vessel segmentation in color images of the retina. *IEEE Trans. Med. Imaging*, 23(4):501–509, 2004.
- [37] Ahmed Taha, Pechin Lo, Junning Li, and Tao Zhao. Kid-net: Convolution networks for kidney vessels segmentation from ct-volumes. In *Proc. MICCAI*, 2018.
- [38] Stavros Tsogkas and Iasonas Kokkinos. Learning-based symmetry detection in natural images. In *Proc. ECCV*, pages 41–54, 2012.
- [39] Chenglong Wang, Yuichiro Hayashi, Masahiro Oda, Hayato Itoh, Takayuki Kitasaka, Alejandro F. Frangi, and Kensaku Mori. Tubular structure segmentation using spatial fully connected network with radial distance loss for 3d medical images. In *Proc. MICCAI*, 2019.
- [40] Yukang Wang, Yongchao Xu, Stavros Tsogkas, Xiang Bai, Sven J. Dickinson, and Kaleem Siddiqi. Deepflux for skeletons in the wild. In *Proc. CVPR*, 2019.
- [41] Yan Wang, Yuyin Zhou, Wei Shen, Seyoun Park, Elliot K. Fishman, and Alan L. Yuille. Abdominal multi-organ segmentation with organ-attention networks and statistical fusion. *Medical Image Analysis*, 55:88–102, 2019.
- [42] Yan Wang, Yuyin Zhou, Peng Tang, Wei Shen, Elliot K. Fishman, and Alan L. Yuille. Training multi-organ segmentation networks with sample selection by relaxed upper confident bound. In *Proc. MICCAI*, 2018.
- [43] Onno Wink, Wiro J. Niessen, and Max A. Viergever. Fast delineation and visualization of vessels in 3d angiographic images. *IEEE Trans. Med. Imaging*, 19(4):337–346, 2000.
- [44] Yingda Xia, Fengze Liu, Dong Yang, Jinzheng Cai, Lequan Yu, Zhuotun Zhu, Daguang Xu, Alan L. Yuille, and Holger Roth. 3d semi-supervised learning with uncertainty-aware multi-view co-training. In *WACV*, 2019.
- [45] Peter J. Yim, Juan R. Cebal, Rakesh Mullick, and Peter L. Choyke. Vessel surface reconstruction with a tubular deformable model. *IEEE Trans. Med. Imaging*, 20(12):1411–1421, 2001.
- [46] Kai Zhao, Wei Shen, Shanghua Gao, Dandan Li, and Ming-Ming Cheng. Hi-fi: Hierarchical feature integration for skeleton detection. In *Proc. IJCAI*, 2018.
- [47] Yuyin Zhou, Yingwei Li, Zhishuai Zhang, Yan Wang, Angtian Wang, Elliot K. Fishman, Alan L. Yuille, and Seyoun Park. Hyper-pairing network for multi-phase pancreatic ductal adenocarcinoma segmentation. In *Proc. MICCAI*, 2019.
- [48] Yuyin Zhou, Zhe Li, Song Bai, Xinlei Chen, Mei Han, Chong Wang, Elliot K. Fishman, and Alan L. Yuille. Prior-aware neural network for partially-supervised multi-organ segmentation. In *Proc. ICCV*, 2019.
- [49] Yuyin Zhou, Lingxi Xie, Wei Shen, Yan Wang, Elliot K. Fishman, and Alan L. Yuille. A fixed-point model for pancreas segmentation in abdominal CT scans. In *Proc. MICCAI*, 2017.
- [50] Zhuotun Zhu, Yingda Xia, Wei Shen, Elliot K. Fishman, and Alan L. Yuille. A 3d coarse-to-fine framework for volumetric medical image segmentation. In *Proc. 3DV*, 2018.
- [51] Zhuotun Zhu, Yingda Xia, Lingxi Xie, Elliot K. Fishman, and Alan L. Yuille. Multi-scale coarse-to-fine segmentation for screening pancreatic ductal adenocarcinoma. In *Proc. MICCAI*, 2019.