

Weakly-supervised Convolutional Neural Networks of Renal Tumor Segmentation in Abdominal CTA Images

Guanyu Yang (✉ ggyang1980@qq.com)

Key Laboratory of Computer Network and Information Integration, Southeast University, Ministry of Education, Nanjing, China <https://orcid.org/0000-0003-3704-1722>

Chuanxia Wang

Key Laboratory of Computer Network and Information Integration, Southeast University, Ministry of Education, Nanjing, China

Jian Yang

Beijing Engineering Research Center of Mixed Reality and Advanced Display, School of Optics and Electronics, Beijing Institute of Technology, Beijing 100081, China

Yang Chen

Key Laboratory of Computer Network and Information Integration, Southeast University, Ministry of Education, Nanjing, China

Lijun Tang

Dept. of Radiology, the First Affiliated Hospital of Nanjing Medical University, Nanjing, China

Pengfei Shao

Dept. of Urology, the First Affiliated Hospital of Nanjing Medical University, Nanjing, China

Jean-Louis Dillenseger

Univ Rennes, Inserm, LTSI-UMR1099, Rennes, F-35000, France

Huazhong Shu

Key Laboratory of Computer Network and Information Integration, Southeast University, Ministry of Education, Nanjing, China

Limin Luo

Key Laboratory of Computer Network and Information Integration, Southeast University, Ministry of Education, Nanjing, China

Research article

Keywords: weakly-supervised, renal tumor segmentation, bounding box, convolutional neural network

Posted Date: March 25th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-17221/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Medical Imaging on April 15th, 2020.

See the published version at <https://doi.org/10.1186/s12880-020-00435-w>.

1 Weakly-supervised Convolutional Neural Networks of Renal
2 Tumor Segmentation in Abdominal CTA Images

3 Guanyu Yang^{1,6*}, Chuanxia Wang¹, Jian Yang², Yang Chen^{1,6}, Lijun Tang³, Pengfei Shao⁴, Jean-
4 Louis Dillenseger^{5,6}, Huazhong Shu^{1,6}, Limin Luo^{1,6}

5 1 LIST, Key Laboratory of Computer Network and Information Integration, Southeast University, Ministry of
6 Education, Nanjing, China

7 2 Beijing Engineering Research Center of Mixed Reality and Advanced Display, School of Optics and Electronics,
8 Beijing Institute of Technology, Beijing 100081, China

9 3 Dept. of Radiology, the First Affiliated Hospital of Nanjing Medical University, Nanjing, China

10 4 Dept. of Urology, the First Affiliated Hospital of Nanjing Medical University, Nanjing, China

11 5 Univ Rennes, Inserm, LTSI - UMR1099, Rennes, F-35000, France

12 6 Centre de Recherche en Information Biomédicale Sino-Français (CRIBs)

13 **Email Address:**

14 Guanyu Yang: gyyang1980@qq.com; yang.list@seu.edu.cn

15 Chuanxia Wang: rachelgrey@foxmail.com

16 Jian Yang: jyangbit@163.com

17 Yang Chen: chenyang.list@seu.edu.cn

18 Lijun Tang: lijun.tang@hotmail.com

1 Pengfei Shao: spf8629@163.com

2 Jean-louis Dillenseger: jean-louis.coatriexu@univ-rennes1.fr

3 Huazhong Shu: shu.list@seu.edu.cn

4 Limin Luo: luo.list@seu.edu.cn

5 **Corresponding Author:**

6 * Guanyu Yang; Email address: gyyang1980@qq.com; yang.list@seu.edu.cn

7 **Mailing address:** Laboratory of Image Science and Technology, School of Computer Science and Engineering,
8 Southeast University, Nanjing, People's Republic of China. Tel: +86 25-83794249; Fax: +86 25-83792628.

9 **Abstract**

10 **Background:** Renal cancer is one of the ten most common cancers in human beings. The
11 laparoscopic partial nephrectomy (LPN) is an effective way to treat renal cancer. Localization and
12 delineation of the renal tumor from pre-operative CT Angiography (CTA) is an important step for
13 LPN surgery planning. Recently, with the development of the technique of deep learning, deep
14 neural networks can be trained to provide accurate pixel-wise renal tumor segmentation in CTA
15 images. However, constructing the training dataset with a large amount of pixel-wise annotations is
16 a time-consuming task for the radiologists. Therefore, weakly-supervised approaches attract more
17 interest in research.

18 **Methods:** In this paper, we proposed a novel weakly-supervised convolutional neural network
19 (CNN) for renal tumor segmentation. A three-stage framework was introduced to train the CNN
20 with the weak annotations of renal tumors, i.e. the bounding boxes of renal tumors. The framework

1 includes pseudo masks generation, group and weighted training phases. Clinical abdominal CT
2 angiographic images of 200 patients were applied to perform the evaluation.

3 **Results:** Extensive experimental results show that the proposed method achieves a higher dice
4 coefficient (DSC) of 0.826 than the other two existing weakly-supervised deep neural networks.
5 Furthermore, the segmentation performance is close to the fully supervised deep CNN.

6 **Conclusions:** The proposed strategy improves not only the efficiency of network training but also
7 the precision of the segmentation.

8 **Keywords:** weakly-supervised; renal tumor segmentation; bounding box; convolutional neural
9 network.

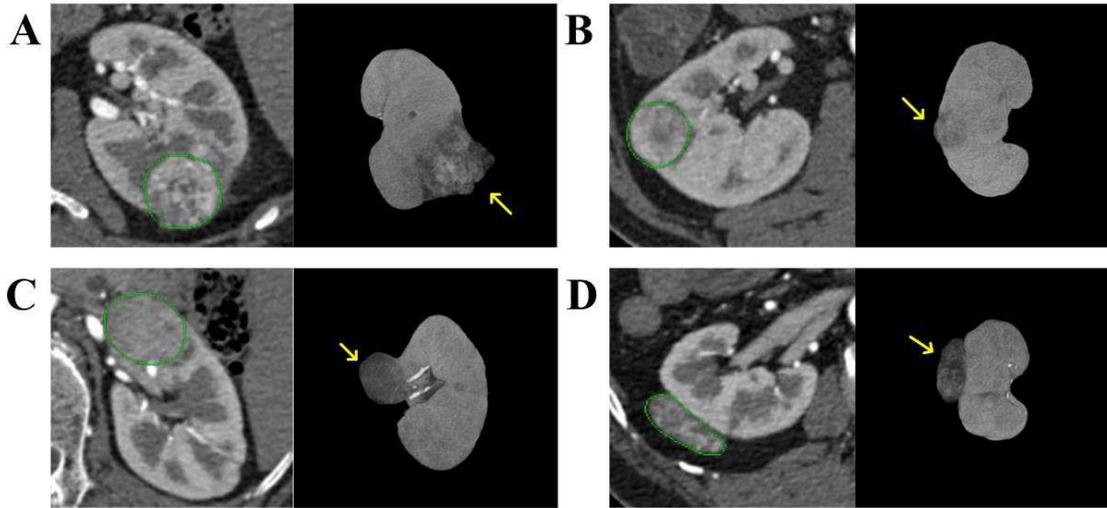
10 **1. Background**

11 Renal cancer is one of the ten most common cancers in human beings. The minimally invasive
12 laparoscopic partial nephrectomy (LPN) is now increasingly used to treat the renal cancer [1]. In
13 the clinical practice, some anatomical information such as the location and the size of the renal
14 tumor is very important for the LPN surgery planning. However, manual delineation of the contours
15 of the renal tumor and kidney in the pre-operative CT images including more than two hundred
16 slices is a time-consuming work. In recent years, deep neural networks have been the widely used
17 for organ and lesion segmentation in medical images [2]. However, fully-supervised deep neural
18 networks were trained by a large number of training images with pixel-wise labels, which take a
19 considerable time for radiologists to build. Thus, weakly supervised approaches attract more interest,
20 especially for medical image segmentation.

1 In recent years, several weakly-supervised CNNs have been developed for semantic segmentation
2 in natural images. According to the weak annotations used for CNN training, these approaches can
3 be divided into four main categories: bounding box [3-6], scribble [7, 8], points [9, 10] and image-
4 level labels [11-17]. However, as far as we know, there are only a few weakly-supervised methods
5 reported for the segmentation tasks in medical images. DeepCut [18] adopted an iterative
6 optimization method to train CNNs for brain and lung segmentation with the bounding-box labels
7 which are determined by two corner coordinates, and the target object is inside the bounding box.
8 In another weakly-supervised scenario [19], fetal brain MR images were segmented using a fully
9 convolutional network (FCN) trained by super-pixel annotations [20] which refer to an irregular
10 region composed of adjacent pixels with similar texture, color, brightness or other features.
11 Kervadec et al. [21] conducted a size loss on CNN, which was used to obtain the segmentation of
12 different organs from the scribbled annotations which annotate different areas and their classes.
13 These weakly learned-based methods have achieved comparable accuracy on normal organs but
14 have not yet been applied to lesions. The approaches for renal tumor segmentation are mainly based
15 on traditional methods such as level-set [22], SVM [23] and fully-supervised deep neural networks
16 [24, 25]. To the best of our knowledge, there is no weakly-supervised deep learning technique
17 reported for renal tumor segmentation.

18 As shown in Fig.1, the precise segmentation of renal tumors is a challenging task because of the
19 large variation of the size, location, intensity and image texture of renal tumors in CTA images. For
20 example, small tumors are often overlooked since they are difficult to be distinguished from the
21 normal tissue, as displayed in Fig.1(b). Different pathological types of renal tumors show varied
22 intensities and textures which increases the difficulty of segmentation [26]. Thus, the segmentation

1 of renal tumors by a weakly-supervised method is still an open problem.



2

3 **Fig. 1. Four contrast-enhanced CT images of different pathological renal tumors. The tumors are marked by**
4 **yellow arrows in 3D views. The manual contours of the renal tumors delineated by a radiologist are displayed**
5 **in 2D slices. The pathological subtypes of the renal tumors are clear cell renal cell carcinoma (RCC) in (a) and**
6 **(b), chromophobe RCC in (c) and angiomyolipoma in (d).**

7 In this paper, bounding boxes of renal tumors are provided as weak annotations to train a CNN
8 which can generate pixel-wise segmentation of renal tumors. Compared to the other types of
9 annotations, the bounding box is a simple way to be defined by radiologists [27]. The main
10 contributions of this paper are as follows:

11 (1) To the best of our knowledge, we proposed a weakly-supervised CNN for renal tumor
12 segmentation for the first time.

13 (2) The proposed method can accomplish network training faster and overcome the under-
14 segmentation problem compared with the iterative training strategy usually adopted by the other
15 weakly-supervised CNNs [18, 28].

1 (3) The experimental results of a 200-patients clinical dataset with different pathological types of
2 renal tumors show that the CNN trained by our method can provide precise renal tumor
3 segmentation.

4 The remaining paper is organized as follows: Section 2 describes the datasets used in this paper.
5 In Section 3 the method is introduced in detail. Experimental results are summarized in Section 4.
6 We give extra discussion in Section 5, a conclusion in Section 6 and abbreviations in Section 7. The
7 last section is the declarations of this paper.

