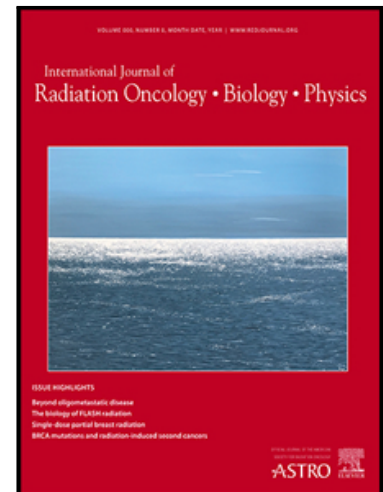


Comprehensive evaluation of a deep learning model for automatic organs at risk segmentation on heterogeneous computed tomography images for abdominal radiotherapy

Wenjun Liao M.D. , Xiangde Luo Ph.D. , Yuan He M.D. ,  
Ye Dong M.D. , Churong Li M.D. , Kang Li Ph.D. ,  
Shichuan Zhang M.D. , Shaoting Zhang , Guotai Wang Ph.D. ,  
Jianghong Xiao Ph.D.

PII: S0360-3016(23)00520-5  
DOI: <https://doi.org/10.1016/j.ijrobp.2023.05.034>  
Reference: ROB 28248



To appear in: *International Journal of Radiation Oncology, Biology, Physics*

Received date: 8 November 2022  
Revised date: 13 March 2023  
Accepted date: 18 May 2023

Please cite this article as: Wenjun Liao M.D. , Xiangde Luo Ph.D. , Yuan He M.D. , Ye Dong M.D. , Churong Li M.D. , Kang Li Ph.D. , Shichuan Zhang M.D. , Shaoting Zhang , Guotai Wang Ph.D. , Jianghong Xiao Ph.D. , Comprehensive evaluation of a deep learning model for automatic organs at risk segmentation on heterogeneous computed tomography images for abdominal radiotherapy, *International Journal of Radiation Oncology, Biology, Physics* (2023), doi: <https://doi.org/10.1016/j.ijrobp.2023.05.034>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

**Title:** Comprehensive evaluation of a deep learning model for automatic organs at risk segmentation on heterogeneous computed tomography images for abdominal radiotherapy

**Running title:** Deep learning for abdominal OARs delineation on CT images

**Authors' names:** Wenjun Liao<sup>1</sup>, M.D., Xiangde Luo<sup>2,6</sup>, Ph.D., Yuan He<sup>3</sup>, M.D., Ye Dong<sup>4</sup>, M.D., Churong Li<sup>1</sup>, M.D., Kang Li<sup>5</sup>, Ph.D., Shichuan Zhang<sup>1</sup>, M.D., Shaoting Zhang<sup>2,6</sup>, Guotai Wang<sup>2,6</sup>, Ph.D., Jianghong Xiao<sup>7\*</sup>, Ph.D.

**Authors' affiliations**

<sup>1</sup> Department of Radiation Oncology, Radiation Oncology Key Laboratory of Sichuan Province, Sichuan Clinical Research Center for Cancer, Sichuan Cancer Hospital & Institute, Sichuan Cancer Center, Affiliated Cancer Hospital of University of Electronic Science and Technology of China, Chengdu, 610041 China

<sup>2</sup> School of Mechanical and Electrical Engineering, University of Electronic Science and Technology of China, Chengdu, 611731, China

<sup>3</sup> Department of Radiation Oncology, The First Affiliated Hospital of USTC, Division of Life Sciences and Medicine, University of Science and Technology of China, Hefei, Anhui, 23000, China

<sup>4</sup> Department of NanFang PET Center, Nanfang Hospital, Southern Medical University, Guangzhou, 510515 China

<sup>5</sup> West China Biomedical Big Data Center, West China Hospital, Sichuan University, Chengdu, 610041, China

<sup>6</sup> Shanghai AI Laboratory, Shanghai, 200030, China

<sup>7</sup> Radiotherapy Physics & Technology Center, Department of Radiation Oncology,

Cancer Center, West China Hospital, Sichuan University, Chengdu 610041, China

### **Corresponding author**

\* Jianghong, Xiao, Ph.D., Radiotherapy Physics & Technology Center, Department of Radiation Oncology, Cancer Center, West China Hospital, Sichuan University, Chengdu 610041, China

Tel.: +86-28 85422909

Fax: +86-28 85422909

E-mail: xiaojh@scu.edu.cn

### **Authors responsible for statistical analyses**

Wenjun Liao, and Xiangde Luo are responsible for statistical analyses.

Wenjun Liao: Tel: +86-18280167590, E-mail: lwjpsy@163.com

Xiangde Luo: Tel: +86-13032863009, E-mail: xiangde.luo@std.uestc.edu.cn

### **Conflict of interest statement**

The authors report no conflicts of interest in this work

### **Funding statement**

This work was supported by the National Natural Science Foundation of China under Grant 82203197

### **Data sharing statement**

The code for neural network training can be obtained at <https://github.com/Luoxd1996/AbsegNet>. Additionally, we will release part of the carefully-annotated dataset with 16 organ annotations (150 volumes from the development cohort and 20 volumes from LiTS2017) to boost this research task.

Other data generated and analyzed during this study can be obtained by contacting the corresponding author with reasonable requirements.

### **Acknowledgements**

None

### **Abstract**

**Purpose:** To develop a deep learning model (AbsegNet) that produces accurate contours of 16 organs at risk (OARs) for abdominal malignancies, as an essential part of a fully automated radiation treatment planning.

**Material and methods:** Three datasets with 544 computed tomography (CT) scans were retrospectively collected. Dataset 1 was split into 300 training cases and 128 test cases (cohort 1) for AbsegNet. Dataset 2, including cohort 2 ( $n = 24$ ) and cohort 3 ( $n = 20$ ), were used to validate AbsegNet externally. Dataset 3, including cohort 4 ( $n = 40$ ) and cohort 5 ( $n = 32$ ) were used to clinically assess the accuracy of AbsegNet-generated contours. Each cohort was from a different center. The Dice similarity coefficient (DSC) and 95th-percentile Hausdorff distance (HD95) were calculated to evaluate the delineation quality for each OAR. Clinical accuracy evaluation was classified into three levels: no revision, minor revisions ( $0\% < \text{volumetric revision degrees (VRD)} \leq 10\%$ ), moderate revisions ( $10\% \leq \text{VRD} < 20\%$ ), and major revisions ( $\text{VRD} \geq 20\%$ ).

**Results:** For all OARs, AbsegNet achieved a mean DSC of 86.73%, 85.65%, and 88.04% in cohort 1, cohort 2, and cohort 3, respectively, and a mean HD95 of

8.92mm, 10.18mm, and 12.40mm, respectively. The performance of AbsegNet outperformed SwinUNETR, DeepLabV3+, Attention-UNet, UNet and 3D-UNet. When experts evaluated contours from cohorts 4 and 5, four OARs (liver, kidney\_L, kidney\_R, and spleen) of all patients were scored as having no revision, and over 87.5% of patients with contours of the stomach, esophagus, adrenals, or rectum were considered as having no or minor revisions. Only 15.0% of patients with contours of the colon and small bowel required major revisions.

**Conclusions:** We propose a novel deep learning model to delineate OARs on diverse datasets. Most contours produced by AbsegNet are accurate and robust, therefore clinically applicable and helpful to facilitate radiotherapy workflow.

#### Keywords

Abdominal radiotherapy; Organ at Risk; deep learning; automatic segmentation; Knowledge distillation

## Introduction

Radiotherapy is one of the most important local treatment modalities for abdominal malignancies, such as cervical, prostate, pancreatic, and hepatic cancers. In the process of radiotherapy administration, delineating abdominal organs at risk (OARs) on computed tomography (CT) images is an essential step. Radiation treatment planning requires accurately calculating radiation dose especially for OARs close to gross tumor volumes. Hence, inaccurate delineation might lead to dose miscalculations and unexpected side effects<sup>[1]</sup>.

In past decades, OAR delineation has been manually conducted by radiation oncologists with slice-by-slice CT images. It is labor-intensive and may take several hours per case. Additionally, manual delineation of OAR often leads to inter- and intra-observer variabilities that can influence treatment outcomes in certain cases<sup>[2-5]</sup>. Therefore, consistent and high-quality abdominal OAR delineation is greatly desired in clinical practice. In this context, full autosegmentation of whole abdominal OARs by deep learning (DL) methods, if feasible, could be more advantageous.

