

# HarmonySeg: Tubular Structure Segmentation with Deep-Shallow Feature Fusion and Growth-Suppression Balanced Loss

Yi Huang<sup>\*1,2,3</sup>, Ke Zhang<sup>\*4</sup>, Wei Liu<sup>1,2</sup>, Yuanyuan Wang<sup>3</sup>, Vishal M. Patel<sup>4</sup>, Le Lu<sup>1</sup>, Xu Han<sup>5</sup>, Dakai Jin<sup>1</sup>, and Ke Yan<sup>†1,2</sup>

<sup>1</sup>DAMO Academy, Alibaba Group

<sup>2</sup>Hupan Lab, Hangzhou, China

<sup>3</sup>Department of Biomedical Engineering, Fudan University

<sup>4</sup>Department of Electrical and Computer Engineering, Johns Hopkins University

<sup>5</sup>Department of Hepatobiliary and Pancreatic Surgery, The First Affiliated Hospital of College of Medicine, Zhejiang University

## Abstract

good generalizability.

## 1. Introduction

Accurate segmentation of tubular structures in medical images, such as vessels and airway trees, is crucial for computer-aided diagnosis, radiotherapy, and surgical planning. However, significant challenges exist in algorithm design when faced with diverse sizes, complex topologies, and (often) incomplete data annotation of these structures. We address these difficulties by proposing a new tubular structure segmentation framework named HarmonySeg. First, we design a deep-to-shallow decoder network featuring flexible convolution blocks with varying receptive fields, which enables the model to effectively adapt to tubular structures of different scales. Second, to highlight potential anatomical regions and improve the recall of small tubular structures, we incorporate vesselness maps as auxiliary information. These maps are aligned with image features through a shallow-and-deep fusion module, which simultaneously eliminates unreasonable candidates to maintain high precision. Finally, we introduce a topology-preserving loss function that leverages contextual and shape priors to balance the growth and suppression of tubular structures, which also allows the model to handle low-quality and incomplete annotations. Extensive quantitative experiments are conducted on four public datasets. The results show that our model can accurately segment 2D and 3D tubular structures and outperform existing state-of-the-art methods. External validation on a private dataset also demonstrates

---

<sup>\*</sup>Both authors contributed equally to this work.

<sup>†</sup>Corresponding author, yanke.yan@alibaba-inc.com

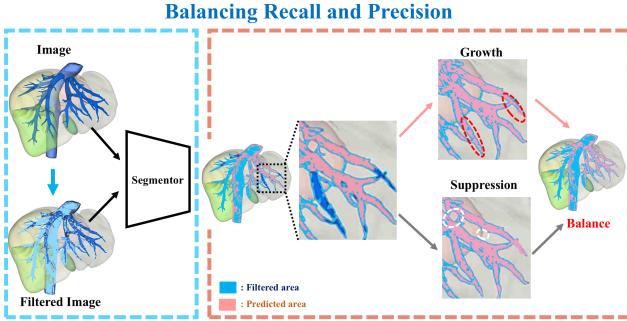


Figure 1. Motivation of the proposed HarmonySeg. We leverage vesselness to improve recall and filter out outliers based on image features to maintain precision, while applying growth-suppression balanced loss to encourage the growth of vessels at a reasonable noise level in the absence of labels.

or even missing. Some works have noticed this problem [1, 15, 31, 36], and they typically focused on maintaining the continuity of structures via skeleton growth. However, they did not adequately balance the growth of skeleton and the suppression of noise resulting from overgrowth. Secondly, tubular structures vary markedly in sizes and shapes. Pioneering studies have demonstrated that models utilizing multi-scale feature extraction are well-suited for this situation [12, 46]. Some studies introduced specialized convolution blocks for snake-shaped tubular structures, achieving good performance but at high computational cost [26]. Others combine traditional techniques, like vesselness filters, with deep learning to improve tubular structure segmentation [8, 10, 11, 38]. However, these models only concatenate the image and filtering results, a shallow integration that overlooks vesselness’s potential to enhance deeper feature extraction for small structures.

To accurately segment tubular structures in medical images, we propose a framework named HarmonySeg. As shown in Figure 1, HarmonySeg leverages vesselness filters to enhance recall in two ways. First, we introduce a deep mutual query (DMQ) module that uses cross-attention between the image and vesselness results to boost deep features, especially for small-scale structures. Second, a deep-to-shallow decoding (D2SD) strategy progressively refines segmentation, preserving multi-scale structures. We replace standard convolutions with flexible blocks featuring diverse receptive fields to capture structures of varying sizes. To address partial annotations, we design loss functions to balance structure growth with noise suppression, reducing false positives and compensating for missing labels. Together, these techniques enable HarmonySeg to effectively capture topology and continuity, demonstrating improved recall and precision. Our contributions are fourfold:

- 1) We introduce a shallow and deep fusion (SADF) module designed to fully harness the potential of vessel-

ness maps for improving recall while simultaneously ensuring precision by filtering out unwarranted vessel candidates based on image features.

- 2) A deep-to-shallow decoding (D2SD) strategy is designed to progressively refine the segmentation results with the enhanced features of SADF, which further align and aggregate the features of vesselness and image at different scales, providing varying sensitivities to target sizes and effectively preserve structures.

- 3) We design loss functions that effectively balance tubular structure growth and noise suppression (GSB). These loss functions compensate for missed labels, enhance recall and reduce false positives.

- 4) Extensive experiments carried out on four public datasets validated the performance of HarmonySeg, which can accurately segment 2D and 3D tubular structures, outperforming existing state-of-the-art methods. An external validation on a private dataset also demonstrates its good generalizability.

## 2. Related Work

In this section, we review the existing approaches that involve tubular structure segmentation in medical images.

**Vesselness Filtering:** Vesselness filters can increase the vessels’ contrast and suppress the signal of non-vessel structures. They are often used as a preprocessing step for tubular structure segmentation [17, 18]. [38] and [10] concatenated images and vesselness filtering results as the model input to highlight the potential vessel regions. Some works also utilized similar strategies to roughly localize potential tubular structures [8, 11]. However, most of these methods simply fused the image and filtered results by concatenating them as input. By contrast, we consider the vesselness filtering result as an independent auxiliary modality, encode it in parallel with the image interactively, and use a novel mutual fusion module to highlight the tubular structure features at different scales.

**Feature Extraction and Fusion:** Flexible convolution with various receptive fields facilitates feature extraction for tubular structures with varying morphology and size [6, 12, 14, 46]. Most of them introduced deformable convolution with flexible receptive fields and enhanced multi-scale feature fusion modules. For example, Qi et al. [26] designed a new dynamic snake convolution to adaptively focus on the slender and tortuous local features of tubular structures. In our framework, we also adopt flexible convolution blocks with diversified receptive fields, incorporated with a multi-scale D2FD strategy, to improve the model’s adaptability for tubular structures of various sizes.

**Topology Exploration and Preservation:** Preserving the complex topology is critical for the segmentation of tubular structures. Improving the model architecture is an effective way to achieve this [19, 42, 43]. Additionally, loss func-

tion constraint is also a useful approach to ensure the connectivity of topology. A new centerline Dice (cIDice) loss was introduced in [31] to measure the similarity between the skeleton of prediction and label to guarantee topology connectivity and consistency. Kirchhoff et al. [15] further proposed skeleton recall loss based on the skeleton, which is computationally efficient and suitable for multi-class tubular structure segmentation tasks. However, these methods are applicable only when the labels are complete. In this paper, we introduce a novel loss function to preserve the topology, ensuring the rationality of skeleton growth by mitigating the noise caused by overgrowth, especially in the context of partial labeling.

### 3. Method

Our approach is designed for tubular structure segmentation in various medical images. In this section, we take hepatic vessel segmentation in computed tomography (CT) as an example task. The overall architecture is shown in Fig. 2.

#### 3.1. Shallow and Deep Fusion (SADF)

The vesselness filter is an effective image preprocessing technique that enhances vessel regions while distinguishes vessels from non-vessel structures<sup>1</sup>. We treat the vesselness map as an auxiliary modality and introduce the SADF module to effectively fuse information at both shallow and deep stages. This module comprises two key components: Deep Mutual Query (DMQ) and Shallow Query (SQ).

**Deep Mutual Query (DMQ)** is conducted on deep features of both CT and vesselness map, i.e.  $F_{C_4}$  and  $F_{V_4}$ , which fulfills the integration of global dependencies between the CT and the vesselness and as a basis for decoding. As shown in of Figure 2(b), it can be described as:

$$DQ_{V2C} = \text{Cross}(Q_V, K_C, V_C) + \text{Self}(F_{C_4}), \quad (1)$$

where  $K_C$  and  $V_C$  are the projected key and value maps from the deep feature  $F_{C_4}$  of CT image, and the query maps  $Q_V$  is projected from the deep features  $F_{V_4}$  of vesselness map. Next, all of them are injected into the cross-attention mechanism (*Cross*) and further combined with the self-attention (*Self*) results of  $F_{C_4}$  to obtain the enhanced tubular structures' features  $DQ_{V2C}$ . In this way, the vesselness highlights the densely vascularized regions of the CT globally, and the global dependence of the CT itself is also protected from being severely affected by noises in the vesselness. Similar processing is done between  $F_{V_4}$  and  $F_{C_4}$  as well to obtain  $DQ_{C2V}$ , which aims to mitigate the negative effects of outliers on vesselness by obvious vessels in CT, such as those located at the liver border.

**Shallow Query (SQ)** is defined as follows:

$$SQ_i = \text{Cat}(F_{S'_1}, F_{S_2}), \quad (2)$$

<sup>1</sup>More details are present in Appendix A.

$$F_{S'_1} = \text{Self}(\text{AvgP}(Q_{S_1}), \text{MaxP}(K_{S_1}), V_{S_1}), \quad (3)$$

in which *Cat*, *AvgP*, and *MaxP* refer to concatenation, average, and max pooling, respectively. SQ is implemented on shallow features of CT and vesselness ( $F_{C_i}$  and  $F_{V_i}$ ,  $i = 1, 2, 3$ ). Shallow feature maps contain more accurate spatial information than deep ones, so we utilize vesselness to help locate the potential vessel in the CT image. First, we fuse the shallow features of CT and vesselness, and then equally split them into  $F_{S_1}$  and  $F_{S_2}$  on the channel dimension, to reduce the computational cost of the self-attention mechanism. Then,  $F_{S_1}$  is fed into *Self*, in which pooling operations are used to enhance sensitivity to tiny targets inspired by [24]. The optimized  $F_{S'_1}$  indicates vessel candidates at a global scale and is further concatenated with  $F_{S_2}$  to recover the original shape. Finally,  $SQ_i$  is added with input features improved by the up-sampled  $DQ_{V2C}$  and  $DQ_{C2V}$  to integrate all vessel information and feed into the parallel multiple decoding stage, see Figure 2(c).