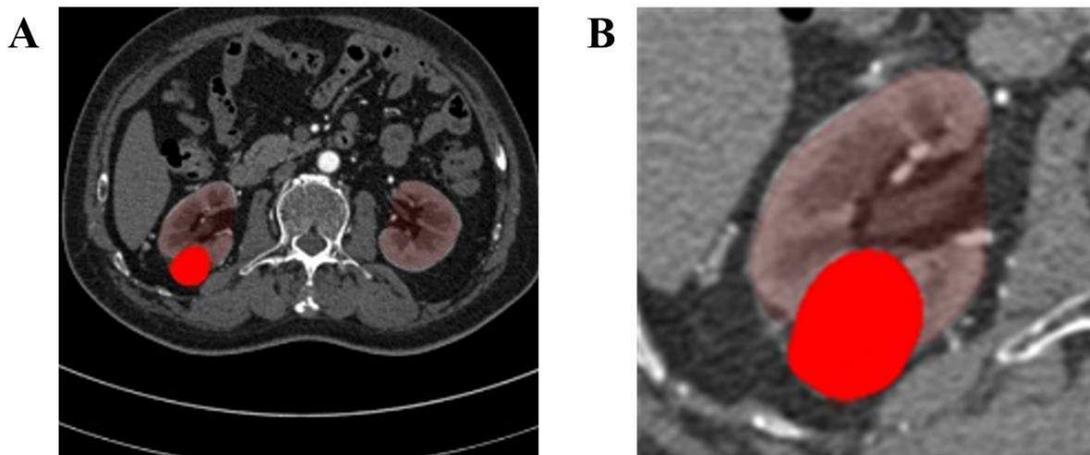
8 **2. Materials**

9 The pre-operative CT images of 200 patients who underwent an LPN surgery were included in
10 this study. The CT images were generated on a Siemens dual-source 64-slice CT scanner. The
11 contrast media was injected during the CT image acquisition. The study was already approved by
12 the institutional review board of Nanjing Medical University. Two scan phases including arterial
13 and excretion phases were performed for data acquisition. In this paper, CT images acquired in
14 arterial phase were used for training and testing. The arterial scan was triggered by the bolus tracking
15 technique after 100ml of contrast injection (Ultravist 370, Schering) in the antecubital vein at a
16 velocity of 5ml/s. Bolus tracking used for timing and scanning was started automatically 6s after
17 contrast enhancement reached 250HU in a region of interest (ROI) placed in the descending aorta.
18 The pixel size of these CT images is between 0.56mm^2 to 0.74mm^2 . The slice thickness and the
19 spacing in z-direction were fixed at 0.75mm and 0.5mm respectively. After LPN surgery,
20 pathological tests were performed to examine the pathological types of renal tumors. Five types of
21 renal tumors were included in this study, i.e. clear cell RCC (172 patients), chromophobe RCC (4

1 patients), papillary RCC (6 patients), oncocytoma (6 patients) and angiomyolipoma (12 patients).

2 The volume of the renal tumors' ranges from 12.21ml to 159.67ml and the mean volume is 42.58ml.

3 As shown in Fig.2(a), each original CT image was resampled to an isotropic volume with the size
4 of axial slice equal to 512*512. The original CT image contained the entire abdomen, whereas only
5 the area of the kidney needed to be considered in this experiment. Thus, the kidneys in the images
6 were firstly segmented by the multi-atlas-based method [29] to define the ROIs of kidneys as shown
7 in Fig.2(b). The multi-atlas-based method just produce initial segmentation of kidneys, two
8 radiologists checked the contours of kidneys and corrected them if necessary. The contours of
9 tumors were drawn manually by one radiologist with 7-years' experience and checked by another
10 radiologist with 15-years' experience in the cross-sectional slices. However, the pixel-wise masks
11 were only used for bounding boxes generation and testing dataset evaluation. Among 200-patient
12 images, 120 patients were selected to build the training dataset and the other 80 patients were used
13 as the testing dataset.



15 **Fig. 2. (a) The original image with labeled kidney and renal tumor. The region in red represents renal tumor.**

16 **(b) The cropped original image with the label for renal tumor segmentation.**

1 **3. Methods**

2 We train our proposed method via bounding boxes of renal tumors to obtain pixel-wise
3 segmentation. Thus, a pre-processing step is performed before the training procedure of weakly-
4 supervised model. In Section 3.1, the pre-processing including normalization and bounding box
5 generation is briefly introduced. Then the proposed weakly-supervised method is illustrated in detail
6 in Section 3.2. Finally, the parameters of training are explained in Section 3.3.

7 **3.1 Pre-processing**

8 **3.1.1 Normalization:**

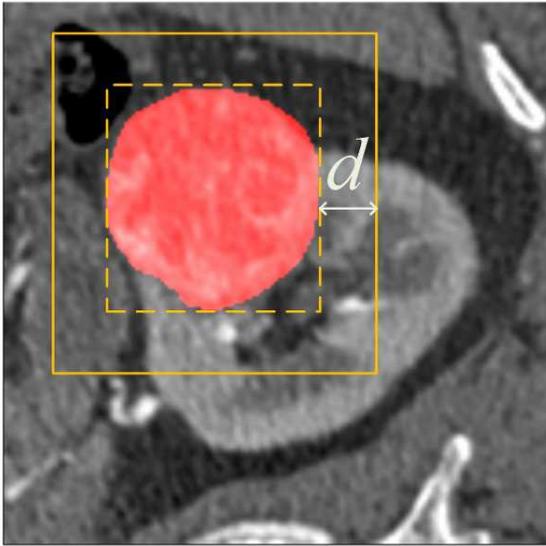
9 As is done in other studies, original CT images should be normalized before fed into the neural
10 network. Due to the existence of bones, contrast media and air in the intestinal tract, CT values in
11 the abdominal CT image or extracted ROIs can range from -1000HU to more than 800HU. Thus,
12 Hounsfield values were clipped to a range of -200 to 500 HU. After thresholding, the pixel values
13 in all images are normalized to 0~1 by Min-Max Normalization:

$$14 \quad X' = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (1)$$

15 **3.1.2 Bounding box generation**

16 In this paper, bounding boxes are generated by ground truth of renal tumors. As shown in Fig.3,
17 the bounding box of ground truth is shown in the dotted line. The parameter d in pixel represents
18 the margin added to the bounding box in our experiment to generate different types of weak
19 annotations. In addition, the reference labels of renal tumors in the training dataset were only used
20 to generate bounding boxes and not used for CNN training, and the reference labels in the testing

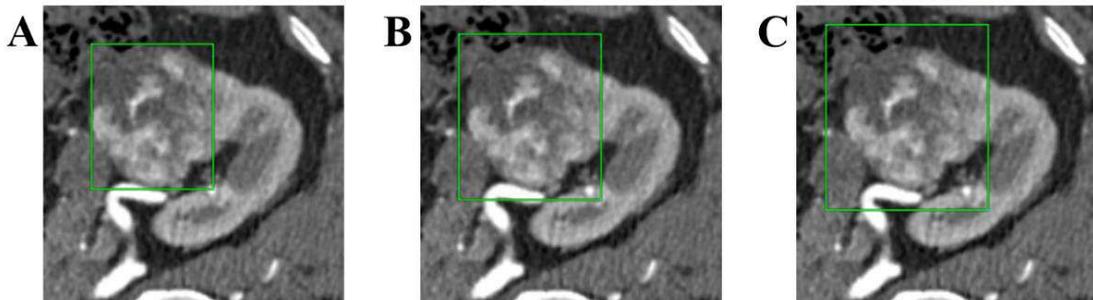
1 dataset were used for quantitative evaluation.



2

3 **Fig. 3. The bounding box with margin d is defined as weak annotations according to the label of renal**
4 **tumors.**

5 The bounding boxes with different margins are defined according to the ground truth and used as
6 weak annotations for CNN training. We set d to be 0, 5 and 10 pixels (Fig.4(a)-(c)) in our study to
7 simulate the manual weak annotations by radiologists. If the bounding boxes with margin d are
8 beyond the range of images, it will be limited in the region of images. As shown in Fig.4, the
9 comparison of bounding boxes with different margin values is given.

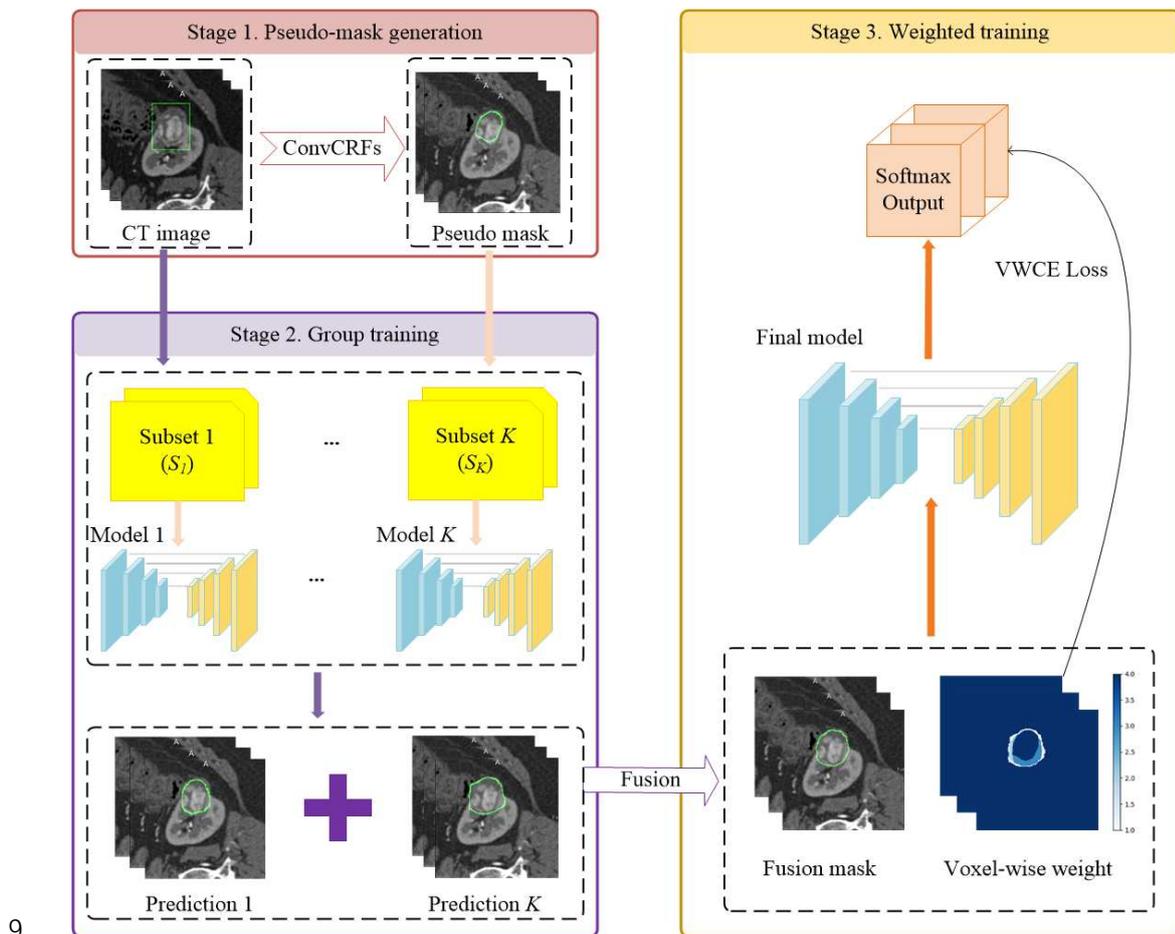


10

11 **Fig. 4. Comparison of bounding boxes with different margins. The 2D image is the maximum slice. Contours**
12 **in green correspond to bounding boxes.**

1 **3.2 Weakly supervised segmentation from bounding box**

2 Three main steps are included in the proposed method as shown in Fig.5. Firstly, we get pseudo
3 masks from bounding boxes by convolutional conditional random fields (ConvCRFs) [30]. Then,
4 in the group training stage, several CNNs are trained by using pseudo masks. Fusion masks and
5 voxel-wise weight map are generated based on the predictions of the CNNs trained in this stage. In
6 the last stage of weighted training, the final CNN is trained by fusion masks and voxel-wise
7 weighted cross-entropy (VWCE) loss function. These three main stages are described in the
8 following Section 3.2.1 to 3.2.3 respectively.