Benefitting from the advantages of feature learning, DL-based automatic delineation has offered a promising solution to solve the problems of manual delineation<sup>[6]</sup>. In recent years, a wide range of DL-based models have been proposed to segment abdominal OARs, and huge progress has been made. For example, in previous studies, spleen and liver segmentation could obtain an accuracy of 96% and 95%, in terms of their Dice similarity coefficients (DSC), respectively<sup>[7, 8]</sup>. In a recent work for kidneys and pancreas segmentation, 93% and 79% DSC were obtained, respectively<sup>[9]</sup>.

However, several limitations still exist in current automatic segmentation models for abdominal OARs. For instance, many abdominal datasets just contain single-institutional, single-scanner, or single-disease patients<sup>[10]</sup>. It is unclear whether the segmentation performance acquired on those datasets might generalize well on more heterogenous data. Huge differences in organ morphological structure, disease status, image appearance, and image quality obtained from various patients by different scanners could affect the accuracy of image segmentation.

The main focus of most studies is on the reporting of segmentation results in their own cohorts, but they fail to perform comprehensive clinical assessment for automatically segmented contours, a critical step in medical image segmentation<sup>[11]</sup>. Van et al reported that it took an average of 34-54 minutes to make the automatic contours suitable for clinical use<sup>[12]</sup>. Other disadvantages should not be ignored. For example, some studies just segmented a few OARs, or in some studies, ground truth OARs utilized networks' predictions as initial results and were refined by experts, which might bring the raters' bias<sup>[10, 13, 14]</sup>.

To address the aforementioned issues, in this paper, we retrospectively collected whole abdominal CT images of 544 patients diagnosed with various tumors from five different centers. Then, we proposed a new DL model, named AbsegNet, that could delineate a comprehensive set of 16 abdominal OARs. The performance of AbsegNet was validated in 172 patients across three different cohorts, and the accuracy was compared to five previous state-of-the-art methods. Moreover, two experienced radiotherapy practitioners were invited to evaluate the accuracy of AbsegNet in another independent two cohorts with 72 patients.

## Material and Methods

### Data summary

Three datasets with a total number of 544 patients in this study were used (Table 1). Dataset 1, collected from XXXX (XXXX), contained 428 planning CT images, which were randomly split into 300 cases for training and 128 cases for internal testing (cohort 1), with a ratio of around 3:1. Dataset 2 included cohort 2 (n = 24) collected

from XXXX (XXXX), and cohort 3 ( $n = 20$ ) collected from a public dataset (LiTS2017)<sup>[15]</sup>. These two cohorts were used to externally validate the performance of AbsegNet. Dataset 3 consisted of cohort 4 ( $n = 40$ ) collected from XXXX (XXXX), and cohort 5 ( $n = 32$ ) from XXXX (XXXX), which were used to clinically assess the accuracy of AbsegNet-generated contours. The flow chart of this study was illustrated in Fig. 1.

Distributions of gender, age, tumor sites, and scanning parameters of these datasets are presented in Table 1. Patients with various tumors with CT images from different scanners using different slices were incorporated into this study. In addition, a comparison of the included datasets with several available public datasets is shown in Table S1. The included datasets mostly covered OARs for abdominal radiotherapy with ample sample sizes from multiple centers.

This retrospective study was approved by the Ethics Committee on Biomedical Research for these hospitals and informed consent was waived.

### Ground truth contours

To ensure data uniformity, abdominal OARs in each case were manually delineated by a radiation oncologist from XXXX with more than eight years of experience in the treatment of abdominal malignancies. After that, another senior oncologist from the same hospital with more than 20 years of experience checked and revised these annotations carefully, and in cases of disagreement, produced consensus annotations. All CT scans of these datasets, with the exception of LiTS2017<sup>[15]</sup>, were exhaustively labeled with 16 anatomical organs, including the liver, spleen, kidneys (left and right),



stomach, gallbladder, esophagus, pancreas, adrenals, duodenum, colon, small bowel, rectum, bladder, and head of the femurs (left and right). Considering the liver was manually previously labeled in LiTS2017, it was excluded and the remaining 15 OARs were labeled in this study. OAR delineation principles were in accordance with relevant radiation guidelines<sup>[16, 17]</sup>. An example CT scan and ground truth contour from dataset 1 is shown in Fig. S1.

All manually delineated contours were performed in the MIM 7.07 Software (Cleveland, OH, USA)<sup>[18, 19]</sup>.

### Segmentation network construction

The AbsegNet framework is shown in the top right of Fig. 1, and detailed architecture parameters are presented in Table S2. In this work, we present a new method to train accurate and robust segmentation networks by employing data augmentation and knowledge distillation, consisting of a teacher-student model. Considering that CT images come from different centers, patients, scanning protocols, tumors, and contrast types may cause data distribution shifts. As expected, there were huge intensity distribution gaps among our datasets (Fig. S2). These distribution shifts could lead to model collapse on unseen centers<sup>[20]</sup>. To boost the network's robustness on unseen datasets, we utilized a wide range of data augmentation strategies to generate different augmented images for network training online, such as intensity transformations (randomly using random noise, sharpening, histogram match, nonlinear transformation and histogram equalization) and spatial transformations (randomly applying rotation, rescaling, elastic deformation). Furthermore, we combined data

augmentation with a general knowledge distillation framework to train segmentation models<sup>[21, 22]</sup>. Specifically, in the training stage, teacher ( $\Theta$ ) and student ( $\Psi$ ) networks take augmented images ( $T^1(i)$  and  $T^2(i)$ ) as inputs and produce corresponding predictions ( $\Theta(T^1(i))$  and  $\Psi(T^2(i))$ ). Here,  $T^1$  is a random noise transformation and  $T^2$  is a random one in the set of intensity and spatial transformations.  $T^1$  and  $T^2$  can be considered weak and strong augmentations, respectively. Then, we encouraged the student to generate predictions based on the teacher via a knowledge distillation loss ( $L_{kd}$ ) according to,

$$L_{kd} = L_{kl}(\Psi(T^2(i)), \Theta(T^1(i))/t)$$

Where  $t$  is a temperature factor that controls the importance of the teacher's predictions and set to 4;  $Kl$  is the Kullback-Leibler divergence function. At the same time, the student network is also supervised by the ground truth ( $gt(i)$ ) via a combination loss as,

$$L_{seg} = 0.5 \times (L_{ce}(\Psi(T^2(i)), gt(i)) + L_{dice}(\Psi(T^2(i)), gt(i)))$$

Where  $ce$  and  $dice$  represent the cross-entropy loss and dice loss, respectively; and the total objective loss function is  $L_{total} = L_{seg} + 0.1 \times L_{kd}$ . Afterwards, the student network updates the parameter by minimizing  $L_{total}$ , and the teacher's parameter is updated as an exponential moving average (EMA) of the student's parameter. Based on the proposed method, the segmentation network is encouraged to learn the anatomical context feature and ignore the intensity distribution to boost the generalization on unseen datasets. In the testing stage, the teacher model was used to produce final results following previous suggestions<sup>[22]</sup>. Different from previous

works<sup>[23-25]</sup>, AbsegNet could be applied on unseen datasets without any finetuning or retraining.

To confirm the usefulness of distillation learning in the segmentation of CT images from various centers, we firstly reported the results of the proposed method with and without knowledge distillation in cohorts 1, 2, and 3, which showed that using knowledge distillation could improve performance in the three cohorts (Table S3). The overall average DSC of AbsegNet with knowledge distillation was significantly higher than that of AbsegNet without knowledge distillation (cohort 1, 86.73% vs. 85.34%; cohort 2, 85.65% vs. 81.98%; cohort 3, 88.04% vs. 84.65%; all  $p$ -values  $< 0.05$ ) (Table S3).

### Pre-processing of images

In the pre-processing phase, all images were reformatted into a standard RAI orientation (Right-to-left, Anterior-to-posterior, and Inferior-to-superior, in the x, y, and z axes, respectively). The intensities (Hounsfield Units) of each image were adjusted based on the grey-level histogram, and cut off intensities outside the 0.5 and 99.5 percentiles. Then, we resampled the images to the fixed resolution of  $0.98 \times 0.98 \times 3.0 \text{ mm}^3$ , which was the medium resolution of the training set. Finally, all images were normalized to zero mean and unit variance. In the post-processing phase, the largest connected component selection and morphological operation were used to refine the network's predictions and generate final results.