#### 3.2. Deep-to-Shallow Decoding (D2SD)

Standard convolution and downsampling operations enrich semantic information while diluting local detail information. In our study, the size of the vessel is varied and the topology of vessels is complex, so consistently preserving the spatial information becomes more critical. Therefore, we design the D2SD strategy. Unlike the common U-Net architecture, HarmonySeg performs decoding progressively at multiple scales from deep to shallow after the complete fusion between CT and vesselness enhanced features achieved by the SADF. Pre-decoding is independently carried out at each shallow scale, taking advantage of the local invariance and detailed spatial information present in the shallow scale features. This allows for the effective localization of vessels that are observable at this scale. Moreover,  $DQ_{V2C}$  and  $DQ_{C2V}$  are also involved in the pre-decoding after up-sampling, so the fused global dependencies from CT and vesselness further assist in the alignment of the two modality features at the shallow-scale decoding process, which leverages the vesselness to highlight potential vessel regions once again in decoding. Finally, the pre-decoded results  $Seg_i$  ( $i = 1, 2, 3$ ) and deep one  $Seg_4$ , with varying sensitivities to vessel sizes at different scales, are further fused to get better predictions through a convolution block<sup>2</sup>.

Besides, flexible convolution blocks, used in pre-decoding and encoding phases, also benefit the model's ability to extract and aggregate multi-scale features for vessels. We parallelize and stack convolution blocks to provide diverse receptive fields, followed by an  $1 \times 1 \times 1$  convolution to adjust the channel number<sup>2</sup>. Compared to the

<sup>2</sup>More details are present in Appendix B and C.

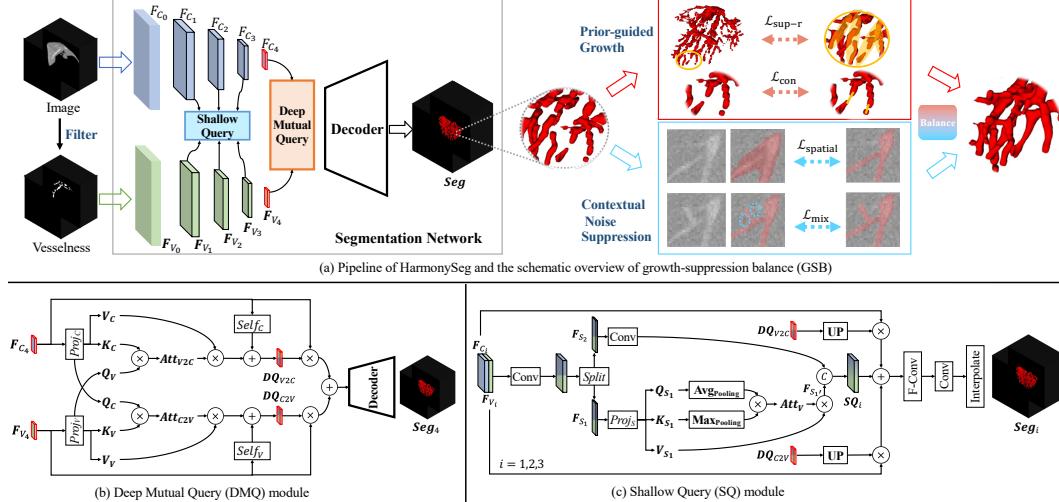


Figure 2. The framework of HarmonySeg. Our model takes an image and its filtered results as input, performing multi-level feature fusion to enhance recall. Growth-suppression loss functions are then applied to improve segmentation precision.

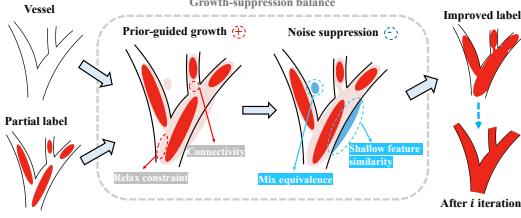


Figure 3. Illustration of growth-suppression balance (GSB).

direct use of dilated convolution, the convolution stacking approach does not generate a serious grid effect, so it can better preserve local details and be more suitable for vessel segmentation with complex local details.

### 3.3. Growth-Suppression Balance (GSB)

In this section, we present the growth-suppression balance strategies for vessel segmentation. As visualized in Figure 3, the process is divided into two stages. In the first stage, we utilize local context and shape priors to expand the vessel area, referred to as prior-aware vessel growth. In the second stage, we enforce the mix-equivalence of the predicted segmentation, which helps suppress contextual noise and ensures more refined results.

**Prior-guided vessel growth:** We leverage shape priors to drive the expansion of vessel segmentation. First, due to incomplete annotations, we relax the constraints of the segmentation loss, allowing the segmentation to extend beyond the annotated regions. We introduce a region of interest (ROI) term,  $\mathbf{R}$ , to define a coarse area encompassing all vessels. This region is represented by a bounding box enclosing all annotated vessels. We assume that pixels outside

the box belong to the background, which are treated as negative samples. The pixels that belong to annotation  $y$  are positive labels. While the remaining pixels represent uncertain vessel regions. The uncertain-aware prediction  $\hat{y}'$  is represented as:

$$\hat{y}' = \underbrace{y\hat{y}}_{\text{positive}} + \beta \underbrace{(y^c\mathbf{R}\hat{y})}_{\text{uncertain}} + \underbrace{y^c\mathbf{R}^c\hat{y}}_{\text{negative}}, \quad (4)$$

where  $\hat{y}$  is model prediction; the uncertainty ratio  $\alpha$  is defined as  $\beta = \frac{1}{\log(\sum y^c / \sum y)}$ ;  $y^c$  and  $\mathbf{R}^c$  represent the complement sets of  $y$  and  $\mathbf{R}$ , respectively. Then, we derive the relaxed supervision loss  $\mathcal{L}_{\text{r-sup}}$  as:

$$\mathcal{L}_{\text{r-sup}} = \mathcal{L}'_{\text{Dice}} + \mathcal{L}_{\text{ce}} = -\frac{y\hat{y}}{y + \hat{y}'} - y \log(\hat{y}'). \quad (5)$$

We further reconstruct missing vessel segments, ensuring the completed regions form a continuous, interconnected structure while maintaining density consistency with the surrounding vessels. First, we apply the soft-skeletonization method from [31], which uses iterative min-and max-pooling operations as proxies for morphological erosion and dilation to capture the vessel's skeleton results  $\hat{y}^s = f_{\text{skeleton}}(\hat{y})$ . Following [7], We identify the largest connected component (CC) from all smaller ones. For each smaller component, we identify and extract all endpoints. Then, we connect the endpoint of each CC to the nearest endpoint by drawing a line between them. The resulting paths are represented as  $\hat{y}^c = f_{\text{connect}}(\hat{y}^s)$ . We then treat these reconnected pixels as pseudo-labels and define a connectivity loss,  $\mathcal{L}_{\text{con}}$ , to encourage the skeleton outputs to align with their reconnected versions.

$$\mathcal{L}_{\text{con}} = -\hat{y}^c \log(\hat{y}^s) = -f_{\text{connect}}(\hat{y}^s) \log(\hat{y}^s), \quad (6)$$

where  $\mathcal{L}_{\text{con}}$  calculate the cross entropy between the predicted skeleton  $\hat{y}^s$  and reconnected skeleton  $\hat{y}^c$ .

**Contextual noise suppression** We utilize spatial priors and mix-equivalence as guidance to suppress noise in the segmentation results. Firstly, shallow feature similarities are employed to control the growth, ensuring that the expanded regions remain consistent with the original vessel areas regarding density distribution and spatial alignment. Inspired by [25], we further extend the spatial regularization from 2D images to 3D volumes. We leverage Gaussian kernel  $k_{ij}$  to measure the shallow feature similarity between pixels at location  $i$  and  $j$ , i.e.,  $k_{ij} = \exp\left(-\frac{(l_i - l_j)^2}{2\sigma_l^2} + \frac{(c_i - c_j)^2}{2\sigma_c^2}\right)$ .  $l$  and  $c$  denote the location and color feature specific to position  $i$  and  $j$ , then we derive the loss of spatial prior with a gated function define the local neighborhood window  $\Omega_r$  with radius  $r$  for each coordinate:

$$\mathcal{L}_{\text{spatial}} = \sum_{i,j \in \Omega_r} k_{ij} \hat{y}_i \hat{y}_j \quad (7)$$

We further generate mixed samples using the MixUp technique as auxiliary inputs to enforce mix-equivalence. Given two input images,  $x_1$  and  $x_2$ , the mixed image is defined as  $x' = \alpha x_1 + (1 - \alpha)x_2$ , where  $\alpha$  is sampled from uniform distribution. The predicted outputs for these mixed inputs are then constrained by the following loss function:

$$\mathcal{L}_{\text{mix}} = -\frac{\hat{y}' \cdot [\alpha y_1 + (1 - \alpha)y_2]}{\|\hat{y}'\| \cdot \|\alpha y_1 + (1 - \alpha)y_2\|} \quad (8)$$

Finally, we derive the optimization objective ( $\mathcal{L}$ ) that balances the opposing forces of growth and suppression,

$$\mathcal{L} = \underbrace{\mathcal{L}_{\text{r-sup}} + \mathcal{L}_{\text{con}}}_{\text{grow}} + \underbrace{\lambda(\mathcal{L}_{\text{spatial}} + \mathcal{L}_{\text{mix}})}_{\text{suppress}}, \quad (9)$$

The first two terms promote vessel growth, while the last two focus on noise suppression, and  $\lambda$  controls the weight for noise suppression. The proposed loss  $\mathcal{L}$  encourages precise boundary delineation while ensuring effective noise suppression in a well-balanced manner.

## 4. Experiments

### 4.1. Dataset

**Hepatic vessel segmentation (HVS).** We curated a hybrid hepatic vessel segmentation dataset using contrast enhanced CT scans from 992 patients to evaluate our method. This dataset combines three public datasets: LiVS [8] (532 volumes), MSD8 [32] (440 volumes), and 3DIRCADb [33] (20 volumes). LiVS has incomplete labels, and MSD8's test set labels are not publicly available. These label gaps hinder topology capture and loss computation for correctly segmented regions without labels. To fix this, we first trained

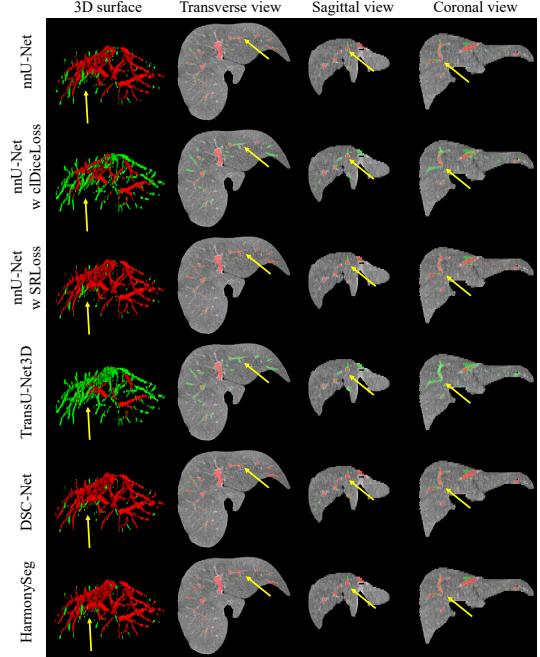


Figure 4. Visualization of hepatic vessel segmentation results, with red indicating the segmentation and green representing the corresponding labels. Yellow arrows highlight the improvements.

a model using MSD8's training set since it has relatively complete labels, and then used it to refine labels on LiVS and MSD8's test set, followed by manual refinement by experts<sup>3</sup>. We applied Hounsfield Unit (HU) value clipping (-100, 200) and liver cropping to the MSD8 and 3DIRCADb volumes, following the preprocessing steps for LiVS. Finally, volumes with refined labels from LiVS and MSD8 were used for training and validation, while 3DIRCADb volumes were used for testing.