10 **Fig. 5. An overview of the proposed weakly-supervised method.**

1 3.2.1 Pseudo masks generation

2 As adopted by other methods [3, 18], the pseudo masks of renal tumors are generated from
3 bounding boxes as initialization for CNN model training. The quality of pseudo masks influences
4 the performance of CNN. Inspired by fully connected conditional random fields (CRFs) [31], this
5 problem can be regarded as maximum a posteriori (MAP) inference in a CRF defined over pixels
6 [5]. The CRF potentials take advantage of the context between pixels and encourage consistency
7 between similar pixels. Suppose an image $X = \{x_1 \dots x_N\}$ and corresponding voxel-wise label
8 $Y = \{y_1 \dots y_N\}$, here $y_i \in \{0,1\}$. $y_i = 0$ means x_i is located outside the bounding box, while $y_i = 1$
9 means x_i is located inside the bounding box. The CRF conforms to the Gibbs distribution. Then, the
10 Gibbs energy can be defined as:

$$11 \quad E(X) = \sum_i U(y_i) + \sum_{i,j} P(y_i, y_j) \quad (2)$$

12 where the first term is unary potential, representing the energy of assigning class y_i to the pixel x_i ,
13 which is given by the bounding box. The latter term represents the pairwise potential, which is used
14 to represent the energy of two pixels x_i and x_j in the image whose label are assigned to y_i and y_j
15 respectively. In the fully connected CRFs, the pairwise potential function is defined as follows:

$$16 \quad P(y_i, y_j) = \mu(y_i, y_j) \sum_{i \neq j \leq N} w \cdot g(f_i, f_j) \quad (3)$$

17 where w is a learnable parameter, g is the gaussian kernel defined by feature vectors f and μ is
18 a label compatibility function.

19 However, because the volumetric image was used in our study, the computation of fully connected
20 CRFs has high time complexity. Thus, inspired by Teichmann et al. [30], ConvCRFs were used for

1 our pseudo masks generation. ConvCRFs adds the assumption of conditional independence into
 2 fully connected CRFs. Here, the matrix of gaussian kernel changes to:

$$3 \quad g(f_i, f_j) = \exp\left(-\sum_{i \neq j \leq D} \frac{f_i - f_j}{2\theta^2}\right) \quad (4)$$

4 where θ is a learnable parameter and D is the Manhattan distance between pixels x_i and x_j , the
 5 pairwise energy is zero when the Manhattan distance exceeds D . The complexity of pairwise potential
 6 is simplified when conditional independence is added.

7 The merged kernel matrix G is calculated by $\sum w \cdot g$, and the inference result is $\sum G \cdot X$ which
 8 is similar to convolutions of CNNs. This assumption makes it possible to reformulate the inference
 9 in terms of convolutions in CRF, which can carry out efficient GPU calculation and complete feature
 10 learning. Thus, we can quickly get pseudo masks of renal tumors by minimizing the object function
 11 defined by Eq. (2).

12 **3.2.2 Group training and fusion mask generation**

13 Once we have generated pseudo masks of renal tumors, these masks are fed into CNN as weak
 14 labels for parameter learning. Most of weakly supervised segmentation methods used iterative
 15 training [5, 7] to optimize the accuracy of the weak labels from coarse to fine. However, the
 16 preliminary results showed that this iterative strategy is hard to improve the accuracy of pseudo
 17 masks due to the difficulties of the renal tumor segmentation mentioned before. To overcome this
 18 problem, we proposed a new CNN training strategy instead of iterative training method.

19 In the group training stage, we have input images $\{X_1 \dots X_M\}$ and pseudo masks $\{I_1 \dots I_M\}$. The
 20 input training dataset is divided into K subsets $\{S_1 \dots S_K\}$. For each subset S_k , a CNN $f(X; \theta_k), X \in$

1 S_k with parameter θ_k is trained. In total, we can get K CNNs trained in this stage. After that, for
 2 each image X_m , we can get K predictions $\{P_m^1 \dots P_m^K\}$ of renal tumors by these CNN models. We
 3 denote that $P_m^k = f(X_m; \theta_k)$. Pseudo code of group training is shown in Algorithm 1.

Algorithm 1. Group training

Input: $\{X_1 \dots X_M\}, \{I_1 \dots I_M\}$

Divide input into K subsets $\{S_1 \dots S_K\}, S_1 \cap \dots \cap S_K = \emptyset$

for $k = 1 : K$ **do**

train CNN $f(X; \theta_k), X \in S_k$

end

for $m = 1 : M$ **do**

for $k = 1 : K$ **do**

obtain $P_m^k = f(X_m; \theta_k)$

end

end

Output: $\{\{P_1^1 \dots P_1^K\} \dots \{P_M^1 \dots P_M^K\}\}$

4 One thing worth to be mentioned is that one image in the training dataset is used to train only one
 5 CNN model in this stage. Once K CNN models are trained successfully, all the images in the training
 6 dataset will be used to test each CNN model and obtain K results for prediction. Thus, the proposed
 7 group training strategy can ameliorate the overfitting of the model. In order to alleviate the under-
 8 segmentation in the K predictions, a mask image is generated by fusing these predictions. The fusion
 9 mask is defined as follows:

$$10 \quad FM_m = ConvCRFs(PM_m \cup P_m^1 \cup \dots \cup P_m^K) \quad (5)$$

11 where FM indicates the fusion masks, and PM indicates pseudo masks generated in Section 3.2.1.
 12 The ConvCRFs is adopted to refine the union of all prediction masks. The outputs of ConvCRFs

1 will be used as the new weak labels for the next weighted training stage. In addition, a weight map
 2 is generated simultaneously which is defined as follows,

$$3 \quad v_m = PM_m + P_m^1 + \dots + P_m^K, v[v = 0] = K + 1 \quad (6)$$

4 When the predicted label of a voxel is renal tumor in one prediction result, its v_m will be an integer
 5 within the range of 1 to $K+1$. When v_m is equal to 0, its value will be reset to $K+1$ to represent the
 6 weight of background.

7 **3.2.3 Training with VWCE loss**

8 After Section 3.2.1 and 3.2.2, the fusion masks of training dataset are generated for the final CNN
 9 model training in this stage. Only the final CNN model will be used for testing dataset evaluation.
 10 In this stage, we train the CNN on the whole training dataset with the fusion masks. In addition, a
 11 new voxel-wise weighted cross-entropy (VWCE) loss function is designed to constrain the CNN
 12 training procedure. The traditional cross-entropy loss is defined as follows:

$$13 \quad L_{CE} = -\frac{1}{M} \sum_{m \in M} \sum_{c \in C} FM_{m,c} \log f(X_{m,c}; \theta) \quad (7)$$

14 where FM are fusion masks defined in Eq. (5), $f(X; \theta)$ are the outputs of CNN, M represents
 15 the number of samples and C represents the number of classes. In Eq. (7), pixels belonging to
 16 different classes have equal weight. In the case of unbalanced datasets, [32] proposed weighted
 17 cross-entropy loss defined as follows:

$$18 \quad L_{WCE} = -\frac{1}{M} \sum_{m \in M} \sum_{c \in C} w_c FM_{m,c} \log f(X_{m,c}; \theta) \quad (8)$$

19 where, w_c represents the weight of class c . Considering the weak annotations used in the training
 20 procedure, the voxel-wise weight map generated in the previous stage represents the probability of

1 the predicted class given in the fusion mask. Thus, the voxel-wise weights obtained in Eq. (6) are
2 introduced into Eq. (8) which is defined as follows:

$$3 \quad L_{VWCE} = -\frac{1}{M} \sum_{m \in M} v_m \sum_{c \in C} w_c F M_{m,c} \log f(X_{m,c}; \theta) \quad (9)$$

4 Finally, we conduct the final CNN model training with VWCE loss function on fusion masks.
5 Our evaluations are all conducted on CNN trained in this stage.

6 **3.3 Training**

7 **3.3.1 Data augmentation**

8 The ROIs of the pathological kidneys were cropped from the original images. The size of ROI is
9 fixed at $150 \times 150 \times N$. Due to limited memory of GPU, the original ROIs were resampled to
10 $128 \times 128 \times 64$ before fed into the network. For each data, random crops and flipping were used for
11 data augmentation. After data augmentation, the original 120 CT images were augmented into 14400
12 images for the CNN training.

13 **3.3.2 Parameter settings**

14 The input are ROIs of kidneys and bounding boxes without any other annotations. Considering
15 that UNet [32] has been widely used for medical image segmentation, we adopted UNet to be the
16 CNN models in stage2 and stage3 in our experiments. The network parameters are updated by means
17 of the back-propagation algorithm using the Adam optimizer. The initial learning rate was set to be
18 0.001 and decreased by $decayed_learning_rate = learning_rate * decay_rate^{\frac{global_step}{decay_steps}}$. In
19 each epoch of training, it takes 3600 iterations to traverse all the training images with the batch size
20 of 4. The class weights of cross-entropy w_c in Eq. (8) and (9) were set to 1.0 and 0.2 for renal tumor

1 and background respectively.

2 In stage2, we set the number of subset K to 3 for the training dataset of 120 CT images. Each
3 subset contains 40 CT images. Three CNN models were trained to generate corresponding
4 predictions of each training image. And fusion masks were generated by these predictions. The loss
5 used in this stage is WCE loss defined in Eq. (8).

6 In stage3, the final CNN is trained by fusion masks as weak annotation labels. We evaluated the
7 performance of the final CNN model with 80 patient images. In order to remove some misclassified
8 outlier voxels, a connected component analysis with an 18-connectivity in 3D was carried out finally.
9 The largest connected component in the output of the final CNN model was extracted as the
10 segmentation results of renal tumors.

11 **3.4. Existing methods**

12 We mainly compared with two weakly-supervised methods, i.e., SDI [5] and constrained-CNN
13 [21]. The SDI method used 2D UNet to generate weak labels from bounding box by recursive
14 training and carry out final segmentation. The weakly-supervised information used in the
15 constrained-CNN method includes scribbles and the volume of target tissue. In this paper, the
16 scribbles annotations used in constrained-CNN were generated by employing binary erosion on
17 ground truth for every slice. Furthermore, the volumetric threshold of renal tumor was used in the
18 loss function of Constrained-CNN. It was set to $[0.9V, 1.1V]$, where V represents the volume of
19 renal tumor in ground truth. As the architecture of UNet was used in [5] and [21], as well as our
20 proposed method, the UNet was trained by all the training dataset with the pixel-wise labels to
21 generate a fully-supervised UNet model for extensive comparison.

1 **4. Results**

2 Our method has been implemented using PyTorch framework in version 1.1.0. The network
3 training and testing experiments were performed on a workstation with: CPU of i7-5930K, 128GB
4 RAM and a GPU card of NVIDIA TITAN Xp of 12GB memory.

