

# 3D RoI-aware U-Net for Accurate and Efficient Colorectal Tumor Segmentation

Yi-Jie Huang<sup>1,2</sup>, Qi Dou<sup>3</sup>, Zi-Xian Wang<sup>4</sup>, Li-Zhi Liu<sup>4</sup>, Ying Jin<sup>4</sup>, Chao-Feng Li<sup>4</sup>,  
Lisheng Wang<sup>1\*</sup>, Hao Chen<sup>2\*</sup>, Rui-Hua Xu<sup>4\*</sup>

<sup>1</sup>Institute of Image Processing and Pattern Recognition, Department of Automation, Shanghai Jiao Tong University, China

<sup>2</sup>ImSight Medical Technology Co. Ltd., China

<sup>3</sup>Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong

<sup>4</sup>Sun Yat-sen University Cancer Center; State Key Laboratory of Oncology in South China;  
Collaborative Innovation Center for Cancer Medicine, Guangzhou, China

## Abstract

Segmentation of colorectal cancerous regions from 3D Magnetic Resonance (MR) images is a crucial procedure for radiotherapy which conventionally requires accurate delineation of tumour boundaries at an expense of labor, time and reproducibility. While deep learning based methods serve good baselines in 3D image segmentation tasks, small applicable patch size limits effective receptive field and degrades segmentation performance. In addition, Regions of interest (RoIs) localization from large whole volume 3D images serves as a preceding operation that brings about multiple benefits in terms of speed, target completeness, reduction of false positives. Distinct from sliding window or non-joint localization-segmentation based models, we propose a novel multi-task framework referred to as 3D RoI-aware U-Net (3D RU-Net), for RoI localization and in-region segmentation where the two tasks share one backbone encoder network. With the region proposals from the encoder, we crop multi-level RoI in-region features from the encoder to form a GPU memory-efficient decoder for detail-preserving segmentation and therefore enlarged applicable volume size and effective receptive field. To effectively train the model, we designed a Dice formulated loss function for the global-to-local multi-task learning procedure. Based on the efficiency gains demonstrated by the proposed method, we went on to ensemble models with different receptive fields to achieve even higher performance costing minor extra computational expensiveness. Extensive experiments were subsequently conducted on 64 cancerous cases with a four-fold cross-validation, and the results showed significant superiority in terms of accuracy and efficiency over conventional state-of-the-art frameworks. In conclusion, the proposed method has a huge potential for extension to other 3D object segmentation tasks from medical images due to its inherent generalizability. The code for the proposed method is publicly available.

3D CNN, region of interest, multi-task learning, tumor segmentation, colorectal cancer.

## 1 Introduction

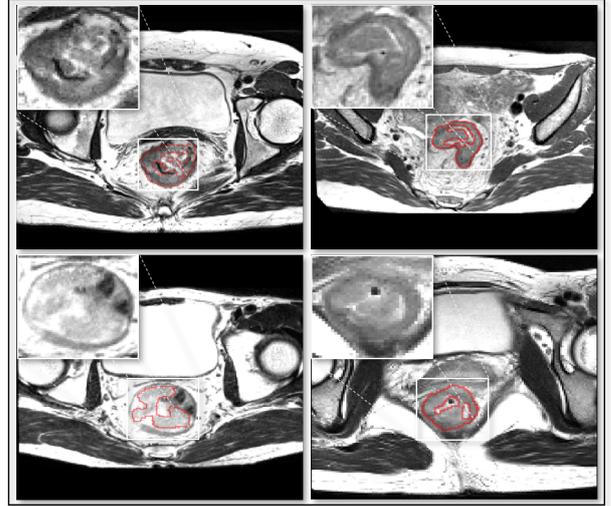


Figure 1: Typical examples of MR slices with colorectal cancer. The cancer regions are delineated with red lines and zoomed in for clear illustration. Clearly, the target areas cannot be well separated by intensity clipping, shape models or positional priors.

Colorectal cancer strikes more than 1.4 million people and accounts for 694,000 deaths globally in 2012 [1]. It is more common in developed countries, for example, in the USA, colorectal cancer is the second leading cause of cancer-related mortalities [2]. In current clinical routine of radiotherapy, colorectal cancer regions are manually recognized and delineated from volumetric images acquired by magnetic resonance (MR) imaging for treatment including surgery and radiation therapy. However, this procedure is laborious, time-consuming and observer-dependent, thus suffers from tedious effort and limited reproducibility. Therefore, automatic colorectal tumor detection and segmentation methods are highly demanded to improve the clinical routine.

Such demand defines a task of automatic detection and segmentation of the targets from whole 3D image volumes. Compared to processing manually selected RoI

patches, the superiority of being fully automatic simplifies the workflow, excludes manual intervention and enables fast processing of large amounts of image volumes. Taking initial works based on super-voxel clustering [3,4] one step further, deep learning based methods dominate the state-of-the-art of detection and segmentation field. However, deep learning based methods for this task are challenged by following factors: weak intensity specificity, absence of shape characteristic, lacking positional priors (as is illustrated in Fig. 1), class imbalance and long processing time of existing methodologies under inferior GPU or CPU-only deployment environments.