**HVS-External.** An external testing set was collected to further assess HarmonySeg's effectiveness, which contains 21 cases from our collaborating hospital with hepatic vessel labels manually delineated by experienced experts and underwent the same preprocessing as HVS. It was tested directly using the models trained on HVS.

**Retinal vessel segmentation (RVS).** The DRIVE dataset [34] consisted of 40 2D digital retinal images. They were equally divided into 20 images for training and 20 for testing. The retinal vessel labels were manually delineated by ophthalmological experts.

**Airway tree segmentation (ATS).** The airway tree segmentation dataset [27, 35] includes 90 chest CT volumes, in which 70 and 20 were collected from the LIDC dataset [2] and the training set of EXACT'09 dataset [20], respectively. These volumes were further split into 50, 20 and 20 volumes for training, validation and testing, respectively.

<sup>3</sup>Details are displayed in the supplementary materials.

Table 1. **Quantitative comparison on HVS task.** The best and second-best results are highlighted in bold and underlined, respectively.

Model	Dataset	Dice(%, $\uparrow$ )	clDice(%, $\uparrow$ )	HD( $\downarrow$ )	ASSD( $\downarrow$ )	F1-score(%, $\uparrow$ )
nnU-Net [13]	HVS	60.15 $\pm$ 9.49	69.41 $\pm$ 5.43	10.00 $\pm$ 4.22	2.23 $\pm$ 0.85	62.04 $\pm$ 12.15
+ clDice $\mathcal{L}$ [31]		40.82 $\pm$ 11.19	41.98 $\pm$ 10.25	11.78 $\pm$ 4.27	6.81 $\pm$ 2.83	42.22 $\pm$ 10.56
+ SRL [15]		61.24 $\pm$ 8.85	<u>71.80</u> $\pm$ 4.54	9.94 $\pm$ 4.18	<b>1.80</b> $\pm$ 0.59	63.28 $\pm$ 12.43
TransU-Net [4]		42.49 $\pm$ 17.84	50.11 $\pm$ 18.54	12.96 $\pm$ 7.14	4.33 $\pm$ 3.69	48.54 $\pm$ 21.66
DSC-Net [26]		63.57 $\pm$ 6.83	70.62 $\pm$ 5.57	9.65 $\pm$ 3.86	2.02 $\pm$ 0.77	65.51 $\pm$ 13.13
HarmonySeg		<b>66.79</b> $\pm$ 6.34	<b>72.04</b> $\pm$ 5.22	<b>9.60</b> $\pm$ 3.98	1.96 $\pm$ 0.73	<b>67.17</b> $\pm$ 14.39
nnU-Net [13]	HVS-External	63.78 $\pm$ 9.27	69.30 $\pm$ 9.52	3.69 $\pm$ 1.17	2.33 $\pm$ 1.69	64.40 $\pm$ 0.75
+ clDice $\mathcal{L}$ [31]		57.31 $\pm$ 7.23	53.52 $\pm$ 8.20	4.27 $\pm$ 1.13	2.76 $\pm$ 0.99	57.85 $\pm$ 5.98
+ SRL [15]		71.64 $\pm$ 9.33	78.62 $\pm$ 9.72	3.19 $\pm$ 1.12	1.32 $\pm$ 1.04	72.47 $\pm$ 3.13
TransU-Net [4]		42.33 $\pm$ 22.25	43.63 $\pm$ 21.36	8.03 $\pm$ 7.54	6.55 $\pm$ 9.33	56.10 $\pm$ 30.53
DSC-Net [26]		75.83 $\pm$ 6.19	77.54 $\pm$ 7.43	3.15 $\pm$ 1.16	1.02 $\pm$ 0.50	76.17 $\pm$ 5.52
HarmonySeg		<b>77.76</b> $\pm$ 6.12	<b>80.61</b> $\pm$ 8.07	3.55 $\pm$ 1.19	1.40 $\pm$ 1.32	<b>77.15</b> $\pm$ 8.48

Table 2. **Dice(%) for different vessel sizes on the HVS task.** Vessel branches were categorized into small ( $<5\text{mm}$ ), medium ( $5\text{--}10\text{mm}$ ), and large ( $\geq 10\text{mm}$ ) based on diameter.

Model	Small	Medium	Large
nnU-Net [13]	24.04 $\pm$ 14.77	39.02 $\pm$ 18.04	54.11 $\pm$ 15.77
+ clDice $\mathcal{L}$ [31]	8.08 $\pm$ 8.27	27.93 $\pm$ 16.88	29.51 $\pm$ 15.65
+ SRL [15]	25.02 $\pm$ 10.93	40.16 $\pm$ 16.76	52.26 $\pm$ 15.47
TransU-Net [4]	11.07 $\pm$ 9.35	25.45 $\pm$ 16.21	35.15 $\pm$ 20.93
DSC-Net [26]	25.66 $\pm$ 8.93	40.85 $\pm$ 8.99	57.64 $\pm$ 11.30
HarmonySeg	<b>27.93</b> $\pm$ 13.59	<b>40.97</b> $\pm$ 7.47	<b>58.46</b> $\pm$ 12.87

Table 3. **Quantitative comparison on the RVS task.**

Model	Dice(%, $\uparrow$ )	clDice(%, $\uparrow$ )	HD( $\downarrow$ )
U-Net [29]	80.73 $\pm$ 1.77	79.66 $\pm$ 4.00	6.86 $\pm$ 0.56
TransU-Net [3]	80.56 $\pm$ 2.14	79.02 $\pm$ 5.05	6.83 $\pm$ 0.52
CS <sup>2</sup> -Net [23]	77.53 $\pm$ 2.94	74.88 $\pm$ 5.27	6.90 $\pm$ 0.48
DCU-Net [39]	80.83 $\pm$ 1.99	80.19 $\pm$ 4.80	<u>6.68</u> $\pm$ 0.49
DSC-Net [26]	<b>81.85</b> $\pm$ 1.74	<u>81.16</u> $\pm$ 4.54	<u>6.68</u> $\pm$ 0.49
nnU-Net [13]	80.13 $\pm$ 1.60	78.82 $\pm$ 3.77	7.61 $\pm$ 0.47
HarmonySeg	81.33 $\pm$ 1.61	<b>82.03</b> $\pm$ 3.58	<b>6.51</b> $\pm$ 0.83

Similar to the HVS task, the HU clip (-1350, 150) was performed for the volumes in the ATS task.

**Coronary artery segmentation (CAS).** The coronary artery segmentation dataset was established based on the ImageCAS dataset [40], which captured data from 1000 patients. The dataset was divided into 700, 50, and 250 cases for training, validation, and testing, respectively. An HU clip (-400, 500) was also conducted to optimize the observability of coronary arteries.

#### 4.2. Implementation Details

We implemented HarmonySeg in PyTorch and ran experiments on an NVIDIA A100 GPU with 80 GB of memory. The training was performed using the Adam optimizer with an initial learning rate of 1e-4 and a polynomial decay strat-

Table 4. **Results on ATS task.** Prec denotes Precision.

Model	BD(%, $\uparrow$ )	TLD(%, $\uparrow$ )	Prec(%, $\uparrow$ )
Juarez et al. [9]	69.2 $\pm$ 25.4	53.5 $\pm$ 20.9	<b>99.9</b> $\pm$ 0.1
WingsNet [45]	89.2 $\pm$ 5.8	77.1 $\pm$ 5.7	99.0 $\pm$ 0.8
CFDA [41]	90.9 $\pm$ 6.7	78.9 $\pm$ 8.1	<b>99.1</b> $\pm$ 0.6
Qin et al. [28]	90.9 $\pm$ 8.8	80.7 $\pm$ 9.9	98.4 $\pm$ 1.0
Zheng et al. [44]	<u>91.1</u> $\pm$ 5.5	80.1 $\pm$ 6.6	98.9 $\pm$ 0.7
nnU-Net [13]	90.0 $\pm$ 30.0	<u>91.0</u> $\pm$ 5.0	88.7 $\pm$ 5.7
HarmonySeg	<b>95.0</b> $\pm$ 21.8	<b>92.3</b> $\pm$ 4.5	91.1 $\pm$ 2.9

egy. The preprocessing followed the same scheme as the nnU-Net [13]. For the vesselness filter, we follow parameters from [17]. The suppression weight  $\lambda$  is empirically set to 1, with its impact validated in the ablation study.

**Metrics:** We used a series of quantitative evaluation metrics based on overlap, distance, and connectivity to comprehensively measure the performance of the model. Concretely, the HVS task employed pixel-wise Dice, centerline Dice (clDice), Hausdorff distance (HD), average symmetric surface distance (ASSD) and F1-score as evaluation metrics [8]. Similarly, the RVA task also used Dice, clDice, and HD [26]. The ATS task focused on connectivity, so branch detected (BD), tree length detected (TLD), and precision, were used as evaluation metrics [20]. The CAS task introduced Dice, HD, and average Hausdorff distance (AHD) to assess the model performance [40]. In addition, we thoroughly compared with various benchmarks published on the dataset of RVS, ATS, and CAS.

**Baselines:** We compare HarmonySeg with three state-of-the-art approaches designed for tubular structure segmentation in medical images and two widely-used general medical segmentation models, including nnU-Net [13], nnU-Net with clDice loss (clDice $\mathcal{L}$  [31]), nnU-Net with skeleton recall loss (SRL [15]), TransU-Net [4] and DSC-Net [26]. nnU-Net is a representative baseline in medical image segmentation, and its cooperation with the clDice and the

Table 5. Quantitative comparison on the CAS task.

Model	Dice(%, $\uparrow$ )	HD( $\downarrow$ )	AHD( $\downarrow$ )
Shen et al. [30]	80.58	28.67	0.85
Chen et al. [5]	72.01	40.96	3.07
Kong et al. [16]	68.78	30.34	1.43
Wolterink et al. [37]	70.61	27.87	1.24
U-Net++ [47]	<u>82.96</u>	<u>27.22</u>	<b>0.82</b>
nnU-Net [13]	75.46	34.36	0.91
HarmonySeg	<b>83.24</b>	<b>26.42</b>	0.88

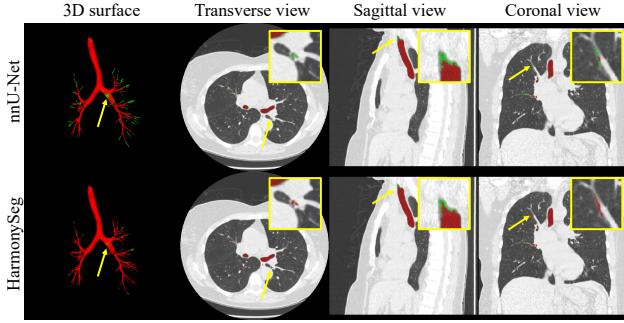


Figure 5. Qualitative comparison of airway tree segmentation.

skeleton recall loss are two effective improved variants for tubular structure segmentation. TransU-Net also performed well in medical image segmentation by integrating the local invariance of convolution and global dependencies of the transformer. DSC-Net utilized a dynamic snake convolution to capture more topological information so that precise vessel segmentation was achieved. The comparison experiment was conducted on all three datasets. It should be noted that our model evaluated on RVS did not involve GSB because the retinal vessel is so well labeled that no further growth is required. Our models evaluated on ATS and CAS did not involve SADF because these two datasets are not suitable for the vessleness filter as adjacent structures can cause serious interference.