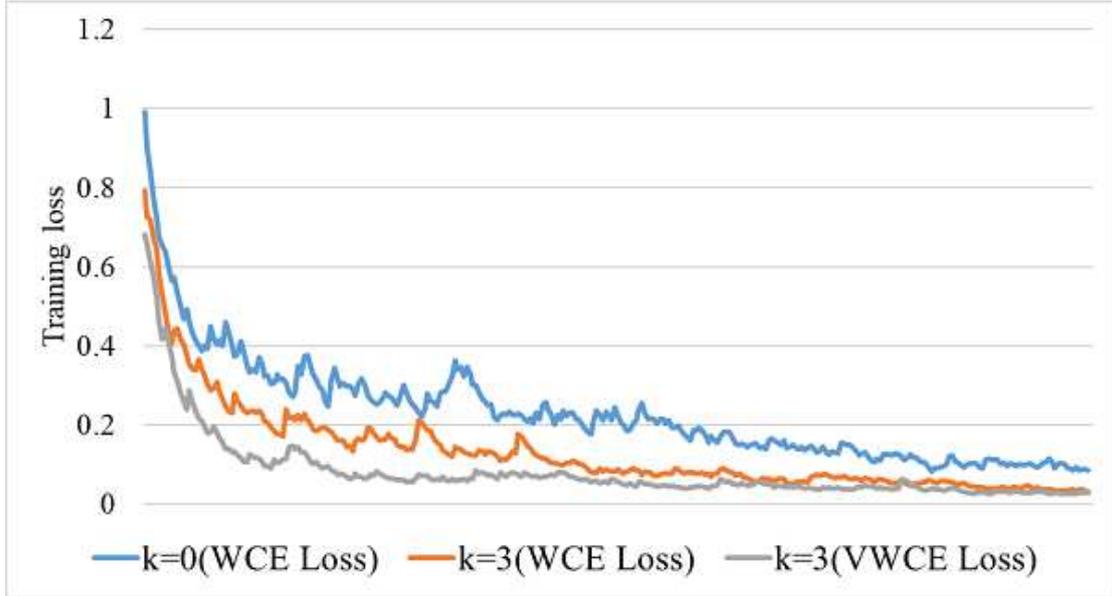
5 **4.1 The comparison of different weak labels and training losses**

6 As shown in Table 1, DSCs between the different masks and the ground truth of the training
7 dataset are displayed. The DSCs of bounding boxes are 0.666, 0.466 and 0.341 respectively when
8 the margins of bounding box were set to 0, 5 and 10 pixels. The DSCs of pseudo masks generated
9 by ConvCRFs can reach 0.862, 0.801 and 0.679. However, the DSCs of fusion masks generated
10 after group training has even higher DSC than pseudo masks. Obviously, the rectangular bounding
11 boxes were improved significantly by the Stage 1 and Stage 2.

12 **Table 1** DSCs between different weak labels and ground truths of the training dataset

	Bounding boxes	Pseudo masks	Fusion masks
$d=0$	0.666	0.862	0.874
$d=5$	0.466	0.801	0.810
$d=10$	0.341	0.679	0.691

13 Furthermore, the improvements of the weak labels contribute to the training of the final CNN
14 model. Fig.6 shows the training loss of the final CNN model with different parameters. Without
15 group training, the training loss shows the slowest rate and the highest loss value during training.
16 Contrarily, the usage of group training and VWCE loss makes the model converges faster and better.



1

2 **Fig. 6. Training losses of the final CNN model in stage3 with different parameters.**

3 **4.2 Evaluation of segmentation results of renal tumors in the testing dataset with different**
 4 **parameters**

5 The DSC, Hausdorff distance (HD) [33] and average surface distance (ASD) were adopted to
 6 evaluate the segmentation results of our proposed method. The segmentation results of renal tumors
 7 in the testing dataset were obtained with different settings of parameters, i.e. number of groups, loss
 8 function and margin of bounding box. The comparison of DSCs in the testing dataset is displayed
 9 in Table 2. $k=0$ means that the procedure of stage2 not used. In this situation, the pseudo masks
 10 generated by ConvCRFs were used as weak labels directly for the final CNN model training in the
 11 stage3. The loss functions used during the final model training is marked in the parentheses. MC
 12 represents the connected component analysis in the post-processing step.

13 **Table 2** Comparison of segmentation results of testing dataset with different margins

	DSC	HD	ASD
--	-----	----	-----

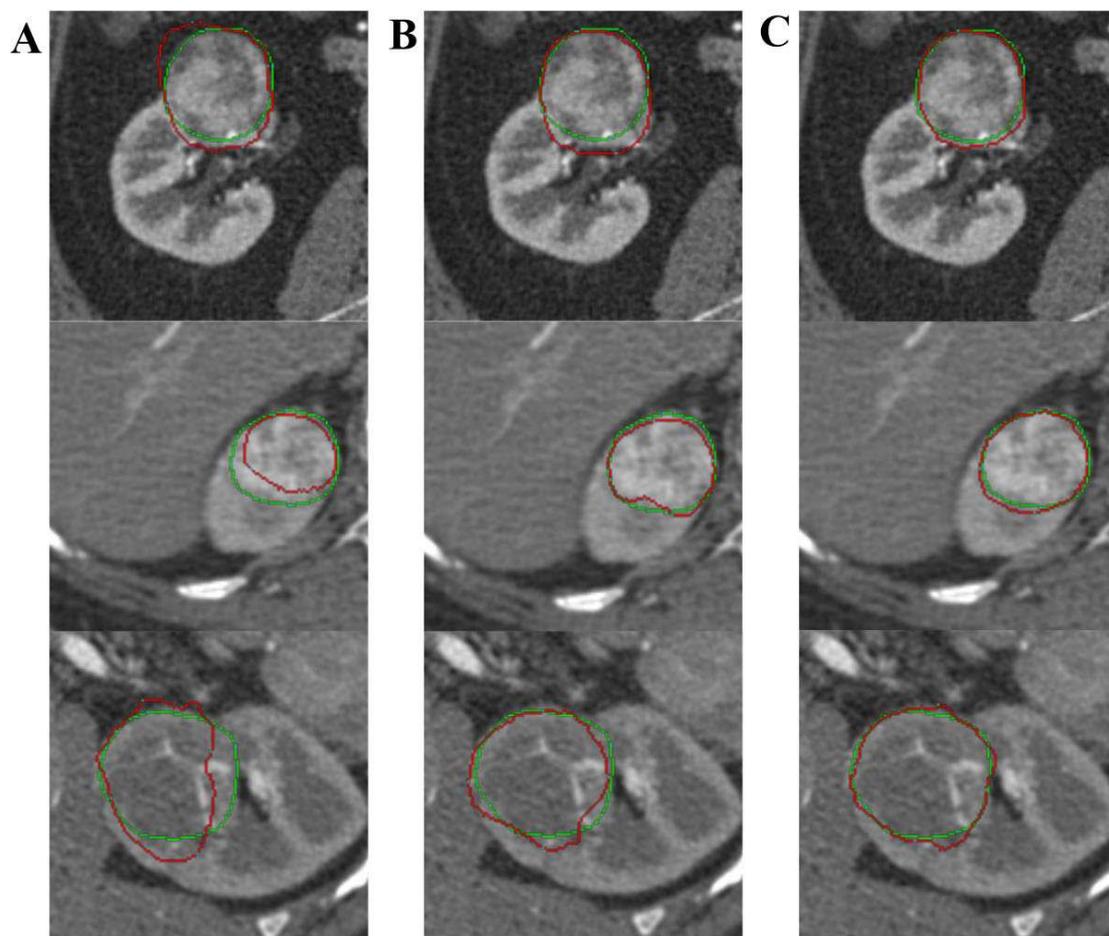
$d=0$	$k=0$ (WCE Loss)	0.788	65.806	6.265
	$k=3$ (WCE Loss)	0.822	34.187	3.889
	$k=3$ (VWCE Loss)	0.834	40.617	3.361
	$k=3$ (VWCE Loss) + 3D MC	0.834	14.346	2.664
$d=5$	$k=0$ (WCE Loss)	0.733	32.459	5.332
	$k=3$ (WCE Loss)	0.784	70.948	7.988
	$k=3$ (VWCE Loss)	0.820	37.633	3.879
	$k=3$ (VWCE Loss) + 3D MC	0.826	15.811	2.838
$d=10$	$k=0$ (WCE Loss)	0.695	58.286	7.499
	$k=3$ (WCE Loss)	0.720	81.611	7.804
	$k=3$ (VWCE Loss)	0.741	36.127	4.672
	$k=3$ (VWCE Loss) + 3D MC	0.742	21.233	4.350

1 **The impact of group training:** According to the values in Table 2, group training can effectively
2 improve the DSC. The DSCs increased by 3.4%, 5.1% and 2.5% when the margin of bounding box
3 was set to 0, 5 and 10 pixels respectively.

4 **The impact of VWCE loss:** The usage of VWCE loss made further improvement of the DSC.
5 The DSCs increased by 1.2%, 3.6%, and 2.1% respectively when the margin of bounding box was
6 set to 0, 5 and 10 pixels. In addition, the application of VWCE loss and MC can alleviate the outliers
7 in the segmentation result. The values of HD and ASD decreased significantly. Finally, the highest
8 DSCs of 0.834, 0.826 and 0.742 can be achieved respectively when different margins of bounding
9 box were set.

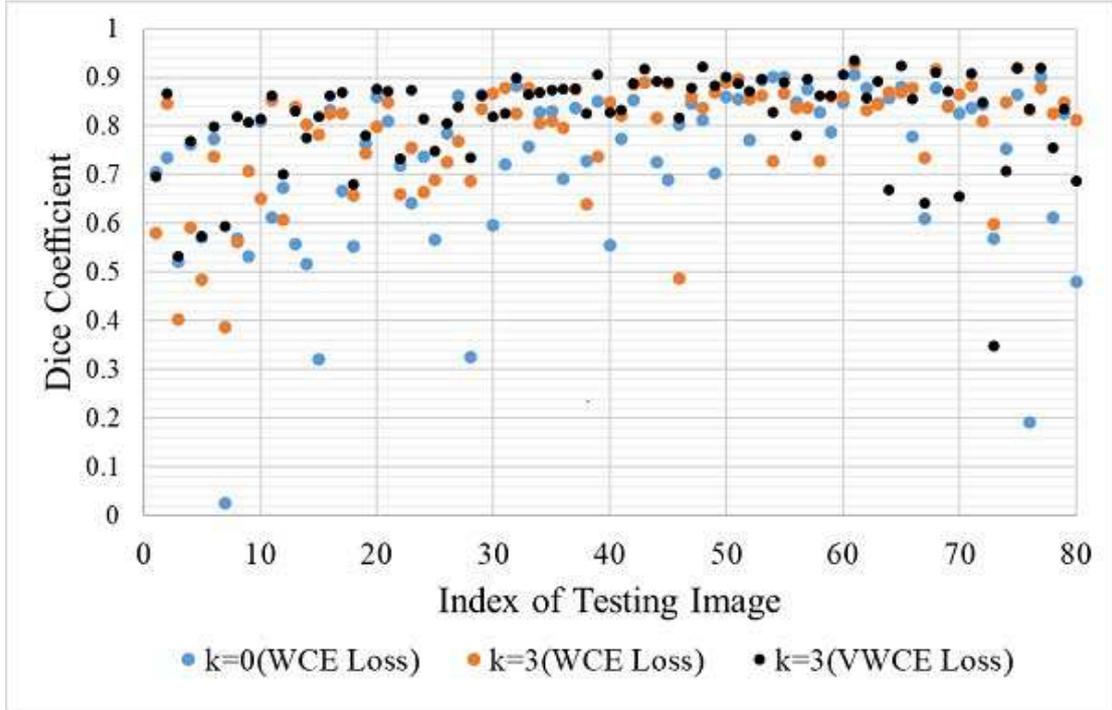
10 Fig.7 shows the 2D visualization of segmentation results with different parameters. Obviously,
11 renal tumors cannot be segmented precisely without group training as shown in Fig.7(a). With the
12 application of group training, the over- or under-segmentation of tumors is significantly improved

1 (Fig.7b). However, the segmentations of the boundary are still imprecise. With the application of
2 group training and VWCE loss function, the best segmentation results have been obtained as shown
3 in Fig.7(c).



4
5 **Fig. 7. The comparison of 2D segmentation results with different parameters: $k=0$ with WCE loss (a), $k=3$**
6 **with WCE loss (b), $k=3$ with VWCE loss (c). Contours in green and red correspond to ground truths and**
7 **segmentation results respectively.**

8 The DSC of each case in the testing dataset with different parameters is shown in Fig.8. For
9 testing dataset, it can be seen that our three-stage training strategy with VWCE loss has significantly
10 improved the segmentation results in most images and achieves the best improvement of DSC.