### Implementation details

The proposed method was implemented by PyTorch on a Ubuntu20.04 desktop with

two NVIDIA 3090 GPUs. A 3D-UNet<sup>[23]</sup> was used as a baseline model. The total epoch was set to 1000, and batch size was 2. The input patch size was randomly cropped from the pre-processed image with a shape of  $64 \times 192 \times 192$ . We employed a set of data augmentation strategies and knowledge distillation to train the network (detailed in Supplementary material and methods). We used the stochastic gradient descent optimizer (weight decay =  $10^{-4}$ , momentum = 0.9) to update network parameters. The initial learning rate was 0.01 and adjusted by a poly learning rate strategy. In the testing stage, we used a sliding window strategy with a stride of  $32 \times 96 \times 96$  to produce final predictions.

### Methods comparison

In this work, AbsegNet was compared with five famous and widely-used methods: (1) UNet, which presents a U-shape encoder-decoder network for biomedical image segmentation and achieves very promising results on many tasks<sup>[26]</sup>; (2) 3D-UNet, an extension of UNet from 2D space to 3D space for volumetric image segmentation<sup>[23]</sup>; (3) DeepLabV3+, an encoder-decoder with an atrous separable convolution network for natural image semantic segmentation<sup>[24]</sup>; (4) Attention-UNet, which extends UNet attention gates to focus on target structures of varying shapes and sizes for better segmentation results<sup>[25]</sup>; and (5) SwinUNETR, a new combination framework that utilizes the merits of Swin Transformers and U-Shape networks for medical image segmentation and achieves encouraging performance<sup>[27]</sup>. To ensure consistency of comparisons, public implementations of the methods were used to directly train the network based on the same training dataset and procedures (Table S4).

## Assessment of AbsegNet-generated contours by experts

Two senior experts (A and B) with more than 15 years of radiation experience from XXXX and XXXX, respectively, were invited to assess the accuracy of AbsegNet-generated contours from cohorts 4 and 5 (Fig. 1). Each expert was required to revise incorrect OAR segmentation when necessary. In the process of correction, experts were blinded to the ground truth contours, and encouraged to obey the same delineation guidelines described above. Considering the heavy burden of annotating with 16 OARs, seven representative OARs, including some solid and gastrointestinal organs, were assigned to expert A, and the other nine OARs, also including some of those organ types, were assigned to expert B. The kidneys (left and right), pancreas, duodenum, bladder, and femurs (left and right) were reviewed by expert A. The liver, spleen, stomach, gallbladder, esophagus, adrenals, colon, small bowel, and rectum were reviewed by expert B.

Next, AbsegNet-predicted contours were compared to their corresponding revised contours to calculate volumetric revision degrees (VRD), which were defined as the volume required to be edited divided by the volume of AbsegNet-generated contours multiplied by 100<sup>[28]</sup>. Accuracy was classified into four levels: no revision (VRD = 0%), minor revisions ( $0 < \text{VRD} \leq 10\%$ ), moderate revisions ( $10\% < \text{VRD} < 20\%$ ), and major revisions ( $\text{VRD} \geq 20\%$ ).

Qualitative analysis was also an important part of our research. Similar to a previous work<sup>[29]</sup>, experts A and B were invited to subjectively evaluate each automatic contouring result together. At first, we randomly selected five cases from

each cohort. Then, autosegmentations from these 25 patients were evaluated by the two experts. We used 3-grade criteria to estimate the degree of clinically acceptable: (1) completely acceptable (the prediction can be used in the treatment planning without any revision), (2) acceptable (the prediction needs a few refinements but has no obvious clinical impact without corrections), and (3) unacceptable (the prediction needs to be substantially revised before treatment planning or needs to be re-delineated manually).

Besides, a study was performed to compare the time spent by expert A in delineating OARs under two modes: with or without assistance from AbsegNet. In the first mode, the contours of 16 OARs produced by AbsegNet were provided to expert A, who would then examine the predictions and revise the incorrect ones when necessary. In the second mode, the OARs delineation was conducted completely manually. The contouring time for each patient includes the time spent on verifying the results, as well as the time spent on revising the model's predictions. We randomly selected 10 cases from the 25 patients mentioned above to conduct this experiment.

### Evaluation Metrics

To evaluate the performance of AbsegNet, the volumetric DSC and 95th percentile Hausdorff distance (HD95) were adopted, which are the most commonly used metrics in this field<sup>[30]</sup>. The DSC measures the volumetric overlap between two contours, and the HD measures the boundaries of two contours. Due to max HD being very sensitive to outliers, HD95, which measures the 95th percentile distance between two

contours, is often used instead<sup>[31]</sup>.

## Statistical analysis

Statistical analysis was performed using an SPSS software package (Version 22.0, IBM SPSS Inc). Numeric variables were denoted as mean  $\pm$  standard deviation, and compared by paired t-test when necessary. A two-tailed  $p$ -value  $< 0.05$  was considered significant.

## Results

### Intraobserver variability examination

To determine the intraobserver variability, ten CT images were randomly selected from dataset 1, and recontoured by the same expert after an interval of two months. The second contours were compared with the first corresponding contours to calculate DSC and HD95. We found that the distribution of mean DSC (mDSC) for 12 out of 16 OARs was around 95%, and the mean HD95 (mHD95) across all OARs was less than 5 mm (Fig. S3). It was suggested that the intraobserver discrepancy was minor, and the annotations were reliable.

### Performance comparison in internal testing cohort

The DSC and HD95 of 16 OARs obtained by six DL algorithms in cohort 1 are listed in Table 2. In terms of mDSC, AbsegNet showed the best performance in 14 out of 16 OARs among six algorithms, and achieved a mDSC greater than 90% for eight out of 16 OARs. Only two OARs (adrenals and duodenum) had a mDSC below 80%. Considering all OARs as a whole, the mDSC was 86.73%, 84.39%, 82.11%, 84.67%, 84.66%, and 82.82% for AbsegNet, SwinUNETR, DeepLabV3+, Attention-UNet,

UNet, and 3D-UNet, respectively. In terms of HD95, AbsegNet showed the best performance for 12 OARs. We also observed that five OARs produced by AbsegNet had a mHD95 below 5 mm, and 11 OARs below 10 mm. The overall mHD95 of AbsegNet, SwinUNETR, DeepLabV3+, Attention-UNet, UNet, and 3D-UNet were 8.92 mm, 12.50 mm, 16.44 mm, 13.12 mm, 13.00 mm, and 23.22 mm, respectively. Furthermore, compared with SwinUNETR (the best one among all previous algorithms), a 29% reduction of mHD95 was observed for AbsegNet. Visualization of one randomly selected CT scan from cohort 1 is illustrated in Fig. 2.

Table S5 contains a summary of previously reported delineation results for multiple abdominal OARs, and comparable accuracy was observed in AbsegNet.

#### Performance comparison in external testing cohorts

Table 3 summarizes the DSC and HD95 in external cohort 2. When it comes to DSC, AbsegNet was prone to produce more accurate contours than other algorithms, showing the best performance at 13 OARs. The overall mDSC of AbsegNet, SwinUNETR, DeepLabV3+, Attention-UNet, UNet, and 3D-UNet were 85.65%, 79.27%, 78.51%, 79.73%, 81.14%, and 77.58%, respectively. In terms of HD95, the accuracy of 13 OARs produced by AbsegNet outperformed five previous algorithms. Furthermore, the advantage of AbsegNet over other algorithms is more obvious when evaluating the mHD95 across all OARs, with the value decreasing nearly 50% compared to SwinUNETR, DeeplabV3+, UNet, and 3D-UNet. Similarly, Visualization of one randomly selected patient from cohort 2 is shown in Fig. 2.

Consistent results were obtained from the public CT images (Table S6).



AbsegNet presented the best accuracy in 13 out of 15 OARs. Overall mDSC of AbsegNet, SwinUNETR, DeepLabV3+, Attention-UNet, UNet, and 3D-UNet was 88.04%, 79.87%, 83.03%, 81.75%, 82.27%, and 82.43%, respectively. An improvement of 5.01 % in mDSC was observed when comparing AbsegNet with DeepLabV3+, the best among the five previous methods. For HD95, AbsegNet showed the most accuracy in 66.7% (10/15) of OARs. The whole mHD95 of AbsegNet, SwinUNETR, DeepLabV3+, Attention-UNet, UNet, and 3D-UNet were 12.40 mm, 23.77 mm, 22.06 mm, 15.72 mm, 24.28 mm, and 25.46 mm, respectively.