#### 4.3. Performance Comparison

**Hepatic vessel segmentation:** Quantitative evaluation metrics are summarized in Table 1. Some visualization examples are given in Figure 4. As indicated in Table 1, HarmonySeg achieves the best performance on both of the two commonly used segmentation evaluation metrics, Dice and HD, with improvements of 5.1% and 0.5% compared to those of the second best model, respectively. Moreover, our model also performs competitively in the comparison of other metrics. Taking F1-score as an example, our model maintains an improved balance between precision and recall, and does not have an obvious drop-off between them like other models. In addition, the Dice comparison in Table 2 further reveals the adaptability of our model for dif-

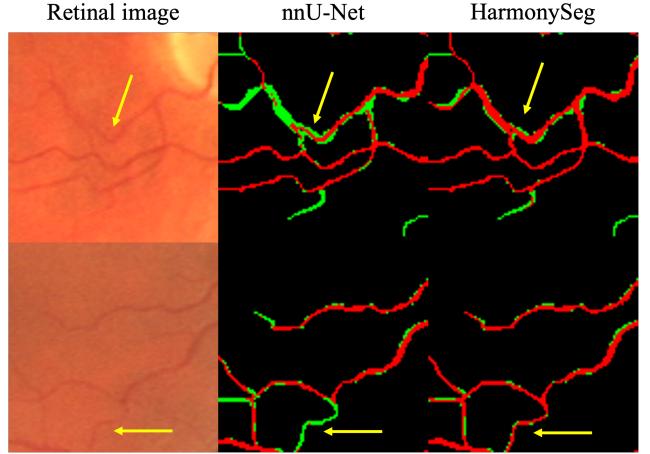


Figure 6. Qualitative comparison on retinal vessel segmentation.

ferent vessel sizes, especially the better accuracy in small vessels. Observing segmentation results in Figure 4, it can be found that our 3D hepatic vessel tree is more complete compared with others. In 2D views, our segmentation is also accurate and has good continuity for both large vessels and tiny branches.

**Results on HVS-External:** The evaluation metrics of HVS-External in Table 1 further reveal HarmonySeg’s generalizability on external data. It achieves the best mean Dice and cIDice, with 2.5% and 4.0% improvements compared with the second. Other metrics are also competitive. This generalizability indicates that our model has potential for application in clinical practice.

**Retinal vessel segmentation:** Table 3 shows the superiority of our framework even without the GSB. The highest cIDice is achieved by HarmonySeg, which has an increment of 1.1% in comparison with the second. This improvement is also visualized in Figure 6. As highlighted in the figure, our model preserves continuity for tiny branches effectively. Further, the HD of HarmonySeg also decreases by 2.5%.

**Airway tree segmentation:** Our model reports the highest BD and TLD with competitive precision in Table 4 for airway tree segmentation. Compared with those of the second-best model, the mean BD and TLD raise 4.3% and 1.4%, respectively. The visualization in Figure 5 also indicates that more airway tree branches are extracted by our model.

**Coronary artery segmentation:** The effectiveness of HarmonySeg in coronary artery segmentation is demonstrated with the best Dice and HD in Table 5. Moreover, it can be observed in the visualization example of Figure 7 that HarmonySeg obtains a more complete coronary artery in 3D views and captures more tiny vessels in 2D views.

Table 6. **Ablation study of HarmonySeg components on the HVS task.** #2 denotes the model using the simple concatenation of CT and vesselness maps as the input, denoted by  $\mathbf{VS}_{\text{Cat}}$ .

Methods	D2SD	SADF	GSB	Dice(%,▶)	clDice(%,▶)	HD( $\downarrow$ )	ASSD( $\downarrow$ )	F1-score(%,▶)
#1	-	-	-	$60.15 \pm 9.49$	$69.41 \pm 5.43$	$10.00 \pm 4.22$	$2.23 \pm 0.85$	$62.04 \pm 12.15$
#2	-	$\mathbf{VS}_{\text{Cat}}$	-	$60.07 \pm 8.78$	$68.63 \pm 5.23$	$10.05 \pm 4.21$	$2.32 \pm 0.83$	$61.82 \pm 11.74$
#3	✓			$63.23 \pm 6.74$	$70.30 \pm 5.79$	$9.74 \pm 3.76$	$2.27 \pm 0.89$	$65.52 \pm 13.97$
#4	✓	✓		$64.20 \pm 6.53$	$73.08 \pm 5.06$	$10.22 \pm 3.91$	$1.79 \pm 0.62$	$66.43 \pm 13.92$
#5			✓	$64.00 \pm 6.10$	$71.65 \pm 4.91$	$9.90 \pm 3.97$	$1.81 \pm 0.61$	$66.18 \pm 13.56$
#6		HarmonySeg		$66.79 \pm 6.34$	$72.04 \pm 5.22$	$9.60 \pm 3.98$	$1.96 \pm 0.73$	$67.17 \pm 14.39$

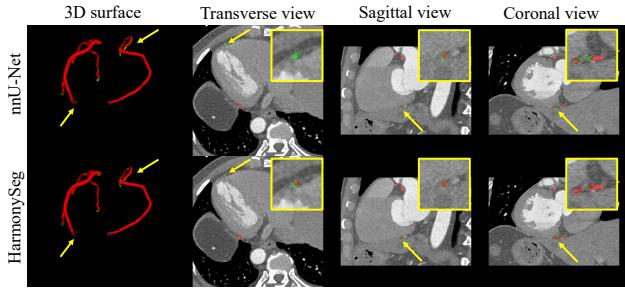


Figure 7. Qualitative comparison of coronary artery segmentation.

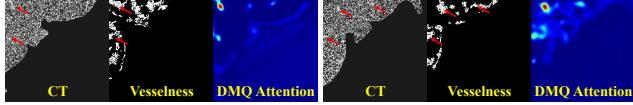


Figure 8. Attention visualization examples of the DMQ module.

#### 4.4. Ablation studies

We conduct four validation studies on the HVS dataset. First, we demonstrate the effectiveness of HarmonySeg’s three key components: D2SD (Deep-to-Shallow Decoding), SADF (Shallow and Deep Fusion), and GSB (Growth-Suppression Balance). Next, we perform ablation studies on the GSB module to evaluate the impact of different loss combinations.

**Ablations of HarmonySeg’s key components:** The results are summarized in Table 6. Using nnUNet as the backbone, the model enhanced with D2SD (Deep-to-Shallow Decoding, Model #3) achieves significant improvements, with a 5.1% increase in Dice and a 2.6% reduction in HD, by leveraging local invariance and detailed spatial information from shallow-scale features. Regarding vesselness utilization, a comparison between Model #1 and Model #2 reveals that simply concatenating CT and vesselness maps ( $\mathbf{VS}_{\text{Cat}}$ ) does not yield notable performance gains, indicating that this fusion method fails to effectively enhance potential vessel regions. To better exploit vesselness information, we introduce SADF (Shallow and Deep Fusion) in Model #4. The

Table 7. **Effectiveness of loss functions.**

$\mathcal{L}_{\text{r-sup}}$	Loss functions			HVS	
	$\mathcal{L}_{\text{con}}$	$\mathcal{L}_{\text{spatial}}$	$\mathcal{L}_{\text{mix}}$	Dice(%,▶)	ASSD( $\downarrow$ )
-	-	-	-	$60.15 \pm 9.49$	$2.23 \pm 0.85$
✓				$63.09 \pm 7.22$	$1.89 \pm 0.70$
	✓			$61.26 \pm 11.1$	$1.76 \pm 0.57$
		✓		$61.82 \pm 11.6$	$1.71 \pm 0.59$
			✓	$62.16 \pm 7.56$	$1.74 \pm 0.57$
✓	✓	✓	✓	$64.00 \pm 6.10$	$1.81 \pm 0.61$

results show improvements in clDice and ASSD, demonstrating that SADF, combined with D2SD, more effectively utilizes vesselness to highlight vessel regions. Through deep querying with CT, the model successfully identifies and focuses on actual vessels. Attention visualization examples are shown in Figure 8. By comparing model #1 and model #5 with GSB in Table 6, GSB improves the Dice score for nnUNet by 6.4%, highlighting its effectiveness in extracting more branches precisely. Visual comparisons are provided in Appendix E.

**Ablations of loss functions:** We evaluate the GSB module by analyzing noise combinations, parameter sensitivity, and computational cost. The robustness of the loss functions and the trade-off between recall and precision are discussed in detail in Appendix G and Appendix H of the supplementary material, respectively. Table 7 show the combination of loss functions:  $\mathcal{L}_{\text{r-sup}}$ ,  $\mathcal{L}_{\text{con}}$ ,  $\mathcal{L}_{\text{spatial}}$ , and  $\mathcal{L}_{\text{mix}}$  individually enhance performance by 2.94%, 1.11%, 1.67%, and 2.01%, respectively. When combined, they synergistically increase the average Dice from 60.15% to 64.0%. We then investigate the impact of the noise suppression weight  $\lambda$  in the combined loss function, as in Eq. (9). As Table 8 shows, Dice improves as  $\lambda$  increases from 0 to 1, but experiences a slight decline when  $\lambda$  exceeds 1. The combined loss computation time ranges from 0.96s to 4.27s per batch, averaging 2.43s. To optimize efficiency, the reconnection loss can be activated after a warm-up phase using other loss functions.

**Table 8. Parameter sensitivity** on the suppression loss weight  $\lambda$ .

$\lambda$	0	0.5	0.75	1	1.5	2
Dice (%)	61.82	62.14	63.45	64.00	63.90	63.64

## 5. Conclusion

In this paper, we propose HarmonySeg for tubular structure segmentation in medical images. Our model used the deep-to-shallow decoding strategy to enhance the model’s adaptability to tubular structures of different sizes. The shallow query and deep mutual query fusion between input images and vesselness filtering results can highlight the potential regions where tubular structures exist. Moreover, we design loss functions to achieve a balance between vessel growth and noise suppression, compensating for the supervision with missing labels. Our model was comprehensively evaluated on four publicly available datasets and the results consistently demonstrated its superiority. A potential improvement is to further integrate the vesselness filter into the network through convolutional operations, combining it with the CT input to form a truly unified entity.