1

2 Fig. 8. DSC of each case in the testing dataset with different parameters. The index of images is ranked
 3 according to the volume of renal tumors.

4 **4.3 Comparison with other methods**

5 Three methods including two weakly-supervised methods (SDI and constrained-CNN) and one
 6 fully-supervised method (UNet) were used to compare with our proposed method. These methods
 7 are briefly summarized in section 3.4. For model training, the computation time of our proposed
 8 method is about 48 hours, the SDI method is about 80 hours, and the Constrained-CNN and fully-
 9 supervised UNet are about 24 hours. For model testing, the computation time of our proposed
 10 method is similar to the fully-supervised method. Our network can generate the segmentation result
 11 of a single image in a few seconds.

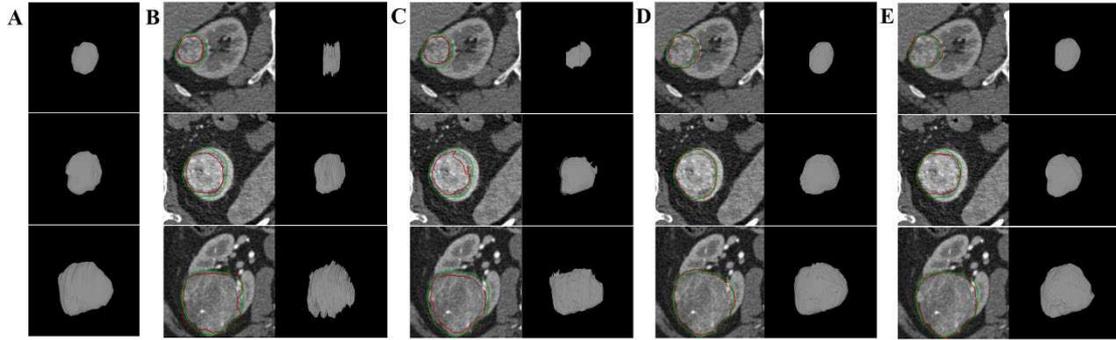
12 Table 3 is the comparison of segmentation results among our method, the other two existing
 13 weakly-supervised methods and fully-supervised method. We only compared the bounding box with

1 $d=5$ for simplicity. Experiments show that our method achieves the best results of DSC, HD and
 2 ASD, which are 0.826, 15.811 and 2.838 respectively. In terms of DSC, neither SDI nor
 3 Constrained-CNN reaches the values higher than 0.8. One thing worth to be mentioned is that the
 4 evaluation metrics are not improved effectively in SDI after MC since we deal with it in 2D situation.
 5 When the margin is lower than 5, the performance of our method is close to the results obtained by
 6 the fully-supervised UNet.

7 **Table 3** Comparison of testing results with different methods

	DSC	HD	ASD
Constrained-CNN [21]	0.705	102.178	8.271
Constrained-CNN [21] + 3D MC	0.712	20.939	5.493
SDI [5]	0.766	73.514	4.639
SDI [5] + 2D MC	0.766	72.368	4.524
Ours ($d=5$)	0.820	37.633	3.879
Ours ($d=5$) + 3D MC	0.826	15.811	2.838
UNet [32] (Fully-supervised)	0.849	84.69	4.886
UNet [32] (Fully-supervised) + 3D MC	0.859	14.252	2.048

8 Fig.9 shows the comparison of segmentation results obtained by different methods. For SDI
 9 method, the shape of the segmented renal tumor in 3D is not continuous as shown in Fig.9 (b).
 10 Furthermore, SDI and Constrained-CNN still suffer from the under-segmentation problem. While,
 11 our proposed method (d) presents better segmentation results which are similar to the fully-
 12 supervised method (e) in visual.



1

2 **Fig. 9. The comparison of the results from three testing images obtained by different methods: 3D ground**
 3 **truth (a), SDI (b), Constrained-CNN(c), the proposed method (d) and fully-supervised method (e). Contours**
 4 **in green and red correspond to ground truth and segmentation results respectively.**

5 **5. Discussion**

6 According to our experimental results, our proposed weakly-supervised method can provide
 7 accurate renal tumor segmentation. The major difficulty for weakly-supervised methods is that
 8 feature maps learned by CNN models can be misled by under- or over-segmentation in the weak
 9 masks. Therefore, the key factor in weakly-supervised segmentation is to generate reliable masks
 10 from the input weak labels. In this paper, the application of pseudo masks generation and group
 11 training improve the quality of the weak masks used for the final CNN model training as shown in
 12 Table 1 and 2.

13 Furthermore, as shown in Fig.8, the DSCs of large and small tumors are relatively low. It is easy
 14 to understand that the DSCs of the small renal tumors are sensitive to the over- or under-
 15 segmentation in the predictions. While in large tumor, the shape and texture of the tumor are
 16 complicated, which leads to the difficulties of the segmentation. Although this problem exists in all

1 three methods, our proposed method shows the most significant improvement compared with the
2 other two methods.

3 Finally, one limitation of this study is the lack of validation of the final CNN model with external
4 datasets. The training and testing datasets in this paper are from the same hospital. Additional
5 validation of the final CNN model with multi-center or multi-vendor images will be performed in
6 the future. Due to the differences in image acquisition protocols or the other factors, the CNN model
7 trained in this paper may not be able to achieve a similar performance on the other datasets. However,
8 the parameters in our model can be optimized by fine-tuning with the external datasets to improve
9 the accuracy. In particular, the main advantage of our method is the use of weak labels for network
10 training, which does not take much time for radiologists to generate bounding-box labels.

11 **6. Conclusion**

12 In this paper we have presented a novel three-stage training method for weakly supervised CNN
13 to obtain precise renal tumor segmentation. The proposed method mainly relies on the group training
14 and weighted training phases to improve not only the efficiency of training but also the accuracy of
15 segmentation. Experimental results with 200 patient images show that the DSCs between ground
16 truth and segmentation results can reach 0.834, 0.826 when the margin of bounding box was set to
17 0 and 5, which are close to the fully-supervised model which is 0.859. The comparison between our
18 proposed method and the other two existing methods also demonstrate that our method can generate
19 a more accurate segmentation of renal tumors than the other two methods.

20 **7. Abbreviations**

- 1 **ASD:** average surface distance
- 2 **CE:** cross-entropy
- 3 **CNN:** convolutional neural network
- 4 **ConvCRFs:** convolutional conditional random fields
- 5 **CRF:** conditional random field
- 6 **CT:** computed tomography
- 7 **CTA:** computed tomographic angiography
- 8 **DSC:** dice coefficient
- 9 **FCN:** fully convolutional network
- 10 **HD:** Hausdorff distance
- 11 **LPN:** laparoscopic partial nephrectomy
- 12 **MAP:** maximum a posteriori
- 13 **MC:** maximum connected component
- 14 **MR:** magnetic resonance
- 15 **RCC:** renal cell carcinoma
- 16 **ROI:** region of interest
- 17 **SVM:** support vector machine

1 **VWCE:** voxel-wise weighted cross-entropy

2 **WCE:** weighted cross-entropy

3 **8. Declarations**

4 **Ethics approval and consent to participate**

5 This study was carried out in accordance with the recommendations of name of the Nanjing Medical
6 University's Committee with written informed consent from all subjects. All subjects gave written
7 informed consent in accordance with the Declaration of Helsinki. The protocol was approved by the
8 name of the Nanjing Medical University's Committee.

9 **Consent for publication**

10 Not applicable.

11 **Availability of data and materials**

12 The clinical data and materials used in this paper are not open to public, but are available from the
13 corresponding author on reasonable request.

14 **Competing interests**

15 Yang Chen, one of the co-authors, is a member of the editorial board (Associate Editor) of this
16 journal. The other authors have no conflicts of interest to disclose.

17 **Funding**

18 This study was funded by a grant from the National Key Research and Development Program of
19 China (2017YFC0107900), National Natural Science Foundation (31571001, 61828101), Key

1 Research and Development Project of Jiangsu Province BE2018749) and the Southeast University-
2 Nanjing Medical University Cooperative Research Project (2242019K3DN08). These funds
3 provided financial support for the research work of our article but had no role in the study.

4 **Authors' contributions**

5 GYY and CXW designed the proposed method and implemented this method. LJT and PFS outlined
6 the data label. JY, YC, JLD, HZS and LML performed the experiments and the analysis of the
7 results. All authors have been involved in drafting and revising the manuscript and approved the
8 final version to be published. All authors read and approved the final manuscript.

9 **Acknowledgements**

10 We acknowledge Key Laboratory of Computer Network and Information Integration, Southeast
11 University, Ministry of Education, Nanjing, People's Republic of China for providing us the
12 computing platform.

13 **Reference**

14 [1] Ljungberg B, Bensalah K, Canfield S, Dabestani S, Hofmann F, Hora M, et al. EAU guidelines on renal cell 569
15 carcinoma 2014 update. *European Urology*. 2015; 67(5): 913-924.

16 [2] Litjens GJ, Kooi T, Bejnordi BE, Setio AA, Ciompi F, Ghahfarooian M, et al. A survey on deep learning in medical
17 image analysis. *Medical Image Analysis*. 2017; 42: 60-88.

18 [3] Dai J, He K, Sun J. BoxSup: Exploiting Bounding Boxes to Supervise Convolutional Networks for Semantic
19 Segmentation. *The IEEE International Conference on Computer Vision*. 2015; pp: 1635-1643.

- 1 [4] Papandreou G, Chen L, Murphy K, Yuille AL. Weakly-and Semi-Supervised Learning of a Deep Convolutional
2 Network for Semantic Image Segmentation. The IEEE International Conference on Computer Vision. 2015; pp:
3 1742-1750.
- 4 [5] Khoreva A, Benenson R, Hosang J, Hein M, Schiele B. Simple Does It: Weakly Supervised Instance and
5 Semantic Segmentation. The IEEE Conference on Computer Vision and Pattern Recognition. 2017; pp: 876-885.
- 6 [6] Hu R, Dollar P, He K, Darrell T, Girshick R. Learning to Segment Every Thing. The IEEE Conference on
7 Computer Vision and Pattern Recognition. 2018; pp: 4233-4241.
- 8 [7] Tang M, Djelouah A, Perazzi F, Boykov Y, Schroers C. Normalized Cut Loss for Weakly-Supervised CNN
9 Segmentation. The IEEE Conference on Computer Vision and Pattern Recognition. 2018; pp: 1818-1827.
- 10 [8] Lin D, Dai J, Jia J, He K, Sun J. ScribbleSup: Scribble-Supervised Convolutional Networks for Semantic
11 Segmentation. The IEEE Conference on Computer Vision and Pattern Recognition. 2016; pp: 3159-3167.
- 12 [9] Maninis K, Caelles S, Pontuset J, Gool L. Deep Extreme Cut: From Extreme Points to Object Segmentation.
13 The IEEE Conference on Computer Vision and Pattern Recognition. 2018; pp: 616-625.
- 14 [10] Bearman A, Russakovsky O, Ferrari V, Fei-Fei L. What's the Point: Semantic Segmentation with Point
15 Supervision. European Conference on Computer Vision. 2016; pp: 549-565.
- 16 [11] Pathak D, Shelhamer E, Long J, Darrell T. Fully Convolutional Multi-Class Multiple Instance Learning. 2014;
17 arXiv: 1412.7144.
- 18 [12] Pinheiro PO, Collobert R. From image-level to pixellevel labeling with convolutional networks. The IEEE
19 Conference on Computer Vision and Pattern Recognition. 2015; pp: 1713-1721
- 20 [13] Saleh FS, Aliakbarian MS, Salzmann M, Petersson L, Gould S, Alvarez JM. Built-in Foreground/Background