#### Performance of AbsegNet in all cohorts

We examined the performance of AbsegNet in two datasets, including cohort 1, cohort 2, and cohort 3. The mDSC and mHD95 of each OAR for all patients are shown in Fig. S4. The mDSC of eight OARs exceeded 90%, while only two OARs had a mDSC less than 80% (adrenals and duodenum). For HD95, nine out of 16 OARs had a mHD95 less than 10 mm, while only three OARs were more than 15 mm (colon, duodenum, and femur\_L).

#### Performance of AbsegNet in different types of organs

The 16 OARs were divided into three groups according to morphological structures and size: solid organs (liver, spleen, kidneys, and pancreas), gastrointestinal organs (esophagus, stomach, duodenum, colon, small bowel, rectum, and bladder), bone tissues (femurs), and small organs (gallbladder and adrenals). Then, the performance of AbsegNet was examined in these four groups. The exact results are listed in Table S7. Of these cohorts, all mDSC were more than 90% for solid organs and bone tissues;

all mDSC exceeded 80% in gastrointestinal organs; and all mDSC were above 75% in small organs.

### Clinical assessment of contours produced by AbsegNet

When using our 4-grading criteria to assess contour accuracy, all patients with AbsegNet-produced contours for kidneys, bladder, and femur\_L were deemed satisfactory by the expert, with no or minor revisions in cohort 4 (Table 4). However, 25% ( $n = 10$ ) and 10% ( $n = 4$ ) of patients with AbsegNet-produced duodenum were considered to have moderate and major revisions, respectively (Table 4). Similarly, in cohort 5, AbsegNet-generated contours for liver and spleen in all patients were not required to be revised; and over 87.5% of patients with contours of the stomach, esophagus, adrenals, and rectum were considered satisfactory, with no or minor revisions (Table 4). Only two (6.2%), five (15.6%), and four (12.5%) patients with autosegmentations of the gallbladder, colon, and small bowel, respectively, required major revisions (Table 4). A visual display of AbsegNet-generated OARs revised by experts is presented in Fig. S5.

In addition, two experts together subjectively evaluated each OARs predicted by AbsegNet from the 25 patients. All patients with AbsegNet-generated predictions for liver, spleen, and kidneys were considered as completely acceptable (Table S8). Only two (8.0%), two (8.0%), five (20.0%), three (12.0%), and two (8.0%) patients with autosegmentations of stomach, adrenals, duodenum, colon, and small bowel, respectively, were considered as clinically unacceptable (Table S8).

We recorded the time spent by expert A to delineate 16 OARs in each of ten

patients. With the assistance of AbsegNet, the expert spent on average  $12.04 \pm 2.93$  min to delineate one patient. However, without assistance from AbsegNet, the delineation time was significantly increased to an average of  $39.73 \pm 3.38$  min per case ( $p < 0.001$ ) (Fig. S6).

## Discussion

In this study, we aimed to develop a novel DL model to segment OARs for abdominal radiotherapy accurately and robustly. At first, large-scale and multi-center CT scans were collected, and a training cohort with high-quality annotations was utilized to train AbsegNet. Then, a comprehensive evaluation of model performance was conducted across five different institutions, including clinical assessment. These results showed that AbsegNet achieved state-of-the-art performance in seen subjects and also generalized well to unseen subjects. Compared with five previous sophisticated DL methods, the accuracy of the majority of OARs produced by AbsegNet was higher. And considering all OARs as a whole, AbsegNet demonstrated the best performance. When experts evaluated these autosegmentations, most OARs were considered as satisfactory with no or minor revisions, suggesting that AbsegNet-generated contours were clinically acceptable.

In recent years, many fully automatic segmentation models have been proposed to delineate OARs in the abdomen<sup>[8, 9, 32, 33]</sup>. Nevertheless, robust segmentation of OARs remains a challenge in real-world clinical scenarios. Most existing datasets using abdominal organ segmentation vary in size (from dozens to hundreds) and number of annotations (single or several)<sup>[10]</sup>. For instance, the BTCV provided only

50 CT scans covering 13 organs<sup>[34]</sup>. Although the AbdomenCT-1K offered more than ten thousand CT scans, only four organs were included<sup>[10]</sup>. Besides, some datasets only included a single disease, such as all patients with gastric cancer<sup>[33]</sup>. Considering underrepresentation problems in those datasets, it is not easy to use those models to directly segment abdominal OARs for radiotherapy. Compared with previous datasets<sup>[1, 9, 14, 35, 36]</sup>, the datasets included in this study were more advantageous from the following aspects, (1) Large-scale: datasets contained over 500 CT scans covering nearly all abdominal and pelvic OARs; (2) Diverse and clinical relevant: the datasets used in this study were collected from real-world clinical settings, for example, these patients were diagnosed with various abdominal primary or metastatic tumors, were scanned by different scanners at different medical centers, and both contrast and non-contrast CT images were included; (3) High-quality annotations: principles for OAR delineation in this study were in line with recommendations of the Radiation Therapy Oncology Group (RTOG)<sup>[16, 17]</sup>, with small intraobserver variability, making them more suitable for radiation treatment planning. Parts of the carefully-annotated dataset of 16 organ annotations (150 volumes from the development cohort and 20 volumes from LiTS2017) will be released to boost this research task. In short, our datasets are more promising to develop a robust segmentation model to clinical application.

To obtain accurate and robust segmentations, we propose combining data augmentation and knowledge distillation for deep network training. First, a series of data augmentation strategies could simulate more challenging scenarios to boost the

robustness of networks. Furthermore, we used knowledge distillation to minimize the difference between the teacher and the student. In general, the input of the student is more challenging than the teacher, and the output of the teacher is more accurate than the student. We encouraged student's outputs to be consistent with their teachers, reminding them to pay more attention to the common anatomical context rather than the variance of the intensity distribution. Based on these approaches, the AbsegNet can learn from one data center and generalize well to many unseen centers. Different from those previous works which focused on improving networks' performance on a single data center<sup>[23-25]</sup>, the AbsegNet considered the domain shift between different centers, and utilized data augmentation and knowledge distillation to boost the models' generalization. In addition, the proposed training strategy can improve performance by a large margin compared with the standard 3D-UNet model<sup>[23]</sup>, further suggesting the efficiency and effectiveness of the proposed approach.

After finishing construction and training, performance of AbsegNet was first validated in cohort 1. It was demonstrated that good segmentation results were acquired for most OARs, with the exception of the duodenum and adrenals. Volumetric overlaps between AbsegNet-generated contours and ground truth contours were around 95% for solid organs (liver, spleen, and kidneys). And for all OARs, overall mDSC reached 87% and overall mHD95 was lower than 10 mm. By contrast to previous DL methods, AbsegNet was prone to generate more accurate contours for most OARs, particularly for gastrointestinal organs that are difficult to delineate due to their huge and variable spatial information.

AbsegNet was tested on two completely unseen data acquired from two hospitals. In such heterogenous data, AbsegNet still acquired comparable segmentation results, with the mDSC approaching that which was obtained in the internal cohort 1 (overall mDSC (cohorts 1, 2 and 3): 86.73% versus 85.65% versus 88.04%), and outperformed five existing methods. On the contrary, the performance of previous algorithms was unstable, affected by data differences. As observed, the mDSC of Attention-UNet in cohort 1 was 84.67%. However, mDSC considerably reduced to 79.63% in cohort 3. Moreover, compared with previous methods, AbsegNet achieved smaller standard deviations for DSC and HD95 overall, suggesting AbsegNet is less affected by individual differences and confirming its robustness for delineating abdominal OARs.

Comparable accuracy was obtained in AbsegNet by contrasting with historical results for multiple abdominal OARs delineation. It should be notice that our results were acquired on heterogeneous CT scans and patients, more approaching to clinical setting, whereas those<sup>[1, 9]</sup> were acquired on more homogeneous data. As indicated, different scanners and CT phases on patients with heterogeneous lesions could lead to obvious variances in organ appearances, resulting in a degradation of model performance. Hence, results showed the powerful generalizability of AbsegNet.

To further confirm the accuracy of AbsegNet clinically, two independent experts were invited to review AbsegNet-generated OARs. AbsegNet-produced contours for solid organs (liver, kidneys, and spleen) of all patients did not need to be modified. Only 15.6%, 12.5%, 10%, and 6.2% of patients with autosegmentations of colon,

small bowel, duodenum, and gallbladder, respectively, required major revisions. For other OARs, only a small portion of patients needed minor or moderate revisions. In subjective evaluation by the two experts together, only several patients with OARs, such as duodenum, adrenals, and colon, were unacceptable, whereas most patients with most OARs were considered as totally acceptable or acceptable. Moreover, we showed that, with the aid of the model, delineating efficiency was substantially improved, saving the delineation time by as much as over 65%. These results indicated that most OARs produced by AbsegNet were clinically applicable, and could be used for radiotherapy.