## References

- [1] Ricardo J Araújo, Jaime S Cardoso, and Hélder P Oliveira. Topological similarity index and loss function for blood vessel segmentation. *arXiv preprint arXiv:2107.14531*, 2021. 2
- [2] Samuel G Armato III, Geoffrey McLennan, Luc Bidaut, Michael F McNitt-Gray, Charles R Meyer, Anthony P Reeves, Binsheng Zhao, Denise R Aberle, Claudia I Henschke, Eric A Hoffman, et al. The lung image database consortium (lidc) and image database resource initiative (idri): a completed reference database of lung nodules on ct scans. *Medical physics*, 38(2):915–931, 2011. 5
- [3] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*, 2021. 6
- [4] Jieneng Chen, Jieru Mei, Xianhang Li, Yongyi Lu, Qihang Yu, Qingyue Wei, Xiangde Luo, Yutong Xie, Ehsan Adeli, Yan Wang, et al. Transunet: Rethinking the u-net architecture design for medical image segmentation through the lens of transformers. *Medical Image Analysis*, page 103280, 2024. 6
- [5] Yo-Chuan Chen, Yi-Chen Lin, Ching-Ping Wang, Chia-Yen Lee, Wen-Jeng Lee, Tzung-Dau Wang, and Chung-Ming Chen. Coronary artery segmentation in cardiac ct angiography using 3d multi-channel u-net. *arXiv preprint arXiv:1907.12246*, 2019. 7
- [6] Shunjie Dong, Zixuan Pan, Yu Fu, Qianqian Yang, Yuanxue Gao, Tianbai Yu, Yiyu Shi, and Cheng Zhuo. Deu-net 2.0: Enhanced deformable u-net for 3d cardiac cine mri segmentation. *Medical Image Analysis*, 78:102389, 2022. 2
- [7] Idris Dulau, Catherine Helmer, Cecile Delcourt, and Marie Beurton-Aimar. Ensuring a connected structure for retinal vessels deep-learning segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2364–2373, 2023. 4
- [8] Zhan Gao, Qiuhan Zong, Yiqi Wang, Yan Yan, Yuqing Wang, Ning Zhu, Jin Zhang, Yunfu Wang, and Liang Zhao. Laplacian salience-gated feature pyramid network for accurate liver vessel segmentation. *IEEE Transactions on Medical Imaging*, 42(10):3059–3068, 2023. 1, 2, 5, 6
- [9] Antonio Garcia-Uceda Juarez, Raghavendra Selvan, Zaigham Saghir, and Marleen de Bruijne. A joint 3d unet-graph neural network-based method for airway segmentation from chest cts. In *Machine Learning in Medical Imaging: 10th International Workshop, MLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Proceedings 10*, pages 583–591. Springer, 2019. 6
- [10] Guillaume Garret, Antoine Vacant, and Carole Frindel. Deep vessel segmentation based on a new combination of vesselness filters. *arXiv preprint arXiv:2402.14509*, 2024. 2
- [11] Jiaxing Huang, Yanfeng Zhou, Yaoru Luo, Guole Liu, Heng Guo, and Ge Yang. Representing topological self-similarity using fractal feature maps for accurate segmentation of tubular structures. *arXiv preprint arXiv:2407.14754*, 2024. 2
- [12] Snawar Hussain, Fan Guo, Weiqing Li, and Ziqi Shen. Dilunet: A u-net based architecture for blood vessels segmentation. *Computer Methods and Programs in Biomedicine*, 218:106732, 2022. 1, 2
- [13] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021. 6, 7
- [14] Qiangguo Jin, Zhaopeng Meng, Tuan D Pham, Qi Chen, Leyi Wei, and Ran Su. Dunet: A deformable network for retinal vessel segmentation. *Knowledge-Based Systems*, 178: 149–162, 2019. 2
- [15] Yannick Kirchhoff, Maximilian R Rokuss, Saikat Roy, Balint Kovacs, Constantin Ulrich, Tassilo Wald, Maximilian Zenk, Philipp Vollmuth, Jens Kleesiek, Fabian Isensee, et al. Skeleton recall loss for connectivity conserving and resource efficient segmentation of thin tubular structures. *arXiv preprint arXiv:2404.03010*, 2024. 2, 3, 6
- [16] Bin Kong, Xin Wang, Junjie Bai, Yi Lu, Feng Gao, Kunlin Cao, Jun Xia, Qi Song, and Youbing Yin. Learning tree-structured representation for 3d coronary artery segmentation. *Computerized Medical Imaging and Graphics*, 80: 101688, 2020. 7
- [17] Jonas Lamy, Odyssée Merveille, Bertrand Kerautret, and Nicolas Passat. A benchmark framework for multiregion analysis of vesselness filters. *IEEE Transactions on Medical Imaging*, 41(12):3649–3662, 2022. 1, 2, 6
- [18] Hao Li, Zeyu Tang, Yang Nan, and Guang Yang. Human treelike tubular structure segmentation: A comprehensive review and future perspectives. *Computers in Biology and Medicine*, 151:106241, 2022. 1, 2
- [19] Ruikun Li, Yi-Jie Huang, Huai Chen, Xiaoqing Liu, Yizhou Yu, Dahong Qian, and Lisheng Wang. 3d graph-connectivity

- constrained network for hepatic vessel segmentation. *IEEE Journal of Biomedical and Health Informatics*, 26(3):1251–1262, 2021. 2
- [20] Pechin Lo, Bram Van Ginneken, Joseph M Reinhardt, Tarunashree Yavarna, Pim A De Jong, Benjamin Irving, Catalin Fetita, Margarete Ortner, Rómulo Pinho, Jan Sijbers, et al. Extraction of airways from ct (exact'09). *IEEE Transactions on Medical Imaging*, 31(11):2093–2107, 2012. 5, 6
- [21] David SK Lu, Steven S Raman, Piyaporn Limanond, Donya Aziz, James Economou, Ronald Busuttil, and James Sayre. Influence of large peritumoral vessels on outcome of radiofrequency ablation of liver tumors. *Journal of vascular and interventional radiology*, 14(10):1267–1274, 2003. 1
- [22] Jian Lu, Xiu-Ping Zhang, Bin-Yan Zhong, Wan Yee Lau, David C Madoff, Jon C Davidson, Xiaolong Qi, Shu-Qun Cheng, and Gao-Jun Teng. Management of patients with hepatocellular carcinoma and portal vein tumour thrombosis: comparing east and west. *The lancet Gastroenterology & hepatology*, 4(9):721–730, 2019. 1
- [23] Lei Mou, Yitian Zhao, Huazhu Fu, Yonghuai Liu, Jun Cheng, Yalin Zheng, Pan Su, Jianlong Yang, Li Chen, Alejandro F Frangi, et al. Cs2-net: Deep learning segmentation of curvilinear structures in medical imaging. *Medical image analysis*, 67:101874, 2021. 6
- [24] Mubashir Noman, Mustansar Fiaz, Hisham Cholakkal, Salman Khan, and Fahad Shahbaz Khan. Elgc-net: Efficient local-global context aggregation for remote sensing change detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. 3
- [25] Anton Obukhov, Stamatis Georgoulis, Dengxin Dai, and Luc Van Gool. Gated crf loss for weakly supervised semantic image segmentation. *arXiv preprint arXiv:1906.04651*, 2019. 5
- [26] Yaolei Qi, Yuting He, Xiaoming Qi, Yuan Zhang, and Guanyu Yang. Dynamic snake convolution based on topological geometric constraints for tubular structure segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6070–6079, 2023. 2, 6
- [27] Yulei Qin, Mingjian Chen, Hao Zheng, Yun Gu, Mali Shen, Jie Yang, Xiaolin Huang, Yue-Min Zhu, and Guang-Zhong Yang. Airwaynet: a voxel-connectivity aware approach for accurate airway segmentation using convolutional neural networks. In *International conference on medical image computing and computer-assisted intervention*, pages 212–220. Springer, 2019. 5
- [28] Yulei Qin, Hao Zheng, Yun Gu, Xiaolin Huang, Jie Yang, Lihui Wang, Feng Yao, Yue-Min Zhu, and Guang-Zhong Yang. Learning tubule-sensitive cnns for pulmonary airway and artery-vein segmentation in ct. *IEEE transactions on medical imaging*, 40(6):1603–1617, 2021. 6
- [29] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III* 18, pages 234–241. Springer, 2015. 6
- [30] Ye Shen, Zhijun Fang, Yongbin Gao, Naixue Xiong, Cengsi Zhong, and Xianhua Tang. Coronary arteries segmentation based on 3d fcn with attention gate and level set function. *IEEE Access*, 7:42826–42835, 2019. 7
- [31] Suprosanna Shit, Johannes C Paetzold, Anjany Sekuboyina, Ivan Ezhov, Alexander Unger, Andrey Zhylka, Josien PW Pluim, Ulrich Bauer, and Bjoern H Menze. cldice-a novel topology-preserving loss function for tubular structure segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16560–16569, 2021. 2, 3, 4, 6
- [32] Amber L Simpson, Michela Antonelli, Spyridon Bakas, Michel Bilello, Keyvan Farahani, Bram Van Ginneken, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, et al. A large annotated medical image dataset for the development and evaluation of segmentation algorithms. *arXiv preprint arXiv:1902.09063*, 2019. 1, 5
- [33] Luc Soler, Alexandre Hostettler, Vincent Agnus, Arnaud Charnoz, Jean-Baptiste Fasquel, Johan Moreau, Anne-Blandine Osswald, Mourad Bouhadjar, and Jacques Marescaux. 3d image reconstruction for comparison of algorithm database. URL: <https://www.ircad.fr/research/datasets/liver-segmentation-3d-ircadb-01>, 2010. 5
- [34] Joes Staal, Michael D Abràmoff, Meindert Niemeijer, Max A Viergever, and Bram Van Ginneken. Ridge-based vessel segmentation in color images of the retina. *IEEE transactions on medical imaging*, 23(4):501–509, 2004. 5
- [35] Puyang Wang, Dazhou Guo, Dandan Zheng, Minghui Zhang, Haogang Yu, Xin Sun, Jia Ge, Yun Gu, Le Lu, Xianhua Ye, et al. Accurate airway tree segmentation in ct scans via anatomy-aware multi-class segmentation and topology-guided iterative learning. *IEEE transactions on medical imaging*, 2024. 5
- [36] Yan Wang, Xu Wei, Fengze Liu, Jieneng Chen, Yuyin Zhou, Wei Shen, Elliot K Fishman, and Alan L Yuille. Deep distance transform for tubular structure segmentation in ct scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3833–3842, 2020. 2
- [37] Jelmer M Wolterink, Tim Leiner, and Ivana Isgum. Graph convolutional networks for coronary artery segmentation in cardiac ct angiography. In *Graph Learning in Medical Imaging: First International Workshop, GLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 17, 2019, Proceedings 1*, pages 62–69. Springer, 2019. 7
- [38] Zhe Xu, Donghuan Lu, Yixin Wang, Jie Luo, Jagadeesan Jayender, Kai Ma, Yefeng Zheng, and Xiu Li. Noisy labels are treasure: mean-teacher-assisted confident learning for hepatic vessel segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24*, pages 3–13. Springer, 2021. 2
- [39] Xin Yang, Zhiqiang Li, Yingqing Guo, and Dake Zhou. Dcunet: A deformable convolutional neural network based on cascade u-net for retinal vessel segmentation. *Multimedia Tools and Applications*, 81(11):15593–15607, 2022. 6