- 1 Prior for Weakly-Supervised Semantic Segmentation. European Conference on Computer Vision. 2016; pp: 413-432.
- 2 [14] Wei Y, Liang X, Chen Y, Shen X, Cheng M, Feng J, et al. STC: A Simple to Complex Framework for Weakly-
3 Supervised Semantic Segmentation. IEEE Transactions on Pattern Analysis and Machine Intelligence. 2017; 39(11):
4 2314-2320.
- 5 [15] Kolesnikov A, Lampert CH. Seed, Expand and Constrain: Three Principles for Weakly-Supervised Image
6 Segmentation. European Conference on Computer Vision. 2016; pp: 695-711.
- 7 [16] Qi X, Liu Z, Shi J, Zhao H, Jia J. Augmented Feedback in Semantic Segmentation under Image Level
8 Supervision. European Conference on Computer Vision. 2016; pp: 90-105.
- 9 [17] Wei Y, Feng J, Liang X, Cheng M, Zhao Y, Yan S. Object Region Mining with Adversarial Erasing: A Simple
10 Classification to Semantic Segmentation Approach. The IEEE Conference on Computer Vision and Pattern
11 Recognition. 2017; pp: 1568-1576.
- 12 [18] Rajchl M, Lee MC, Oktay O, Kamnitsas K, Passerat-Palmbach J, Bai W, et al. DeepCut: object segmentation
13 from bounding box annotations using convolutional neural networks. IEEE Transactions on Medical Imaging. 2017;
14 36(2): 674-683.
- 15 [19] Rajchl M, Lee MC, Schrans F, Davidson A, Passerat-Palmbach J, Tarroni G, et al. Learning under distributed
16 weak supervision. 2016; arXiv: 1606.01100.
- 17 [20] Achanta R, Shaji A, Smith K, Lucchi A, Fua P, Süsstrunk S. Slic superpixels compared to state-of-the-art
18 superpixel methods. IEEE Transactions on Pattern Analysis and Machine Intelligence. 2012; 34 (11): 2274–2282.
- 19 [21] Kervadec H, Dolz J, Tang M, Granger E, Boykov Y, Ayed IB. Constrained-CNN losses for weakly supervised
20 segmentation. Medical Image Analysis. 2019; 54: 88-99.

- 1 [22] Linguraru MG, Yao J, Gautam R, Peterson J, Li Z, Linehan WM, et al. Renal tumor quantification and
2 classification in contrast-enhanced abdominal CT. *Pattern Recognition*. 2009; 42(6): 1149-1161.
- 3 [23] Linguraru MG, Wang S, Shah F, Gautam R, Peterson J, Linehan WM, et al. Automated noninvasive
4 classification of renal cancer on multiphase CT. *Medical Physics*. 2011; 38(10): 5738-5746.
- 5 [24] Yang G, Li G, Pan T, Kong y, Wu J, Shu H, et al. Automatic Segmentation of Kidney and Renal Tumor in CT
6 Images Based on 3D Fully Convolutional Neural Network with Pyramid Pooling Module. *International Conference*
7 *on Pattern Recognition*. 2018; pp: 3790-3795.
- 8 [25] Yu Q, Shi Y, Sun J, Gao Y, Zhu J, Dai Y. Crossbar-Net: A Novel Convolutional Neural Network for Kidney
9 Tumor Segmentation in CT Images. *IEEE Transactions on Image Processing*. 2019; 28(8): 4060-4074.
- 10 [26] Zhang J, Lefkowitz RA, Ishill NM, Wang L, Moskowitz CS, Russo P, et al. Solid Renal Cortical Tumors:
11 Differentiation with CT. *Radiology*. 2007; 244(2): 494-504.
- 12 [27] Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, et al. Microsoft Coco: Common objects in context.
13 *European Conference on Computer Vision*. 2014; pp: 740–755.
- 14 [28] Wang X, You S, Li X, Ma H. Weakly-Supervised Semantic Segmentation by Iteratively Mining Common Object
15 Features. *The IEEE Conference on Computer Vision and Pattern Recognition*. 2018; pp: 1354-1362.
- 16 [29] Yang G, Gu J, Chen Y, Liu W, Tang L, Shu H, et al. Automatic kidney segmentation in CT images based on
17 multi-atlas image registration. *Annual International Conference of the IEEE Engineering in Medicine and Biology*
18 *Society*. 2014; pp: 5538-5541.
- 19 [30] Teichmann M, Cipolla R. Convolutional CRFs for Semantic Segmentation. 2018; arXiv: 1805.04777.
- 20 [31] Krahenbuhl P, Koltun V. Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials. *Advances*

1 in Neural Information Processing Systems. 2011; pp: 109-117.

2 [32] Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image Segmentation.

3 International Conference on Medical Image Computing and Computer Assisted Intervention. 2015; pp: 234-241.

4 [33] Huttenlocher DP, Klanderman GA, Rucklidge WJ. Comparing images using the Hausdorff distance. IEEE

5 Transactions on Pattern Analysis and Machine Intelligence. 1993; 15(9): 850-863.

6 **Figure Legends**

7 Figure. 1. Four contrast-enhanced CT images of different pathological renal tumors. The tumors are

8 marked by yellow arrows in 3D views. The manual contours of the renal tumors delineated by a

9 radiologist are displayed in 2D slices. The pathological subtypes of the renal tumors are clear cell

10 renal cell carcinoma (RCC) in (a) and (b), chromophobe RCC in (c) and angiomyolipoma in (d).

11 Figure. 2. (a) The original image with labeled kidney and renal tumor. The region in red represents

12 renal tumor. (b) The cropped original image with the label for renal tumor segmentation.

13 Figure. 3. The bounding box with margin d is defined as weak annotations according to the label of

14 renal tumors.

15 Figure. 4. Comparison of bounding boxes with different margins. The 2D image is the maximum

16 slice. Contours in green correspond to bounding boxes.

17 Figure. 5. An overview of the proposed weakly-supervised method.

18 Figure. 6. Training losses of the final CNN model in stage3 with different parameters.

19 Figure. 7. The comparison of 2D segmentation results with different parameters: $k=0$ with WCE

1 loss (a), k=3 with WCE loss (b), k=3 with VWCE loss (c). Contours in green and red correspond to
2 ground truths and segmentation results respectively.

3 Figure. 8. DSC of each case in the testing dataset with different parameters. The index of images is
4 ranked according to the volume of renal tumors.

5 Figure. 9. The comparison of the results from three testing images obtained by different methods:
6 3D ground truth (a), SDI (b), Constrained-CNN(c), the proposed method (d) and fully-supervised
7 method (e). Contours in green and red correspond to ground truth and segmentation results
8 respectively.

Figures

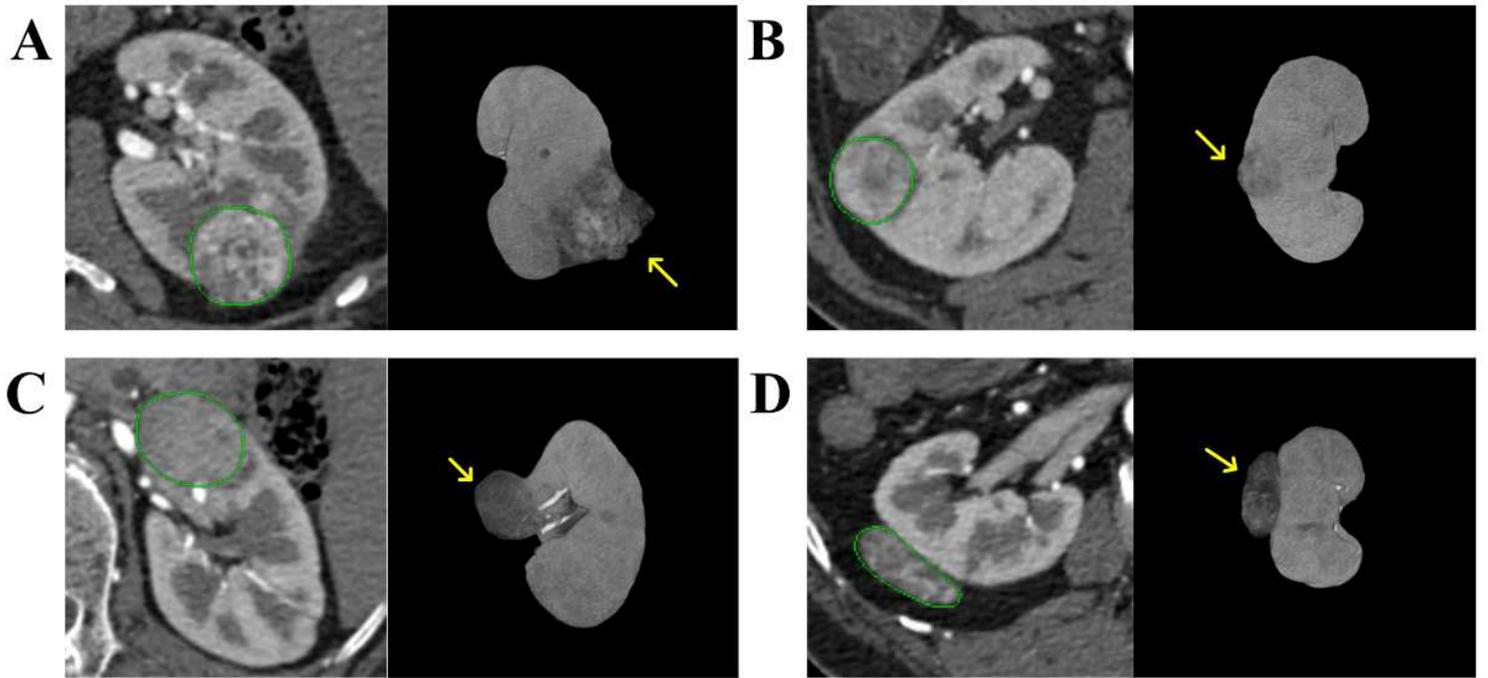


Figure 1

Four contrast-enhanced CT images of different pathological renal tumors. The tumors are marked by yellow arrows in 3D views. The manual contours of the renal tumors delineated by a radiologist are displayed in 2D slices. The pathological subtypes of the renal tumors are clear cell renal cell carcinoma (RCC) in (a) and (b), chromophobe RCC in (c) and angiomyolipoma in (d).

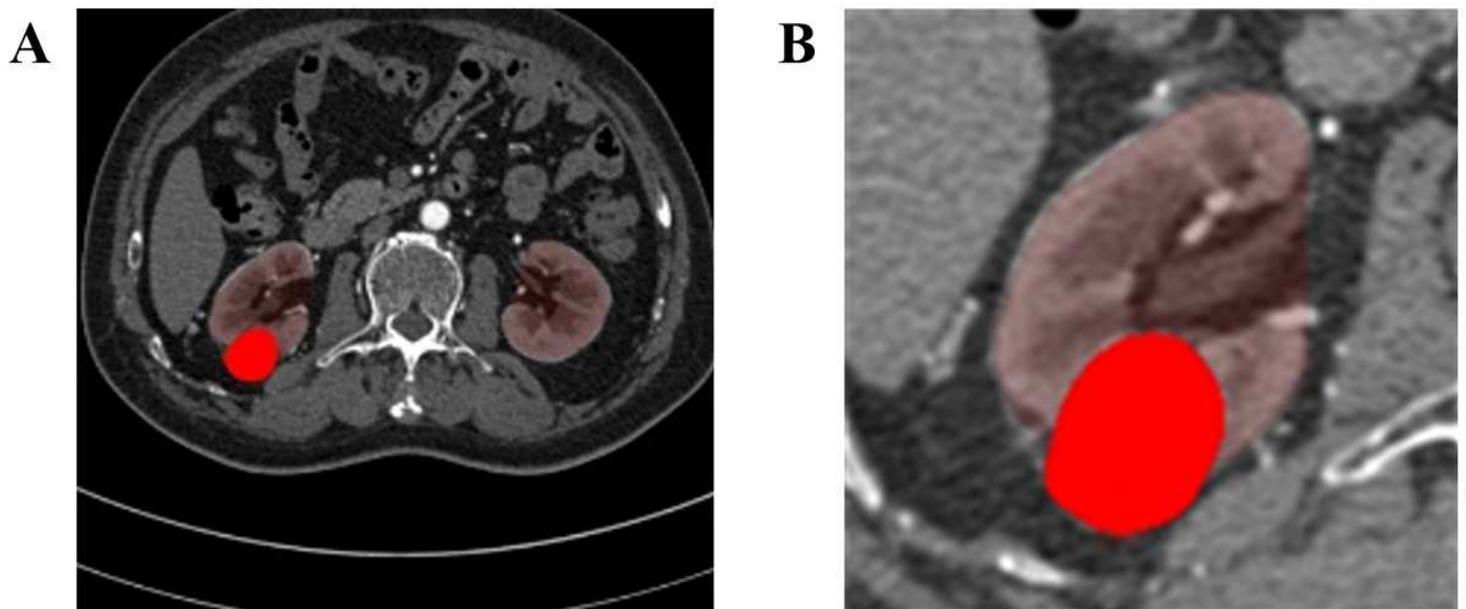


Figure 2

(a) The original image with labeled kidney and renal tumor. The region in red represents renal tumor. (b) The cropped original image with the label for renal tumor segmentation.