This study had several limitations. First, in order to improve the consistency of OAR delineation, only one experienced expert was invited to delineate these contours, which introduced potential subjective variation. Second, although most OARs obtained high accuracy, the performance of AbsegNet on duodenum was not satisfactory. The possible reason might be that the duodenum has the most complex anatomical structure, made up of four parts: the superior, the descending, the horizontal, and the ascending. Hence, the organ volume and location can vary dramatically, posing a great challenge for DL models. Similar results were also observed in the study of Gibson et al<sup>[14]</sup>, with a mDSC of 63%. There is much room to improve the segmentation accuracy for this organ. In the future study, we are going to experiment with different loss functions to optimize our model, hoping to obtain better results on lower performing structures. Third, when performing clinical evaluation, two experts were required to review complementary OARs (seven OARs

for expert A, and nine OARs for expert B) in two cohorts rather than each checking all OARs in each cohort due to the heavy workload and time requirements. This might have had a certain impact on the evaluative accuracy of AbsegNet. Alternatively, pelvic bone delineation was not involved in this study, though it is also an important OAR for abdominal radiotherapy<sup>[37]</sup>. Future research proposes to incorporate it.

## Conclusions

In summary, we proposed a novel fully automatic DL model to delineate whole abdominal OARs. Despite heterogeneous CT scans and individual differences, our findings showed that the vast majority of OARs produced by AbsegNet were accurate and robust. It is clinically applicable and helpful to facilitate radiation treatment planning and workflow with ongoing efforts.

## References

- [1] Chen X, Sun S, Bai N, et al. A deep learning-based auto-segmentation system for organs-at-risk on whole-body computed tomography images for radiation therapy. *Radiother Oncol* 2021; 160: 175-184.
- [2] Joskowicz L, Cohen D, Caplan N, et al. Inter-observer variability of manual contour delineation of structures in CT. *Eur Radiol* 2019; 29: 1391-1399.
- [3] Peng YL, Chen L, Shen GZ, et al. Interobserver variations in the delineation of



target volumes and organs at risk and their impact on dose distribution in intensity-modulated radiation therapy for nasopharyngeal carcinoma. *Oral Oncol* 2018; 82: 1-7.

[4] Spoelstra FO, Senan S, Le Pécoux C, et al. Variations in target volume definition for postoperative radiotherapy in stage III non-small-cell lung cancer: analysis of an international contouring study. *Int J Radiat Oncol Biol Phys* 2010; 76: 1106-1113.

[5] Vinod SK, Jameson MG, Min M, et al. Uncertainties in volume delineation in radiation oncology: A systematic review and recommendations for future studies. *Radiother Oncol* 2016; 121: 169-179.

[6] Litjens G, Kooi T, Bejnordi BE, et al. A survey on deep learning in medical image analysis. *Med Image Anal* 2017; 42: 60-88.

[7] Zhou SK, Greenspan H, Davatzikos C, et al. A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises. *Proceeding of the IEEE* 2021; 109: 820-838.

[8] Humpire-Mamani GE, Bukala J, Scholten ET, et al. Fully Automatic Volume Measurement of the Spleen at CT Using Deep Learning. *Radiol Artif Intell* 2020; 2: e190102.

[9] Weston AD, Korfiatis P, Philbrick KA, et al. Complete abdomen and pelvis segmentation using U-net variant architecture. *Med Phys* 2020; 47: 5609-5618.

[10] Ma J, Zhang Y, Gu S, et al. AbdomenCT-1K: Is Abdominal Organ Segmentation A Solved Problem. *IEEE Trans Pattern Anal Mach Intell* 2021, Pp.

[11] Cardenas CE, Beadle BM, Garden AS, et al. Generating High-Quality Lymph

Node Clinical Target Volumes for Head and Neck Cancer Radiation Therapy Using a Fully Automated Deep Learning-Based Approach. *Int J Radiat Oncol Biol Phys* 2021; 109: 801-812.

[12] van Dijk LV, Van den Bosch L, Aljabar P, et al. Improving automatic delineation for head and neck organs at risk by Deep Learning Contouring. *Radiother Oncol* 2020; 142: 115-123.

[13] Rister B, Yi D, Shivakumar K, et al. CT-ORG, a new dataset for multiple organ segmentation in computed tomography. *Sci Data* 2020; 7: 381.

[14] Gibson E, Giganti F, Hu Y, et al. Automatic Multi-Organ Segmentation on Abdominal CT With Dense V-Networks. *IEEE Trans Med Imaging* 2018; 37: 1822-1834.

[15] Antonelli M, Reinke A, Bakas S, et al. The medical segmentation decathlon. *Nat Commun* 2022; 13: 1-13.

[16] Gay HA, Barthold HJ, O'Meara E, et al. Pelvic normal tissue contouring guidelines for radiation therapy: a Radiation Therapy Oncology Group consensus panel atlas. *Int J Radiat Oncol Biol Phys* 2012; 83: e353-362.

[17] Jabbour SK, Hashem SA, Bosch W, et al. Upper abdominal normal organ contouring guidelines and atlas: a Radiation Therapy Oncology Group consensus. *Pract Radiat Oncol* 2014; 4: 82-89.

[18] Pukala J, Johnson PB, Shah AP, et al. Benchmarking of five commercial deformable image registration algorithms for head and neck patients. *J Appl Clin Med Phys* 2016; 17: 25-40.

- [19] Nakajima Y, Kadoya N, Kanai T, et al. Evaluation of the effect of user-guided deformable image registration of thoracic images on registration accuracy among users. *Med Dosim* 2020; 45: 206-212.
- [20] Zhang L, Wang X, Yang D, et al. Generalizing deep learning for medical image segmentation to unseen domains via deep stacked transformation. *IEEE Trans Med Imaging* 2020; 39: 2531-2540.
- [21] Gou J, Yu B, Maybank SJ, et al. Knowledge distillation: A survey. *Int J Comput Vis* 2021; 129: 1789-1819.
- [22] Tarvainen A, Valpola H. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in Neural Information Processing Systems (NIPS)* 2017, 30.
- [23] Ballestar LM, Vilaplana V. MRI brain tumor segmentation and uncertainty estimation using 3D-UNet architectures. *International MICCAI Brainlesion Workshop: Springer*, 2020:376-390.
- [24] Chen LC, Zhu Y, Papandreou G, et al. Encoder-decoder with atrous separable convolution for semantic image segmentation. *Proceedings of the European conference on computer vision (ECCV)*, 2018:801-818.
- [25] Oktay O, Schlemper J, Folgoc LL, et al. Attention u-net: Learning where to look for the pancreas. *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [26] Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image Segmentation. *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015.

- [27] Tang Y, Yang D, Li W, et al. Self-supervised pre-training of swin transformers for 3d medical image analysis. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022:20730-20740.
- [28] Liang S, Tang F, Huang X, et al. Deep-learning-based detection and segmentation of organs at risk in nasopharyngeal carcinoma computed tomographic images for radiotherapy planning. *Eur Radiol* 2019; 29: 1961-1967.
- [29] Liu Z, Liu X, Guan H, et al. Development and validation of a deep learning algorithm for auto-delineation of clinical target volume and organs at risk in cervical cancer radiotherapy. *Radiother Oncol* 2020; 153: 172-179.
- [30] Vrtovec T, Močnik D, Strojanić P, et al. Auto-segmentation of organs at risk for head and neck radiotherapy planning: From atlas-based to deep learning methods. *Med Phys* 2020; 47: e929-e950.
- [31] Tang H, Chen X, Liu Y, et al. Clinically applicable deep learning framework for organs at risk delineation in CT images. *Nat Mach Intell* 2019; 1: 480-491.
- [32] Qayyum A, Lalonde A, Meriaudeau F. Automatic segmentation of tumors and affected organs in the abdomen using a 3D hybrid model for computed tomography imaging. *Comput Biol Med* 2020; 127: 104097.
- [33] Roth HR, Oda H, Hayashi Y, et al. Hierarchical 3D fully convolutional networks for multi-organ segmentation. *Computer Vision and Pattern Recognition (CVPR)* 2017.
- [34] Landman B, Xu Z, Igelsias J, et al. Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge. *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial*

Vault—Workshop Challenge, 2015:12.