- [40] An Zeng, Chunbiao Wu, Guisen Lin, Wen Xie, Jin Hong, Meiping Huang, Jian Zhuang, Shanshan Bi, Dan Pan, Najeeb Ullah, et al. Imagecas: A large-scale dataset and benchmark for coronary artery segmentation based on computed tomography angiography images. *Computerized Medical Imaging and Graphics*, 109:102287, 2023. 6
- [41] Minghui Zhang, Hanxiao Zhang, Guang-Zhong Yang, and Yun Gu. Cfda: Collaborative feature disentanglement and augmentation for pulmonary airway tree modeling of covid-19 cts. In *International conference on medical image computing and computer-assisted intervention*, pages 506–516. Springer, 2022. 6
- [42] Xiao Zhang, Jingyang Zhang, Lei Ma, Peng Xue, Yan Hu, Dijia Wu, Yiqiang Zhan, Jun Feng, and Dinggang Shen. Progressive deep segmentation of coronary artery via hierarchical topology learning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 391–400. Springer, 2022. 2
- [43] Gangming Zhao, Kongming Liang, Chengwei Pan, Fandong Zhang, Xianpeng Wu, Xinyang Hu, and Yizhou Yu. Graph convolution based cross-network multiscale feature fusion for deep vessel segmentation. *IEEE transactions on medical imaging*, 42(1):183–195, 2022. 2
- [44] Hao Zheng, Yulei Qin, Yun Gu, Fangfang Xie, Jiayuan Sun, Jie Yang, and Guang-Zhong Yang. Refined local-imbalance-based weight for airway segmentation in ct. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 410–419. Springer, 2021. 6
- [45] Hao Zheng, Yulei Qin, Yun Gu, Fangfang Xie, Jie Yang, Jiayuan Sun, and Guang-Zhong Yang. Alleviating class-wise gradient imbalance for pulmonary airway segmentation. *IEEE transactions on medical imaging*, 40(9):2452–2462, 2021. 6
- [46] Xiang Zhong, Hongbin Zhang, Guangli Li, and Donghong Ji. Do you need sharpened details? asking mmdc-net: multi-layer multi-scale dilated convolution network for retinal vessel segmentation. *Computers in Biology and Medicine*, 150: 106198, 2022. 1, 2
- [47] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings* 4, pages 3–11. Springer, 2018. 7

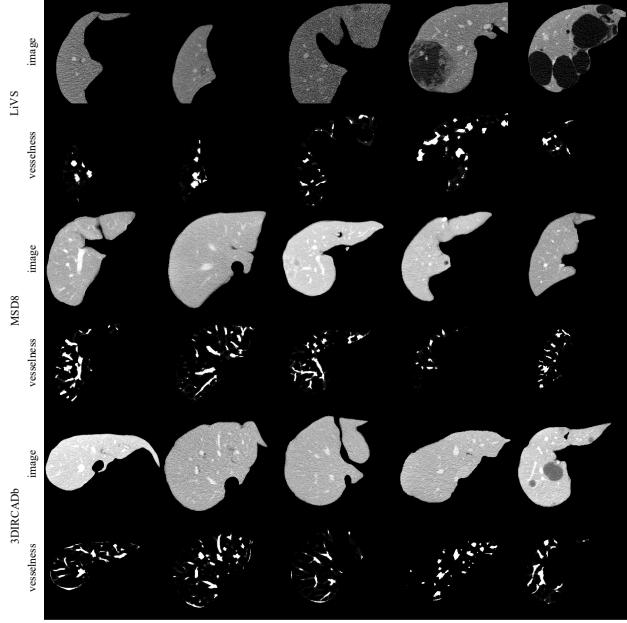


Figure 1. Visualization of images and their corresponding vesselness filtering results. In the CT images, regions with high intensities represent hepatic vessels, while dark regions indicate tumors. In the paired vesselness filtering results, high-intensity patches correspond to vessel candidates, with noise visibly present along the liver border.

Table 1. Enhanced ratio of labeled slices on the HVS task brought by refined labels (%).

Hepatic vessel segmentation (HVS)			
Label	LiVS	MSD8	3DIRCADb (Testing)
Original	$20.26 \pm 12.77$	$71.96 \pm 12.29$	$78.87 \pm 7.27$
Refined	$70.26 \pm 11.75$	$72.54 \pm 11.90$	

In the supplementary material, we provide detailed explanations of the vesselness filter (Appendix A), the flexible convolution block (Appendix B), the segmentation fusion in D2SD (Appendix C), along with additional experimental and visualization results, including refined hepatic vessel labels (Appendix D), ablation studies (Appendix E), our curated HVS-External dataset (Appendix F), the robustness of loss functions (Appendix G), and the trade-off between precision and recall of loss functions (Appendix H).

Table 2. Quantitative improvement on the HVS task brought by refined labels. The segmentation model is nnU-Net.

Dataset	Label	Dice(%, $\uparrow$ )	HD( $\downarrow$ )
HVS	Original	$56.69 \pm 8.42$	$10.83 \pm 4.28$
	Refined	$60.15 \pm 9.49$	$10.00 \pm 4.22$
HVS-External	Original	$63.09 \pm 10.90$	$4.02 \pm 1.43$
	Refined	$63.78 \pm 9.27$	$3.69 \pm 1.17$

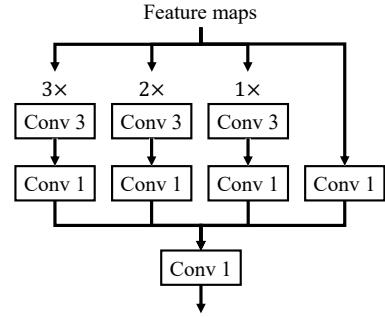


Figure 2. Details of flexible convolution block with diverse receptive fields.

## A. Vesselness Filter

In this section, we elaborate on the details of vesselness filters and present several vesselness maps for visualization. Image derivatives, including first-order derivatives for border detection and second-order derivatives for shape extraction, are commonly used to highlight vascular structures in images [1]. Hessian matrix analysis is a representative method based on second-order derivatives that can distinguish rounded, tubular, and planar structures [3]. Vesselness filters also employ eigen-decomposition of the Hessian matrix to measure tubularity and enhance vessel regions [8]. Let  $H$  be the hessian matrix of a voxel in CT volume, and  $e_1$ ,  $e_2$  and  $e_3$  be the three eigenvectors of  $H$  with corresponding eigenvalues of  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  ( $|\lambda_1| \leq |\lambda_2| \leq |\lambda_3|$ ). The tubularity is defined as [9]:

$$|\lambda_1| \approx 0, \lambda_2 \approx \lambda_3 \ll 0. \quad (1)$$

Based on this, the Jerman [6] vesselness filter used in our framework further regularizes  $\lambda_3$  to reduce the sensitivity

Table 3. Quantitative comparison on the HVS-External test set stratified by diseases of the subjects.