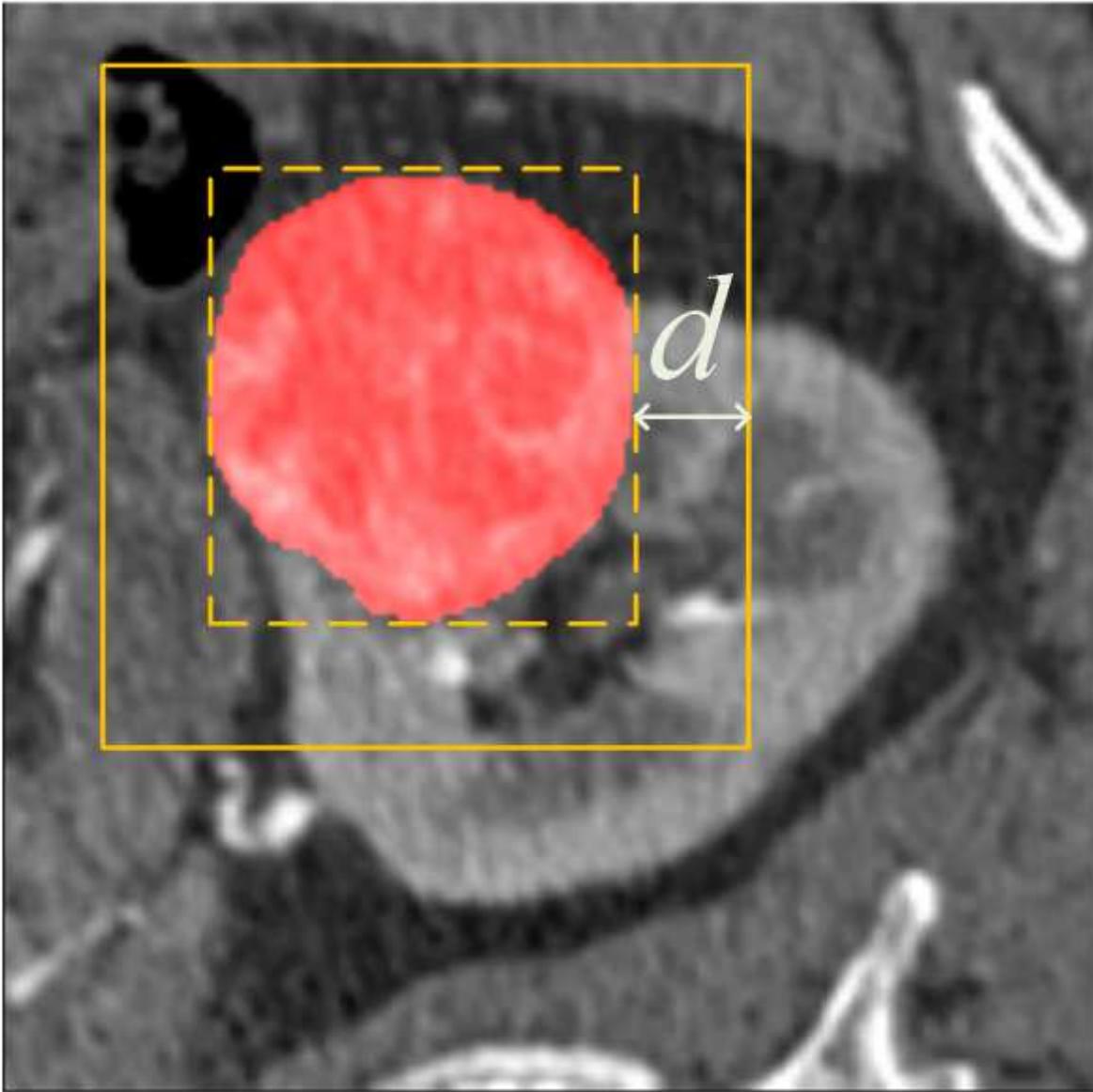


Figure 3

The bounding box with margin d is defined as weak annotations according to the label of renal tumors.

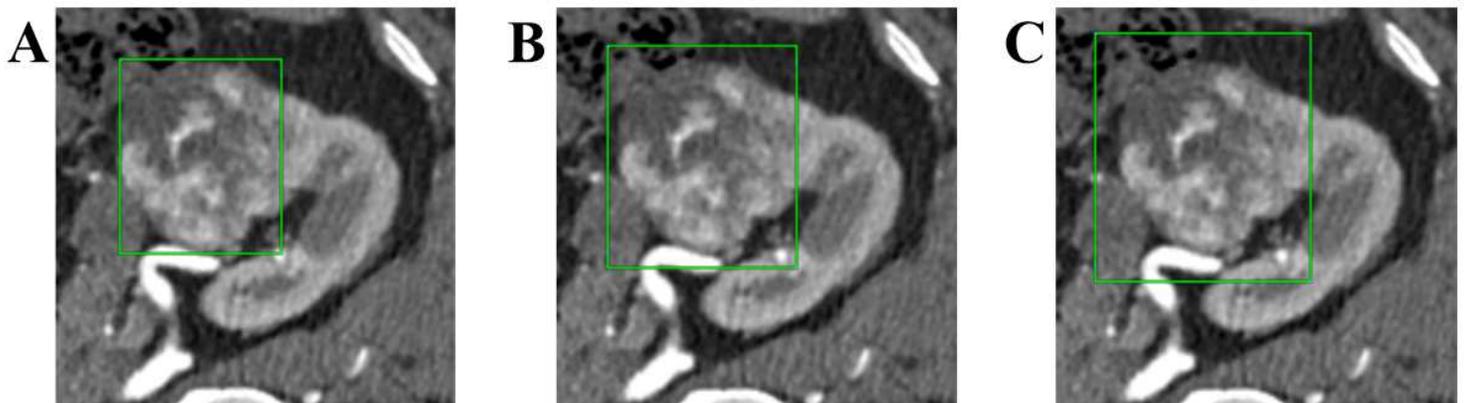


Figure 4

Comparison of bounding boxes with different margins. The 2D image is the maximum slice. Contours in green correspond to bounding boxes.

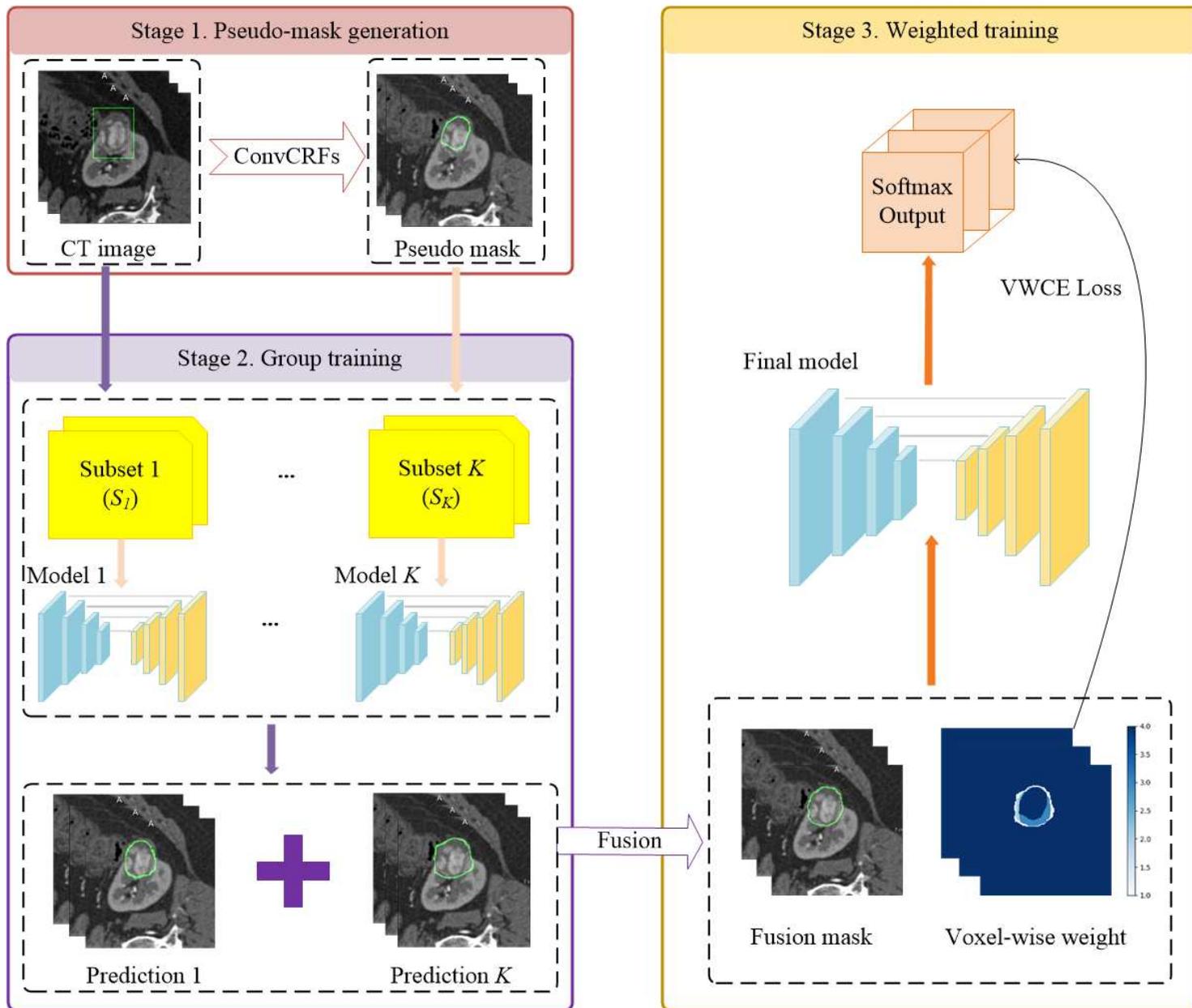


Figure 5

An overview of the proposed weakly-supervised method.