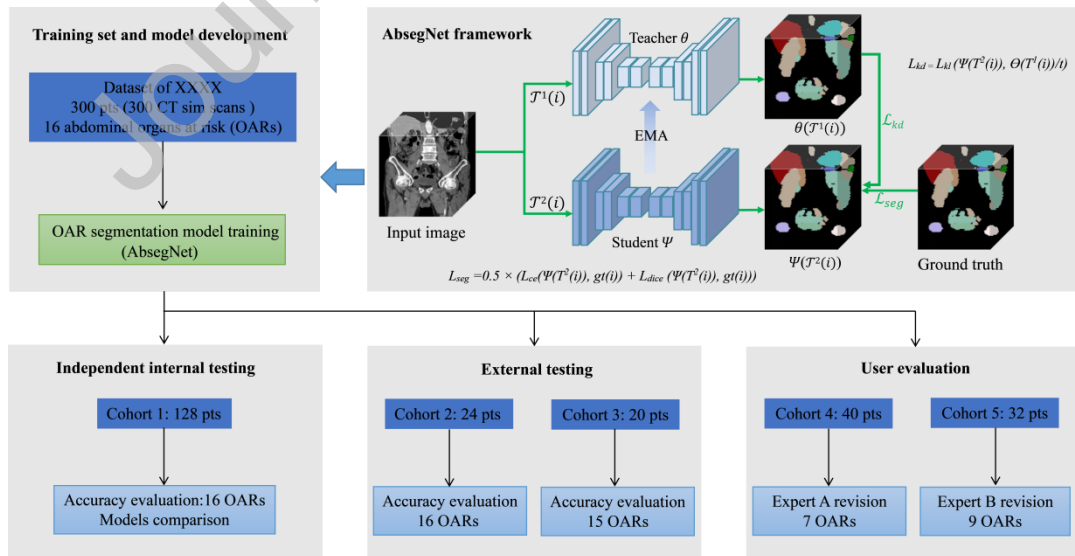
[35] Kavur AE, Gezer NS, Barış M, et al. CHAOS Challenge - combined (CT-MR)

healthy abdominal organ segmentation. Med Image Anal 2021; 69: 101950.

[36] Heller N, Sathianathen N, Kalapara A, et al. The kits19 challenge data: 300 kidney tumor cases with clinical context, ct semantic segmentations, and surgical outcomes 2019.

[37] Albuquerque K, Giangreco D, Morrison C, et al. Radiation-related predictors of hematologic toxicity after concurrent chemoradiation for cervical cancer and implications for bone marrow-sparing pelvic IMRT. Int J Radiat Oncol Biol Phys 2011; 79: 1043-1047.

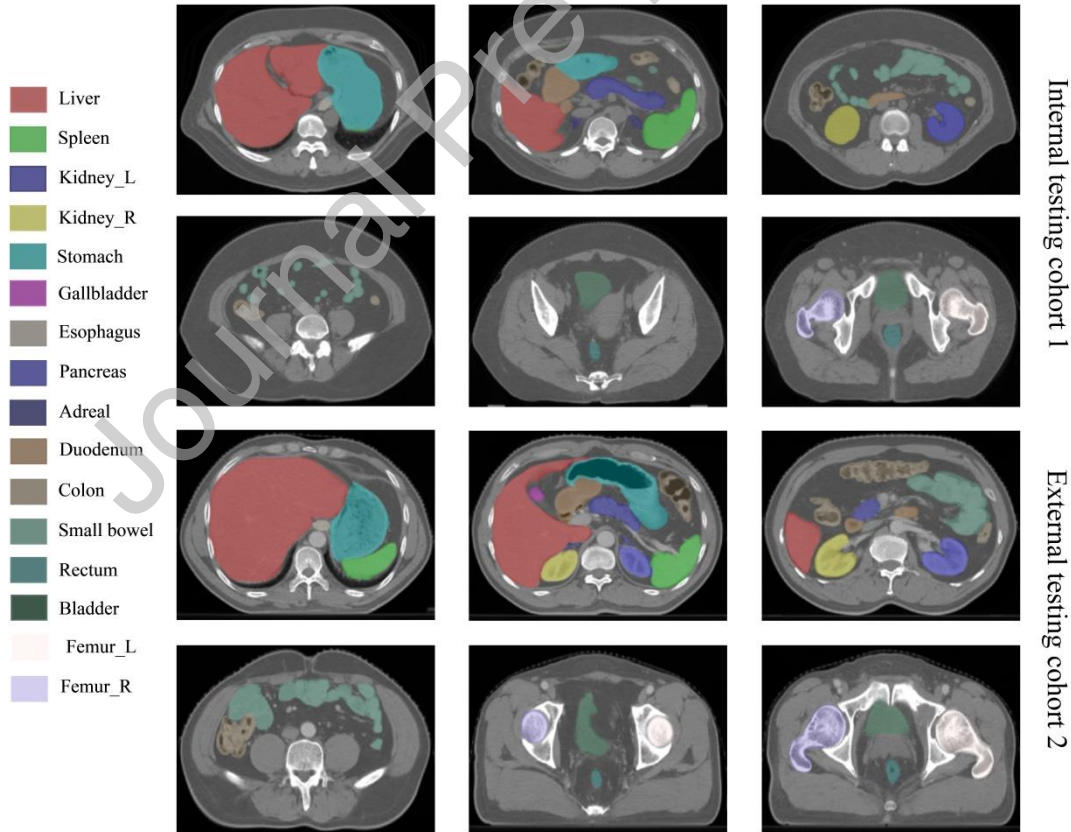
## Figure captions



**Fig. 1.** The flow chart of this study. The top right of the figure is the AbsegNet

framework, which consists of a teacher-student model. In the training stage, the teacher and the student take images with different augmentation strategies (the teacher with weak augmentation and the student with strong augmentation as inputs) and then employ the teacher's output to teach the student to be more robust. The student network updates parameter by minimizing loss functions, and the teacher parameter is updated as an EMA of the student's parameter. In the testing stage, the teacher model was used to produce the final segmentation results. The mathematical definitions are presented in the section of segmentation network construction.

XXXX = XXXX; EMA = exponential moving average; Pts = patients.



**Fig. 2.** Visualization of two randomly selected patients from internal testing cohort 1 and external testing cohort 2, respectively.

Journal Pre-proof

**Table 1** Baseline characteristics and scanning parameters

Characteristics	Dataset 1 (n = 428)		Dataset 2 for external testing (n = 44)		Dataset 3 for clinical evaluation (n = 72)	
	Training cohort	Internal testing cohort	Cohort 2	Cohort 3	Cohort 4	Cohort 5
	XXXX (n = 300, %)	XXXX (n = 128, %)	XXXX (n = 24, %)	Public <sup>#</sup> (n = 20)	XXXX (n = 40, %)	XXXX (n = 32, %)
Gender						
Male	182 (60.7)	76 (59.4)	10 (41.7)	NA	0	0
Female	118 (39.3)	52 (40.6)	14 (58.3)	NA	40 (100.0)	32 (100.0)
Age (median)	47 (17 - 75)	47 (20 - 72)	49 (36 - 72)	NA	55 (35 - 62)	53 (46 - 70)
Tumor site						
Rectal cancer	143 (47.7)	60 (46.9)	5 (20.8)	NA	0	0
Prostate cancer	39 (13.0)	18 (14.1)	7 (29.2)	NA	0	0
Gynecological*	34 (11.3)	11 (8.6)	12 (50.0)	NA	40 (100.0)	32 (100.0)
Bladder cancer	11 (3.6)	5 (3.9)	0	NA	0	0
Metastatic tumor	44 (14.7)	20 (15.6)	0	NA	0	0
Others <sup>¶</sup>	29 (9.7)	14 (10.9)	0	NA	0	0
OAR types annotated	16	16	16	15 <sup>¶</sup>	16	16
Scanning parameters						
Total slice (median)	200 (123 - 368)	201 (145 - 436)	168 (118 - 248)	501 (276 - 842)	176 (152 - 201)	172 (103 - 260)
Thickness (median, mm)	3.0 (0.60 - 0.98)	3.0 (2.5 - 3.0)	3.0 (3.0 - 3.0)	1.0 (0.70 - 1.5)	3.0 (3.0 - 3.0)	3.0 (3.0 - 5.0)
In plane spacing	0.98 (0.60 - 0.98)	0.98 (0.78 - 1.27)	0.98 (0.95 - 0.98)	0.74 (0.60 - 0.90)	0.96 (0.90 - 1.04)	1.04 (0.91 - 1.17)
Manufacture	Siemens	Siemens	Philips	NA	Philips	Philips

Abbreviation: XXXX = XXXX; XXXX = XXXX; XXXX = XXXX; XXXX = XXXX; OAR = organ at risk; NA = not applicable.