Model	HVS-External							
	Fatty liver		Cirrhosis		Tumor		Healthy	
	Dice(%,▶)	HD(▶)	Dice(%,▶)	HD(▶)	Dice(%,▶)	HD(▶)	Dice(%,▶)	HD(▶)
nnU-Net [5]	53.97 <sub>&amp;#00b1;1.97</sub>	3.90 <sub>&amp;#00b1;0.09&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;69.83&lt;sub&gt;&amp;#00b1;9.15&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;2.78&lt;sub&gt;&amp;#00b1;0.64&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;62.88&lt;sub&gt;&amp;#00b1;10.14&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.91&lt;sub&gt;&amp;#00b1;1.29&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;62.09&lt;sub&gt;&amp;#00b1;2.75&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.94&lt;sub&gt;&amp;#00b1;1.08&lt;/sub&gt;&lt;/td&gt;&lt;/tr&gt; &lt;tr&gt; &lt;td&gt;nnU-Net w clDiceLoss [11]&lt;/td&gt;&lt;td&gt;50.80&lt;sub&gt;&amp;#00b1;3.10&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;4.16&lt;sub&gt;&amp;#00b1;0.11&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;63.48&lt;sub&gt;&amp;#00b1;2.72&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.51&lt;sub&gt;&amp;#00b1;0.18&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;57.18&lt;sub&gt;&amp;#00b1;7.33&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;4.40&lt;sub&gt;&amp;#00b1;1.31&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;49.61&lt;sub&gt;&amp;#00b1;8.00&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;4.96&lt;sub&gt;&amp;#00b1;0.72&lt;/sub&gt;&lt;/td&gt;&lt;/tr&gt; &lt;tr&gt; &lt;td&gt;nnU-Net w SRLoss [7]&lt;/td&gt;&lt;td&gt;58.88&lt;sub&gt;&amp;#00b1;1.69&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;4.08&lt;sub&gt;&amp;#00b1;1.23&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;81.05&lt;/b&gt;&lt;sub&gt;&amp;#00b1;2.22&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.21&lt;sub&gt;&amp;#00b1;0.92&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;70.28&lt;sub&gt;&amp;#00b1;9.60&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.35&lt;sub&gt;&amp;#00b1;1.16&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;68.73&lt;sub&gt;&amp;#00b1;4.72&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.56&lt;sub&gt;&amp;#00b1;1.16&lt;/sub&gt;&lt;/td&gt;&lt;/tr&gt; &lt;tr&gt; &lt;td&gt;TransU-Net3D [2]&lt;/td&gt;&lt;td&gt;51.68&lt;sub&gt;&amp;#00b1;7.78&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;6.03&lt;sub&gt;&amp;#00b1;2.99&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;38.18&lt;sub&gt;&amp;#00b1;31.76&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;4.43&lt;sub&gt;&amp;#00b1;1.73&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;39.93&lt;sub&gt;&amp;#00b1;21.86&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;9.76&lt;sub&gt;&amp;#00b1;9.17&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;56.60&lt;sub&gt;&amp;#00b1;12.92&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;5.23&lt;sub&gt;&amp;#00b1;1.67&lt;/sub&gt;&lt;/td&gt;&lt;/tr&gt; &lt;tr&gt; &lt;td&gt;DSC-Net [10]&lt;/td&gt;&lt;td&gt;&lt;u&gt;71.44&lt;/u&gt;&lt;sub&gt;&amp;#00b1;4.51&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;3.10&lt;/b&gt;&lt;sub&gt;&amp;#00b1;0.39&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;80.50&lt;sub&gt;&amp;#00b1;3.68&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;u&gt;2.73&lt;/u&gt;&lt;sub&gt;&amp;#00b1;0.66&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;u&gt;74.86&lt;/u&gt;&lt;sub&gt;&amp;#00b1;7.04&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;3.18&lt;/b&gt;&lt;sub&gt;&amp;#00b1;1.34&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;u&gt;71.26&lt;/u&gt;&lt;sub&gt;&amp;#00b1;7.90&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.69&lt;sub&gt;&amp;#00b1;1.10&lt;/sub&gt;&lt;/td&gt;&lt;/tr&gt; &lt;tr&gt; &lt;td&gt;HarmonySeg&lt;/td&gt;&lt;td&gt;&lt;b&gt;73.15&lt;/b&gt;&lt;sub&gt;&amp;#00b1;1.83&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.74&lt;sub&gt;&amp;#00b1;0.11&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;80.59&lt;sub&gt;&amp;#00b1;4.61&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;2.20&lt;/b&gt;&lt;sub&gt;&amp;#00b1;0.55&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;75.67&lt;/b&gt;&lt;sub&gt;&amp;#00b1;7.36&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;3.61&lt;sub&gt;&amp;#00b1;1.31&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;74.75&lt;/b&gt;&lt;sub&gt;&amp;#00b1;6.05&lt;/sub&gt;&lt;/td&gt;&lt;td&gt;&lt;b&gt;3.37&lt;/b&gt;&lt;sub&gt;&amp;#00b1;1.18&lt;/sub&gt;&lt;/td&gt;&lt;/tr&gt; &lt;/tbody&gt; &lt;/table&gt; &lt;/div&gt; &lt;div data-bbox="95 271 478 503" data-label="Image"&gt; &lt;img alt="Figure 3: Visualization of original and refined hepatic vessel labels in LiVS. The figure shows four rows of liver segmentation results. The first row is labeled 'Original labels' and the second row is labeled 'Refined labels'. Each row contains five liver volumes, each with yellow and cyan vessels overlaid. The yellow vessels represent the original labels, and the cyan vessels represent the refined labels. The refined labels appear more accurate and complete than the original labels, especially in areas with complex vessel structures."/&gt; &lt;/div&gt; &lt;div data-bbox="92 514 483 556" data-label="Caption"&gt; &lt;p&gt;Figure 3. Visualization of original and refined hepatic vessel labels in LiVS. The liver is rendered in gray, while the original and refined vessel labels are denoted in yellow and cyan, respectively.&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="95 574 478 622" data-label="Image"&gt; &lt;img alt="Figure 4: Visualization of predicted hepatic vessel labels in the test set of MSD8. The figure shows five liver volumes with cyan vessels overlaid. The cyan vessels represent the predicted labels from the model. The labels are relatively accurate, showing the major vascular structures of the liver."/&gt; &lt;/div&gt; &lt;div data-bbox="92 634 483 687" data-label="Caption"&gt; &lt;p&gt;Figure 4. Visualization of predicted hepatic vessel labels in the test set of MSD8. The gray color shows the liver and the cyan color denotes the labels. Note that the original label of MSD8’s test set is unavailable.&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="92 720 294 734" data-label="Text"&gt; &lt;p&gt;for those low-contrast regions:&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="103 734 483 790" data-label="Equation-Block"&gt; &lt;math display="block"&gt;F = \begin{cases} 0, &amp;amp; \lambda_2 \leq 0 \text{ or } \lambda_p \leq 0, \\ 1, &amp;amp; \lambda_2 \geq \frac{\lambda_p}{2} &amp;gt; 0, \\ \lambda_2^2(\lambda_p - \lambda_2)(\frac{3}{\lambda_p + \lambda_2})^3, &amp;amp; \text{otherwise,} \end{cases} \quad (2)&lt;/math&gt; &lt;/div&gt; &lt;div data-bbox="92 790 158 804" data-label="Text"&gt; &lt;p&gt;in which:&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="97 802 483 870" data-label="Equation-Block"&gt; &lt;math display="block"&gt;\lambda_p = \begin{cases} \lambda_3, &amp;amp; \lambda_3 &amp;gt; \tau \max_x \lambda_3(x), \\ \lambda_3 &amp;gt; \tau \max_x \lambda_3(x), &amp;amp; 0 &amp;lt; \lambda_3 \leq \tau \max_x \lambda_3(x), \\ 0, &amp;amp; \text{otherwise,} \end{cases} \quad (3)&lt;/math&gt; &lt;/div&gt; &lt;div data-bbox="92 870 483 900" data-label="Text"&gt; &lt;p&gt;where &lt;math&gt;\tau \in [0, 1]&lt;/math&gt;. Benefiting from this regularization, the Jerman vesselness filter becomes robust even when facing&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="511 274 920 348" data-label="Text"&gt; &lt;p&gt;non-homogeneous vessel intensity. To reveal the effectiveness of the vesselness filter, we give examples of paired vesselness filtering results in Figure 1. As shown in the figure, the vesselness filter can highlight liver vessel candidates of different sizes, even for cases in which tumors exist.&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="511 365 770 382" data-label="Section-Header"&gt; &lt;h2&gt;B. Flexible Convolution Block&lt;/h2&gt; &lt;/div&gt; &lt;div data-bbox="511 393 920 632" data-label="Text"&gt; &lt;p&gt;Diversifying the receptive fields of convolutions is an effective way to adapt models to targets of different sizes [15]. In our study, the sizes of liver vessels are also various, so diverse receptive fields are beneficial in enhancing the model capability. The flexible convolution block we designed is shown in Figure 2. To avoid the gridding effect of dilated convolution for extracting local details of vessels, our flexible convolution block provides different receptive fields by stacking the convolutions in parallel rather than using dilated convolution. After the feature maps are fed into this block, they are further encoded by parallel stacked convolutions with different receptive fields [12]. Then a &lt;math&gt;1 \times 1 \times 1&lt;/math&gt; convolution integrates all features and compresses the channel for output. Flexible convolution blocks are used at the encoder and the shallow query module (F-Conv in (c) of Figure 2 in the manuscript).&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="511 647 796 664" data-label="Section-Header"&gt; &lt;h2&gt;C. Segmentation Fusion in D2SD&lt;/h2&gt; &lt;/div&gt; &lt;div data-bbox="511 674 920 900" data-label="Text"&gt; &lt;p&gt;Vessels of varying sizes exhibit distinct feature representations at different scales. Larger vessels can be effectively reflected in multi-scale feature maps. Yet, for smaller vessels, the information loss caused by successive convolutions and pooling tends to impair their feature representation, which is also one of the motivations why skip connections have been introduced. To mitigate this, the D2SD strategy uses low-cost pre-decoders at multiple scales to capture scale-specific information and aggregate multi-scale outputs for final segmentation, as shown in Figure 5. It is important to clarify that the D2SD is distinct from the deep supervision, which does not compute loss for each decoded result. Concretely, the pre-decoder further facilitates the alignment and aggregation between features of vesselness and CT across different scales within the &lt;math&gt;SQ_i&lt;/math&gt; through the&lt;/p&gt; &lt;/div&gt; &lt;div data-bbox="490 921 506 937" data-label="Page-Footer"&gt; &lt;p&gt;2&lt;/p&gt; &lt;/div&gt;</sub>						

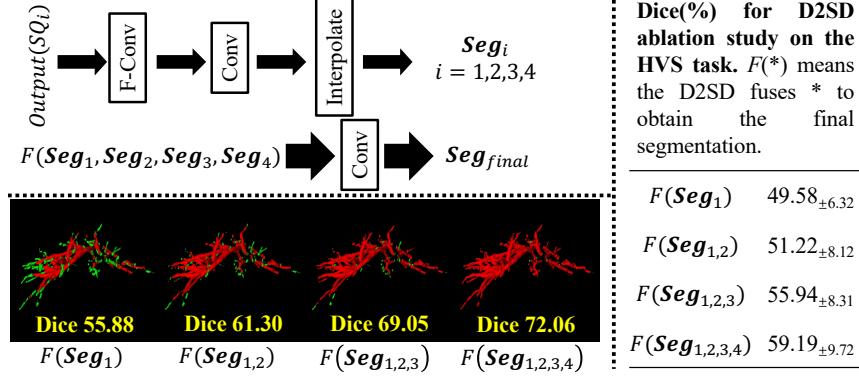


Figure 5. Qualitative and quantitative analysis of segmentation fusion in D2SD: red/green indicating the segmentation and the corresponding labels.

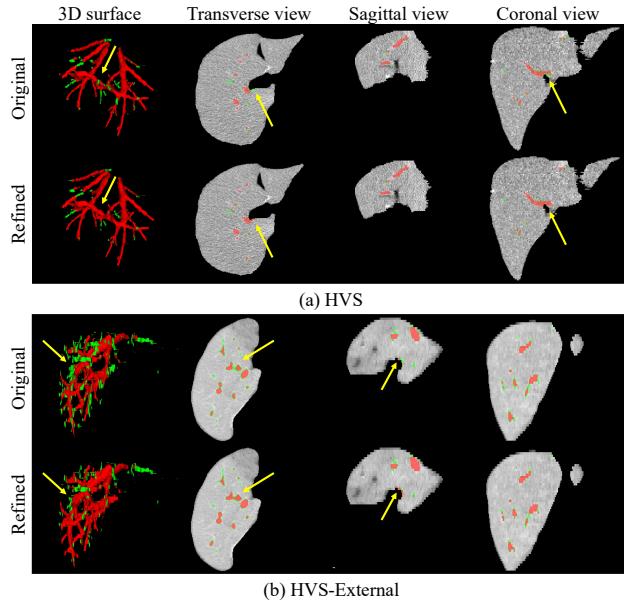


Figure 6. Visualization of hepatic vessel segmentation results using nnU-Net trained on both original and refined labels, in which the red indicates the segmentation and the green indicates the corresponding labels. Improvements are highlighted with yellow arrows.

use of F-Conv, and then, the interpolation is conducted to obtain pre-decoded results that exhibit varying sensitivities to vessel sizes at different scales, which reduces the impact of varying liver vessel sizes. These results are further fused to become the final segmentation, upon which the loss function is calculated.

## D. Refined Hepatic Vessel Labels

In this paper, we use a combined liver vessel segmentation dataset called the HVS dataset. It is based on three publicly available datasets, including LiVS [4], MSD8 [13], and 3DIRCADb [14]. 532, 440, and 20 cases are available for the three datasets, respectively. The three publicly available datasets have made an impressive contribution to developing hepatic vessel segmentation models. However, some slices of the LiVS dataset are insufficiently labeled, and the

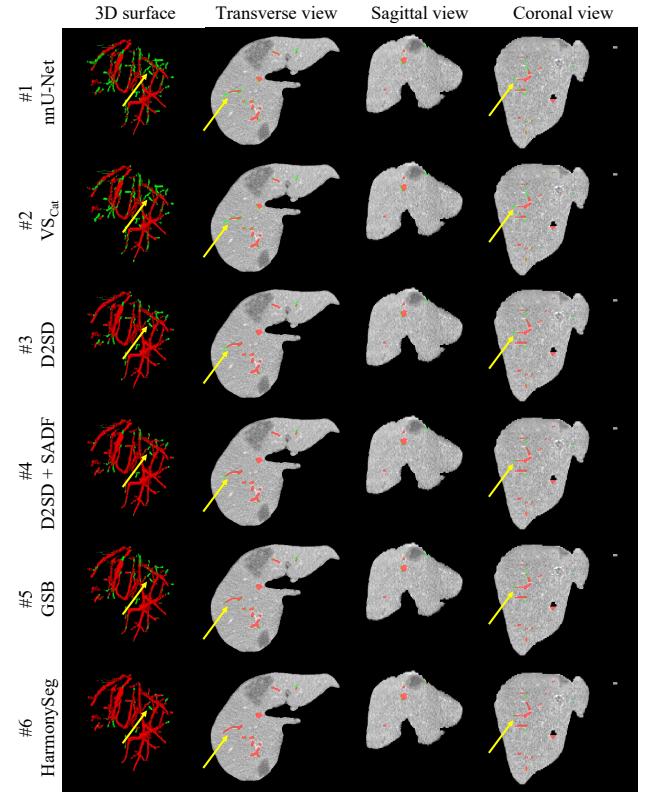


Figure 7. Visualization of hepatic vessel segmentation results in ablation study, with red indicating the segmentation and green representing the corresponding labels. Improvements are highlighted with yellow arrows.