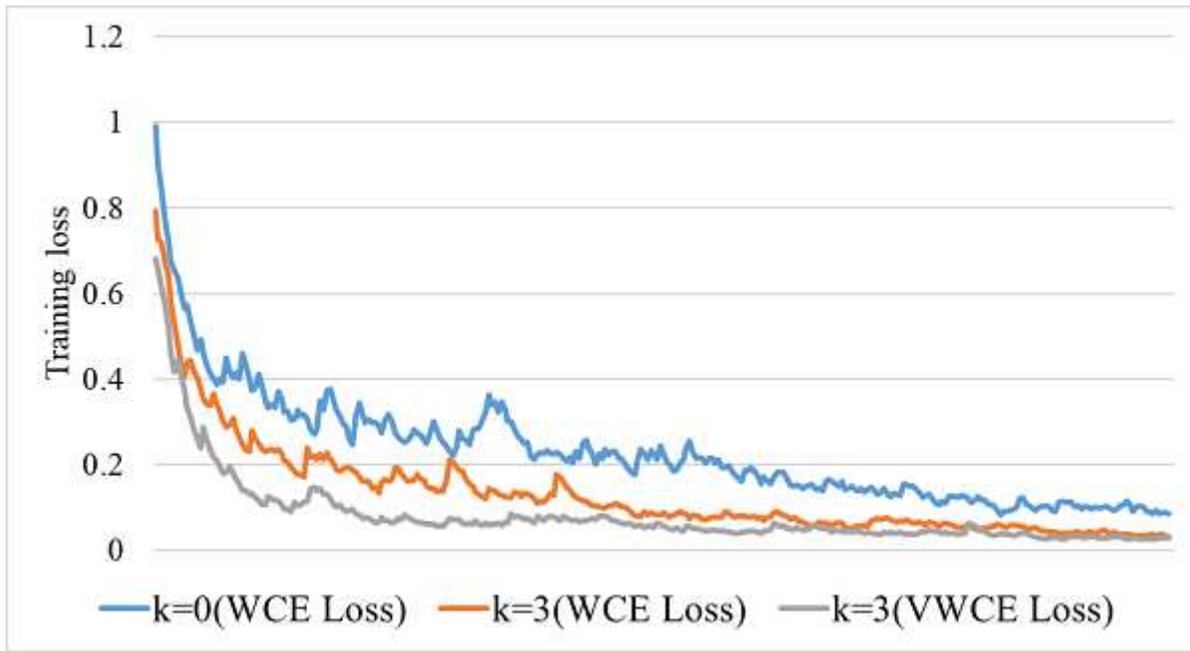


Figure 6

Training losses of the final CNN model in stage3 with different parameters.

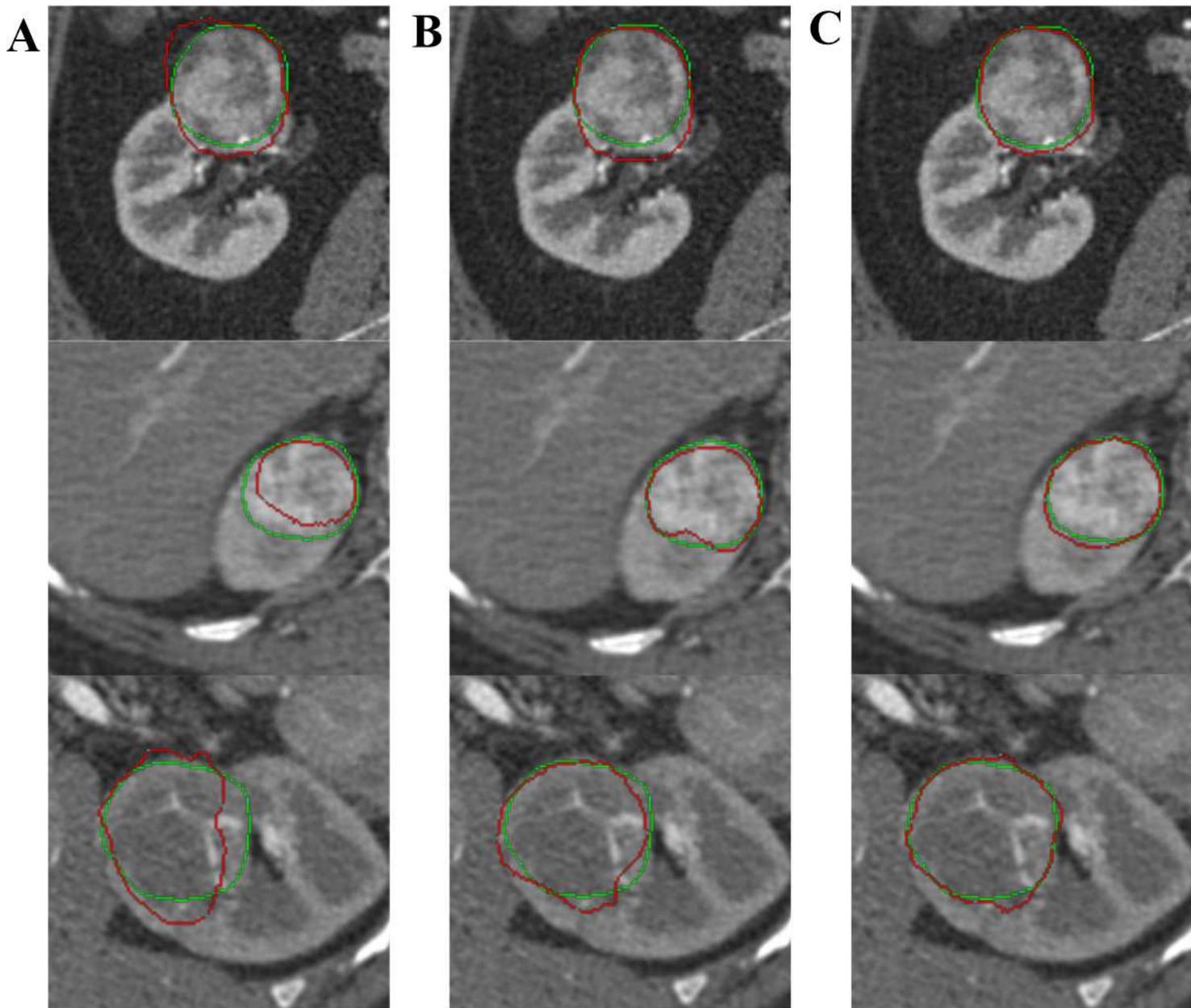


Figure 7

The comparison of 2D segmentation results with different parameters: $k=0$ with WCE loss (a), $k=3$ with WCE loss (b), $k=3$ with VWCE loss (c). Contours in green and red correspond to ground truths and segmentation results respectively.

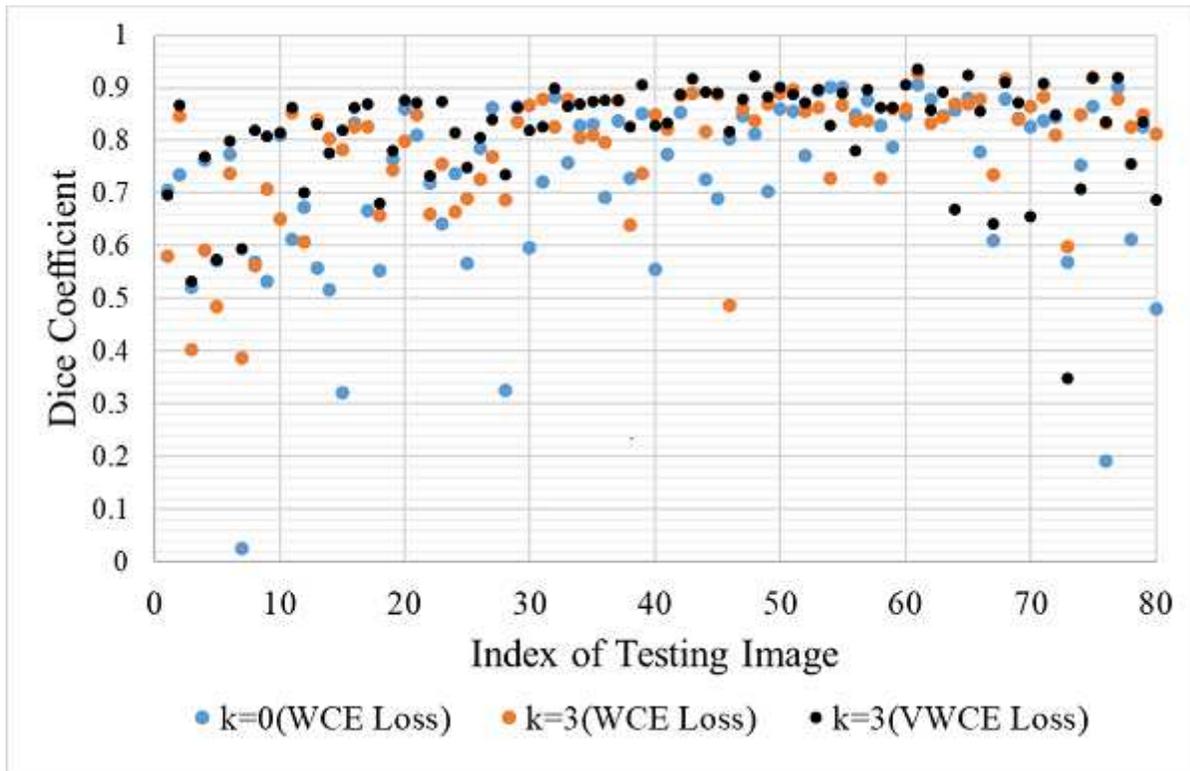


Figure 8

DSC of each case in the testing dataset with different parameters. The index of images is ranked according to the volume of renal tumors.

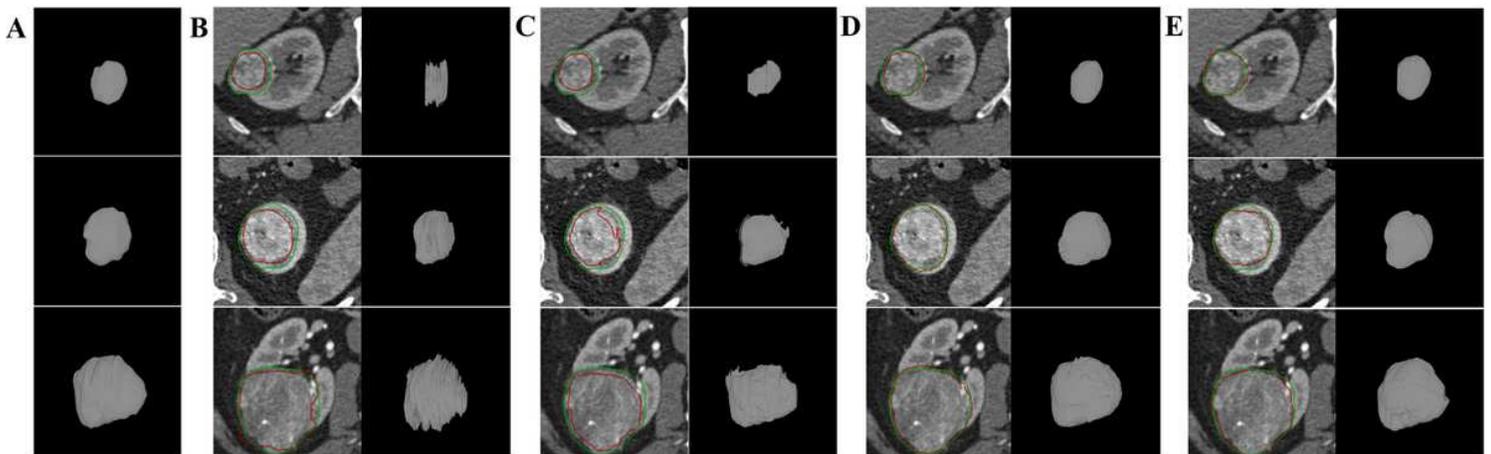


Figure 9

The comparison of the results from three testing images obtained by different methods: 3D ground truth (a), SDI (b), Constrained-CNN(c), the proposed method (d) and fully-supervised method (e). Contours in green and red correspond to ground truth and segmentation results respectively.