\* Included cervical cancer and endometrial cancer; <sup>¶</sup> Included liver cancer, pancreatic cancer, kidney cancer, sarcoma, and testicular cancer.

<sup>#</sup> These patients were collected from a public dataset, named LiTS2017 (ref.15). <sup>¶</sup> Considering the liver was manually previously labeled in LiTS2017, it was excluded and the remaining 15 OARs



were labeled in this study.

**Table 2** Accuracy comparison in the cohort 1 (n = 128)

Variable		DSC (%)					HD95 (mm)					
OAR	Abseg	Swin	Deep	Att	UNet	3D-U	Abseg	Swin	Deep	Att	UNet	3D-U
	Net	UNE	LabV			Net	Net	UNE	LabV			Net
		TR	3+					TR	3+			
Liver	<b>96.40</b>	95.94	95.00	95.61 ±	95.49	95.25	<b>4.26 ±</b>	4.63 ±	6.27 ±	5.04 ±	5.55 ±	5.58 ±
	± <b>1.13</b>	± 1.58	± 2.14	2.06***	± 1.44	± 1.79	<b>5.66</b>	8.27	9.60	4.22	4.20	4.46
Spleen	<b>95.13</b>	94.17	93.81	94.51 ±	94.46	93.77	<b>3.38 ±</b>	4.07 ±	4.31 ±	3.61 ±	6.79 ±	3.98 ±
	± <b>5.32</b>	± 6.31	± 3.66	5.39***	± 3.45	± 4.26	<b>7.71</b>	7.25	5.92	5.46	19.15	5.21
Kidney_L	<b>95.60</b>	94.36	94.31	94.52 ±	94.86	94.16	<b>2.54 ±</b>	8.79 ±	3.18 ±	4.09 ±	3.03 ±	3.78 ±
	± <b>1.07</b>	± 2.93	± 2.68	6.29*	± 2.93	± 3.08	<b>0.72</b>	23.72	1.84	9.63	1.82	4.01
Kidney_R	<b>95.60</b>	94.29	94.11	94.72 ±	94.97	94.53	<b>2.41 ±</b>	7.97 ±	3.59 ±	3.07 ±	3.63 ±	3.26 ±
	± <b>1.31</b>	± 2.76	± 2.98	3.60**	± 2.30	± 2.32	<b>0.83</b>	23.01	3.77	2.43	4.77	3.26
Stomach	<b>91.46</b>	89.11	88.08	88.95 ±	89.41	86.90	<b>9.87 ±</b>	13.48	18.21	13.59	14.30	18.24
	± <b>4.08</b>	± 7.53	± 8.00	7.38***	± 5.95	± 9.38	<b>12.46</b>	± 20.1	±	±	±	±
									42.38	19.91	14.85	38.04
Gallbladder	<b>80.74</b>	71.26	74.42	74.56 ±	76.02	71.79	<b>6.93 ±</b>	13.37	10.86	13.46	10.22	10.95
	±	±	±	20.29**	±	±	<b>9.56</b>	±	±	±	±	±
Intestine	<b>15.17</b>	25.41	17.79	*	17.09	21.59		19.34	12.38	18.53	13.32	15.33
Esophagus	<b>81.36</b>	78.81	73.24	77.54 ±	76.96	73.00	<b>5.19 ±</b>	5.26 ±	6.15 ±	31.91	5.64 ±	7.28 ±
	± <b>6.25</b>	± 6.33	±	8.79***	± 8.41	±	<b>4.46</b>	2.89	3.97	±	3.16	4.84
			11.91			12.51				91.13		
Pancreas	<b>82.71</b>	80.46	77.64	79.82 ±	79.82	76.90	<b>7.64 ±</b>	8.82 ±	10.56	8.62 ±	9.99 ±	12.08
	± <b>7.54</b>	± 7.92	± 8.33	8.43***	± 8.00	± 9.63	<b>7.25</b>	9.66	± 8.79	7.72	8.20	±
												22.83
Adrenals	<b>71.34</b>	65.16	49.30	68.13 ±	67.58	67.71	<b>6.58 ±</b>	22.4 ±	65.22	8.18 ±	8.66 ±	8.39 ±
	±	±	±	12.99**	±	±	<b>5.28</b>	17.49	±	9.59	7.06	8.66
	<b>11.63</b>	16.36	10.72	*	12.06	12.23		***	10.81			
Duodenum	67.77	<b>83.46</b>	61.21	65.51 ±	64.64	62.45	25.85	<b>16.78</b>	23.14	24.31	21.61	23.65
	±	± <b>8.99</b>	±	17.04*	±	±	±	±	±	±	±	±
	16.28		17.15		17.77	17.65	35.60	<b>17.34</b>	16.43	18.97	16.50	17.64
Colon	<b>85.74</b>	84.42	80.03	82.91 ±	81.82	80.45	14.14	<b>9.22 ±</b>	19.34	17.11	25.65	19.70
	± <b>9.51</b>	± 7.71	±	9.36***	± 9.39	±	±	<b>8.55</b>	±	±	±	±
			10.94			10.21	15.44		16.98	16.49	36.03	17.20
Small Bowel	<b>86.42</b>	68.54	82.21	84.56 ±	82.70	82.00	8.49 ±	<b>7.28 ±</b>	11.44	9.66 ±	11.48	11.10
	± <b>8.36</b>	±	± 8.33	8.09***	± 8.66	± 9.36	9.43	<b>6.88</b>	± 8.70	9.07	± 9.04	± 9.43
1		11.88										

Rectum	80.06 ± 12.43	78.08 ± 12.24	78.19 ± 11.59	78.85 ± 12.35*	<b>80.51</b> ± <b>10.12</b>	76.92 ± 11.02	14.06 ± 12.36	14.12 ± 11.71	15.28 ± 15.79	14.40 ± 11.75	<b>13.29</b> ± <b>10.46</b>	17.18 ± 29.77
Bladder	<b>93.14</b> ± <b>7.51</b>	91.1 ± 9.85	90.33 ± 12.45	91.22 ± 11.15**	91.46 ± 11.82	90.48 ± 11.06	<b>4.56</b> ± <b>7.80</b>	14.87 ± 53.42	25.97 ± 73.62	7.00 ± 13.48	19.96 ± 62.55	16.69 ± 53.41
Femur_L	<b>91.87</b> ± <b>4.03</b>	90.46 ± 9.36	90.84 ± 4.22	91.59 ± 4.49	91.86 ± 4.59	88.96 ± 5.60	<b>15.50</b> ± <b>56.57</b>	23.6 ± 85.71	25.06 ± 91.52	17.45 ± 70.20	21.28 ± 81.89	132.8 ± 171.6
Femur_R	<b>92.40</b> ± <b>4.56</b>	90.56 ± 9.2	91.00 ± 5.01	91.72 ± 4.37	92.02 ± 4.48	89.91 ± 5.70	<b>11.32</b> ± <b>30.85</b>	25.32 ± 88.26	25.38 ± 92.43	28.50 ± 95.49	25.99 ± 89.30	88.89 ± 150.7
Average	<b>86.73</b>	84.39	82.11	84.67**	84.66	82.82	<b>8.92</b>	12.50	16.44	13.12	13.00	23.22

Abbreviation: DSC = Dice similarity coefficient; HD95 = 95th percentile Hausdorff distance; Att = attention-UNet; OAR = organ at risk.

Data were denoted as mean ± standard deviation. Bold numbers represented the best results.

P - values were obtained by comparing our method with the best one among the five previous methods according to the overall average DSC; \*,  $P < 0.05$ ; \*\*,  $P < 0.01$ ; \*\*\*,  $P < 0.001$ . # Four patients underwent gallbladder resection, so the number of gallbladders were 124.