lively available datasets, including LiVS [4], MSD8 [13], and 3DIRCADb [14]. 532, 440, and 20 cases are available for the three datasets, respectively. The three publicly available datasets have made an impressive contribution to developing hepatic vessel segmentation models. However, some slices of the LiVS dataset are insufficiently labeled, and the

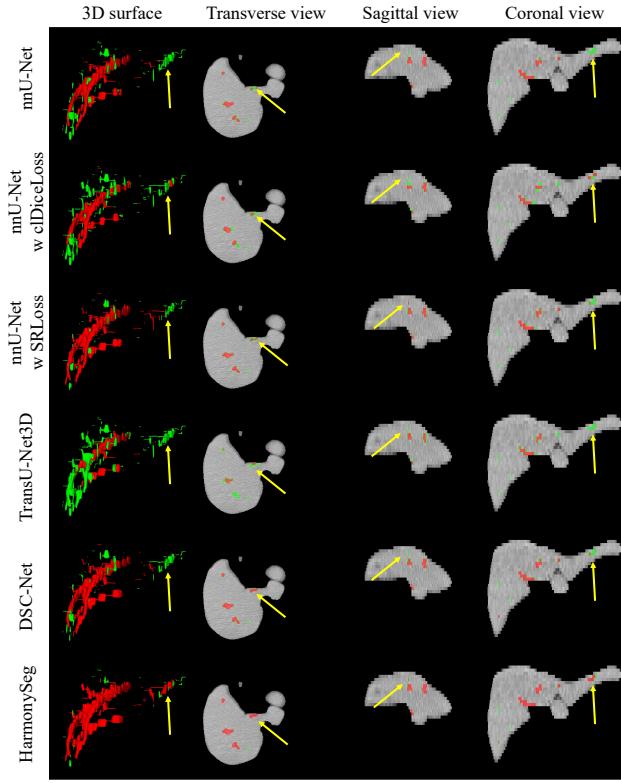


Figure 8. Visualization of hepatic vessel segmentation results in the HVS-External, with red indicating the segmentation and green representing the corresponding labels. Improvements are highlighted with yellow arrows.

labels of the test set of MSD8 are unavailable. The mentioned situations are reflected by labeled ratios (defined by the labeled slices divided by the total slice number of a 3D volume) in Table 1. Thus, to develop our model, we aim to make the best use of the data and refine these hepatic vessel labels. Fortunately, our clinical cooperator, after carefully checking the labels of the MSD8 training set, considered them to be relatively well labeled. Inspired by this, we first trained our model based on the training set of MSD8 and then used it to infer hepatic vessels of the LiVS dataset and the test set of MSD8. Subsequently, the pseudo labels were fused with the raw labels. Fused labels were checked again and manually corrected by a clinician, to serve as the final hepatic vessel labels in the HVS task. From Table 1, it can be found that the ratio of labeled slices has been significantly improved after our optimization, especially the LiVS dataset. Moreover, more visualization examples are given in Figures 3 and 4. Due to the cropping of the CT volume by the organizers of the dataset, the presence of some lesions, such as tumors, and the slice thickness, the continuity of the refined vessel labels is not fully ensured. Still, they are significantly improved compared to the original ones. We also

compare the baseline performance trained by the original labels and the refined ones. As indicated by the evaluation metrics in Table 2 and the visualization examples in Figure 6, the baseline trained by the refined labels performs better in the HVS task.

## E. Ablation studies

Some visualization examples in ablation studies are shown in Figure 7, it can be seen that our D2SD strategy can extract vessels with diverse sizes more effectively compared to the baseline. Moreover, liver vessel segmentation can not benefit from the simple concatenation fusion between the images and corresponding vesselness filtering results. Instead, our SADF fusion module can better utilize the vesselness filtering result to improve the segmentation accuracy. Besides, it can be observed that the GSB further preserves a reasonable continuity of the vessel tree.

## F. Analysis on HVS-External

In the HVS-External, we included cases with various liver diseases, including two cases of fatty liver, four cases of cirrhosis, twelve cases of liver tumors, and three healthy livers. We analyze the results of HVS-External based on the disease stratification, as shown in Table 3. It can be found that the HarmonySeg achieves the highest mean Dice for patients with fatty liver, tumors, and healthy individuals, and mean HDs are competitive compared with other methods. Furthermore, visualization examples are demonstrated in Figure 8. The results indicate the robustness of HarmonySeg to various liver diseases and the potential to be applied in clinical practices.

## G. Robustness discussion

We recognize that the reconnection loss may introduce noise. To address this, we observe that the performance gains of our loss functions following this order:  $\mathcal{L}_{\text{sup-r}}(2.94\%) > \mathcal{L}_{\text{mix}}(2.01\%) > \mathcal{L}_{\text{spatial}}(1.67\%) > \mathcal{L}_{\text{con}}(1.01\%)$ . The first three losses ( $\mathcal{L}_{\text{sup-r}}$ ,  $\mathcal{L}_{\text{mix}}$ ,  $\mathcal{L}_{\text{spatial}}$ ) are robust and applicable to various scenarios. In contrast, the reconnection loss  $\mathcal{L}_{\text{con}}$  is specifically designed to address missing vessel segments. To enhance its robustness, we employ two strategies: (a) We perform skeletonization on the defined reconnect branches, reducing their pixel width to 1, as shown in Eq.(6) of manuscript. Consequently, the loss applied to these pixels remains slight on average. (b) If incorrect pixels are mistakenly defined for reconnection, they can be effectively suppressed by the strong regularization from spatial relationships and mix augmentation invariance. Thus, we incorporate the reconnection loss as an additional strategy tailored for vessel segmentation tasks.

Table 4. **Ablations on recall and precision trade-off:**  $\mathcal{L}^+$  for growth,  $\mathcal{L}^-$  for suppression.

$\mathcal{L}_{\text{r-sup}}^+$	$\mathcal{L}_{\text{con}}^+$	$\mathcal{L}_{\text{spatial}}^-$	$\mathcal{L}_{\text{mix}}^-$	Recall (%)	Precision (%)	F1-score (%)
-	-	-	-	49.14	<b>84.12</b>	62.04
✓				55.57	79.12	63.09
✓	✓			<b>64.35</b>	71.20	<u>65.01</u>
✓		✓		50.33	79.92	59.68
✓			✓	53.58	<u>80.15</u>	62.15
✓	✓	✓	✓	59.93	73.89	<b>66.18</b>

## H. Trade-off between precision and recall

We also analyze the recall-precision trade-off. The recall rates for different loss combinations on HVS are presented in Table 4. Recall is enhanced through relaxed supervision ( $\mathcal{L}_{\text{r-sup}}$ ) and branch reconnection ( $\mathcal{L}_{\text{con}}$ ), while noise is reduced via spatial consistency ( $\mathcal{L}_{\text{spatial}}$ ) and mix equivalence ( $\mathcal{L}_{\text{mix}}$ ). We achieved the best trade-off and the highest F1-score when combining all loss functions.

## References

- [1] Gady Agam, Samuel G Armato, and Changhua Wu. Vessel tree reconstruction in thoracic ct scans with application to nodule detection. *IEEE transactions on medical imaging*, 24(4):486–499, 2005. 1
- [2] Jieneng Chen, Jieru Mei, Xianhang Li, Yongyi Lu, Qihang Yu, Qingyue Wei, Xiangde Luo, Yutong Xie, Ehsan Adeli, Yan Wang, et al. Transunet: Rethinking the u-net architecture design for medical image segmentation through the lens of transformers. *Medical Image Analysis*, page 103280, 2024. 2
- [3] Alejandro F Frangi, Wiro J Niessen, Koen L Vincken, and Max A Viergever. Multiscale vessel enhancement filtering. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI’98: First International Conference Cambridge, MA, USA, October 11–13, 1998 Proceedings 1*, pages 130–137. Springer, 1998. 1
- [4] Zhan Gao, Qiuaho Zong, Yiqi Wang, Yan Yan, Yuqing Wang, Ning Zhu, Jin Zhang, Yunfu Wang, and Liang Zhao. Laplacian salience-gated feature pyramid network for accurate liver vessel segmentation. *IEEE Transactions on Medical Imaging*, 42(10):3059–3068, 2023. 3
- [5] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):203–211, 2021. 2
- [6] Tim Jerman, Franjo Pernuš, Boštjan Likar, and Žiga Špiclin. Enhancement of vascular structures in 3d and 2d angiographic images. *IEEE transactions on medical imaging*, 35(9):2107–2118, 2016. 1
- [7] Yannick Kirchhoff, Maximilian R Rokuss, Saikat Roy, Balint Kovacs, Constantin Ulrich, Tassilo Wald, Maximilian Zenk, Philipp Vollmuth, Jens Kleesiek, Fabian Isensee, et al. Skeleton recall loss for connectivity conserving and resource efficient segmentation of thin tubular structures. *arXiv preprint arXiv:2404.03010*, 2024. 2
- [8] Jonas Lamy, Odyssée Merveille, Bertrand Kerautret, and Nicolas Passat. A benchmark framework for multiregion analysis of vesselness filters. *IEEE Transactions on Medical Imaging*, 41(12):3649–3662, 2022. 1
- [9] Cristian Lorenz, I-C Carlsen, Thorsten M Buzug, Carola Fassnacht, and Jürgen Weese. Multi-scale line segmentation with automatic estimation of width, contrast and tangential direction in 2d and 3d medical images. In *International Conference on Computer Vision, Virtual Reality, and Robotics in Medicine*, pages 233–242. Springer, 1997. 1
- [10] Yaolei Qi, Yuting He, Xiaoming Qi, Yuan Zhang, and Guanyu Yang. Dynamic snake convolution based on topological geometric constraints for tubular structure segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6070–6079, 2023. 2
- [11] Suprosanna Shit, Johannes C Paetzold, Anjany Sekuboyina, Ivan Ezhov, Alexander Unger, Andrey Zhylka, Josien PW Pluim, Ulrich Bauer, and Bjoern H Menze. cldice-a novel topology-preserving loss function for tubular structure segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16560–16569, 2021. 2
- [12] Karen Simonyan. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [13] Amber L Simpson, Michela Antonelli, Spyridon Bakas, Michel Bilello, Keyvan Farahani, Bram Van Ginneken, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, et al. A large annotated medical image dataset for the development and evaluation of segmentation algorithms. *arXiv preprint arXiv:1902.09063*, 2019. 3
- [14] Luc Soler, Alexandre Hostettler, Vincent Agnus, Arnaud Charnoz, Jean-Baptiste Fasquel, Johan Moreau, Anne-Blandine Osswald, Mourad Bouhadjar, and Jacques Marescaux. 3d image reconstruction for comparison of algorithm database. URL: <https://www.ircad.fr/research/datasets/liver-segmentation-3d-ircadb-01>, 2010. 3
- [15] Xiang Zhong, Hongbin Zhang, Guangli Li, and Donghong Ji. Do you need sharpened details? asking mmdc-net: multi-layer multi-scale dilated convolution network for retinal vessel segmentation. *Computers in Biology and Medicine*, 150:106198, 2022. 2