**Table 3** Accuracy comparison in the cohort 2 (n = 24)

Variable	DSC (%)						HD95 (mm)					
	Abseg	Swin	Deep	Att	UNet	3D-U	Abseg	Swin	Deep	Att	UNet	3D-U
	Net	UNE	LabV			Net	Net	UNE	LabV			Net
	TR	3+	3+				TR	3+	3+			
Liver	<b>96.75</b> ± <b>1.21</b>	90.52 ± 15.86	92.94 ± 5.93	92.71 ± 9.72	91.79 ±	89.53 ±	<b>4.14</b> ± <b>3.93</b>	15.32 ±	10.75 ±	11.16 ±	12.14 ±	15.86 ±
Spleen	<b>91.50</b> ±	86.81 ±	89.78 ±	88.81 ±	90.04 ±	90.13 ± 7.52	<b>12.92</b> ±	33.54 ±	13.45 ±	13.98 ±	19.64 ±	13.41 ±
Kidney_L	<b>94.91</b> ± <b>2.14</b>	89.81 ±	91.76 ±	92.94 ± 6.69	92.30 ±	91.91 ± 7.26	5.60 ± 12.44	32.27 ±	<b>4.02</b> ± <b>3.66</b>	3.92 ± 3.17	4.76 ± 5.49	10.83 ±
Kidney_R	<b>95.13</b> ± <b>1.15</b>	91.23 ±	92.78 ± 8.40	93.24 ± 5.11	93.43 ± 7.21	93.36 ± 4.13	<b>2.65</b> ± <b>0.76</b>	31.25 ±	3.61 ± 2.22	6.67 ± 15.44	5.41 ± 10.90	4.67 ± 6.07
Stomach	<b>86.45</b> ±	67.62 ±	75.47 ±	67.11 ±	72.81 ±	61.32 ±	<b>12.43</b> ±	31.54 ±	22.64 ±	27.98 ±	25.31 ±	32.96 ±
	<b>14.42</b>	27.35	22.71	27.16	25.16**	29.86	<b>14.03</b>	35.77	18.22	22.66**	16.09	28.19

Gallb	<b>87.85</b>	74.23	79.13	78.03	77.43	72.51	<b>6.93</b> ±	23.42	10.33	10.28	10.32	11.02
ladde	± <b>6.36</b>	±	±	±	±	±	<b>14.15</b>	±	±	±	±	±
r		28.67	16.45	20.53	20.47**	23.77		76.39	16.78	16.39	14.95	16.73
Esop	<b>79.90</b>	78.23	71.24	74.50	75.58	69.89	<b>5.35</b> ±	16.1	11.54	45.34	10.68	9.36 ±
hagus	± <b>4.56</b>	±	±	± 8.97	±	±	<b>3.01</b>	±	±	±	±	6.70
		5.82	12.05		10.00**	13.44		49.74	12.89	90.36*	13.17	
Pancr	<b>81.32</b>	75.51	74.12	73.40	76.05	72.47	<b>7.91</b> ±	17.58	12.93	15.05	13.48	19.05
eas	± <b>9.68</b>	±	±	±	±	±	<b>7.36</b>	±	± 8.15	±	±	±
		13.43	14.14	18.10	14.89*	16.54		19.11		13.85**	10.86	18.91
Adre	<b>70.87</b>	59.32	46.11	62.93	68.00	66.26	<b>6.19</b> ±	30.31	65.06	7.42 ±	6.87 ±	6.79 ±
nals	±	±	±	±	±	±	<b>3.56</b>	±	± 7.98	4.52	3.59	4.06
	<b>11.66</b>	19.3	13.05	17.22	11.96*	11.92		17.8				
Duod	63.40	<b>81.84</b>	54.75	58.03	60.80	56.88	23.73	28.76	26.46	34.23	<b>23.02</b>	42.34
enum	±	±	±	±	±	±	±	±	±	±	±	±
	19.62	<b>8.39</b>	20.29	21.19	19.43	20.42	13.52	44.95	18.91	23.70*	<b>15.72</b>	68.51
Colo	<b>85.82</b>	66.79	78.91	79.98	81.86	78.33	<b>12.24</b>	23.24	18.76	21.18	29.87	19.41
n	± <b>5.64</b>	±	± 8.47	± 9.06	±	±	± <b>8.17</b>	±	± 8.70	±	±	±
		18.77			5.54***	14.94		16.16		12.74**	39.42	11.51
										*		
Small	<b>82.41</b>	68.18	63.80	67.87	68.73	63.92	<b>14.33</b>	7.0 ±	20.52	17.63	26.74	19.34
Bowe	± <b>9.50</b>	±	±	±	±	±	±	4.23	± 9.44	± 9.05	±	± 9.92
l		10.31	17.85	16.39	14.88**	17.30	<b>13.73</b>				51.21	
					*							
Rectu	83.02	82.06	81.74	81.58	<b>83.30</b>	79.57	<b>10.53</b>	15.66	11.00	10.91	32.77	11.29
m	± 6.42	± 6.6	± 7.77	±	± <b>6.48</b>	±	± <b>7.69</b>	±	±	±	±	±
				13.24		15.72		19.66	10.08	11.14	105.80	10.62
Bladd	<b>87.74</b>	81.48	84.97	82.71	84.72	85.08	<b>6.51</b> ±	33.36	25.86	10.55	41.72	62.64
er	±	±	±	±	±	±	<b>8.07</b>	±	±	±	±	±
	<b>18.25</b>	26.98	17.67	25.46	19.75	20.98		120.1	79.16	14.56	94.01	144.60
								2				
Femu	<b>93.14</b>	89.09	90.75	91.01	91.68	82.15	<b>13.31</b>	39.62	37.26	22.15	28.69	45.19
r_L	± <b>2.79</b>	±	± 4.11	± 7.30	±	±	±	±	±	±	±	±
		6.45			19.75	16.76	<b>46.50</b>	53.69	54.82	63.90	44.01	45.11
Femu	90.22	85.52	87.76	<b>90.85</b>	89.71	87.95	18.04	48.34	37.69	<b>11.02</b>	30.62	32.06
r_R	±	±	± 7.25	±	± 6.21	±	±	±	±	±	±	±
	11.00	17.93		<b>10.68</b>		15.50	43.91	58.75	52.65	<b>20.71</b>	40.68	48.48
Aver	<b>85.65</b>	79.27	78.51	79.73	81.14**	77.58	<b>10.18</b>	26.71	20.62	16.84**	20.13	22.26
age					*					*		

Abbreviation: DSC = Dice similarity coefficient; HD95 = 95th percentile Hausdorff distance; Att = attention-UNet; OAR = organs at risk.

Data were denoted as mean ± standard deviation. Bold numbers represented the best results.

P - values were obtained by comparing our method with the best one among the five previous methods according to the whole average DSC; \*,  $P < 0.05$ ; \*\*,  $P < 0.01$ ; \*\*\*,  $P < 0.001$ .

**Table 4** Clinical accuracy evaluation for each OAR in cohorts 4 and 5

Cohort 4 (n = 40)	No revision (n, %)	Minor revision (n, %)	Moderate revision (n, %)	Major revision (n, %)
Kidney_L	40 (100.0%)	0 (0.0%)	0 (0.0%)	0 (0.0%)
Kidney_R	40 (100.0%)	0 (0.0%)	0 (0.0%)	0 (0.0%)
Pancreas	32 (80.0%)	6 (15.0%)	2 (5.0%)	0 (0.0%)
Duodenum	7 (17.5%)	19 (47.5%)	10 (25.0%)	4 (10.0%)
Bladder	38 (95.0%)	2 (5.0%)	0 (0.0%)	0 (0.0%)
Femur_L	39 (97.5%)	1 (2.5%)	0 (0.0%)	0 (0.0%)
Femur_R	35 (87.5%)	3 (7.5%)	1 (2.5%)	0 (0.0%)
Cohort 5 (n = 32)				
Liver	32 (100.0%)	0 (0.0%)	0 (0.0%)	0 (0.0%)
Spleen	32 (100.0%)	0 (0.0%)	0 (0.0%)	0 (0.0%)
Stomach	23 (71.9%)	6 (18.8%)	3 (9.3%)	0 (0.0%)
Gallbladder	21 (65.7%)	6 (18.8%)	3 (9.3%)	2 (6.2%)
Esophagus	27 (84.4%)	2 (6.3%)	3 (9.3%)	0 (0.0%)
Adrenals	13 (40.6%)	15 (46.9%)	4 (12.5%)	0 (0.0%)
Colon	16 (50.0%)	7 (21.9%)	4 (12.5%)	5 (15.6%)
Small bowel	23 (71.9%)	4 (12.5%)	1 (3.1%)	4 (12.5%)
Rectum	24 (75.0%)	4 (12.5%)	4 (12.5%)	0 (0.0%)

Abbreviation: OAR = organ at risk.