Apart from aforementioned challenges, a vital 3D image specific problem is not fully tackled by the community. Among existing methods for fully automatic image segmentation [5–12], though a plausible performance can be achieved by utilizing multi-level features (*e.g.* use skip connections) to gather fine grained details that are lost in the down-sampling process, the merit of maintaining a global understanding represented by deep features with large receptive field is not fully enjoyed due to patch size limitation of GPU memory. As is supported by many researches for 2D image processing, *e.g.*, dilated convolutions [13] and pyramid pooling schemes [14], enlarging receptive fields enables wide-range context utilization and makes further performance breakthroughs. In medical applications, global understanding is even more important since that the targets and the background are highly correlated.

Generally, existing methods for lesion detection and segmentation from 3D images can be divided into part based models and non-joint localization-segmentation based models.

Initially, as naive practices, part based FCNs learn from local parts of 2D slices [7, 15, 16], 2.5D slices [17, 18] or small 3D patches [10, 19] and perform (often overlapped) part-sliding for whole volume inference, which is slow and prone to false positives and target incompleteness related failures. More importantly, part based methods suffer from limited effective receptive fields. V-Net [9], for example, claimed  $551 \times 551 \times 551$  designed receptive field but used  $64 \times 128 \times 128$  patch sliding scheme, making the large designed receptive field not fully effective. To enlarge the effective receptive field under current part based frameworks, Crossbar-Net [20] proposed to train segmentation networks using non-squared patches with different aspect ratios to add more global contexts to local details.

More recently, trends highlight potential accuracy and speed benefits of adding RoI localization modules prior to FCNs. As a common practice, the RoI localization modules are individually designed as a standalone part of a pipeline. Conventionally, RoIs are localized using prior knowledge such as multi-atlas registration, which is often used to localize normal organs [21, 22]. Apart from their inappropriateness for lesion localization, they are relatively slow. As is reported in [23], registration takes at least 20 seconds per patient using GPUs and typically tens of minutes per patient using CPUs. Learning based RoI localization decouples RoI localization from prior knowledge [24–28]. Some of the related practices [24, 29] ex-

tract region proposals using external modules such as Selective Search [30] or Multiscale Combinatorial Grouping (MCG) [31], which are also well-known speed bottlenecks as is pointed out in [32] and replacing them with RPN accelerated a network from 0.5 fps to 5 fps. Later works adopt light CNN models such as 2D CNNs for RoI localization and 3D FCNs for in-region segmentation [27, 33, 34]. Compared to part based methods, these works tackle the tasks in more graceful manners. Still, using a standalone FCN for RoI segmentation requires repeated extraction of low-level features without possible feature sharing, yet feature sharing is reported in [35] to produce 213X acceleration for object detection, given large numbers of target candidates. Nonetheless, using a patch-based FCN for RoI segmentation leaves the problem of limited effective receptive fields unsolved.

As a promising development, joint RoI localization-segmentation models such as Multi-task Network Cascades (MNC) [36] and Mask R-CNN [37] further eliminate redundant feature extraction and achieve better speed and accuracy by sharing a backbone network across the sub-nets for region proposal, region classification and in-region segmentation. Mask R-CNN employs Feature Pyramid Network (FPN) [38], which is encoder-decoder-skip connection formulated, and used scale-specific feature maps for better segmentation details. An apparent drawback of Mask R-CNN is that using scale-specific feature maps and RoIAlign’s bin-fitting scheme for segmentation are still detail-losing, though better than a non-FPN version; To tackle this issue, PA-Net [39] added another bottom-up path for better segmentation detail, which is even more costly for a 3D application. Another drawback of direct extending it to 3D lies on the need of forming anchor boxes defied by additional aspect ratios along the Z axis. Fitting a small amount of 3D objects to more anchor boxes is prone to bad-shaped bounding box prediction.

Apart from the way whole volume predictions are generated, recent works propose some strategies to further boost the performance of volumetric tasks. Firstly, V-Net [9] adopts parameter-free Dice coefficient [40] loss to harness the class-imbalance issue. Secondly, inspired by the success of multi-task learning [41, 42], Deep Contour-aware Networks (DCAN) [43] and Boundary-aware FCN [44] employ contour-aware loss functions for better discrimination between boundaries and the background. In addition, Multilevel Contextual 3D CNNs [45], DeepMedic [46], Orchestral Fully Convolutional Networks (OFCNs) [47] and Hybrid Loss guided Fully Convolutional Networks (HL-FCNs) [48] adopt model ensemble for better robustness.

A part based initial work to automatically segment colorectal cancer regions was published in ISBI [48]. As a step further, in this paper, we propose a novel joint RoI localization-segmentation framework named as 3D RoI-aware U-Net (3D RU-Net) to enjoy the benefits of fast RoI localization, target completeness and large effective receptive field of joint detection-segmentation frameworks while maintaining the easy-to-train and detail-preserving merits of popular end-to-end and volume-to-volume seg-

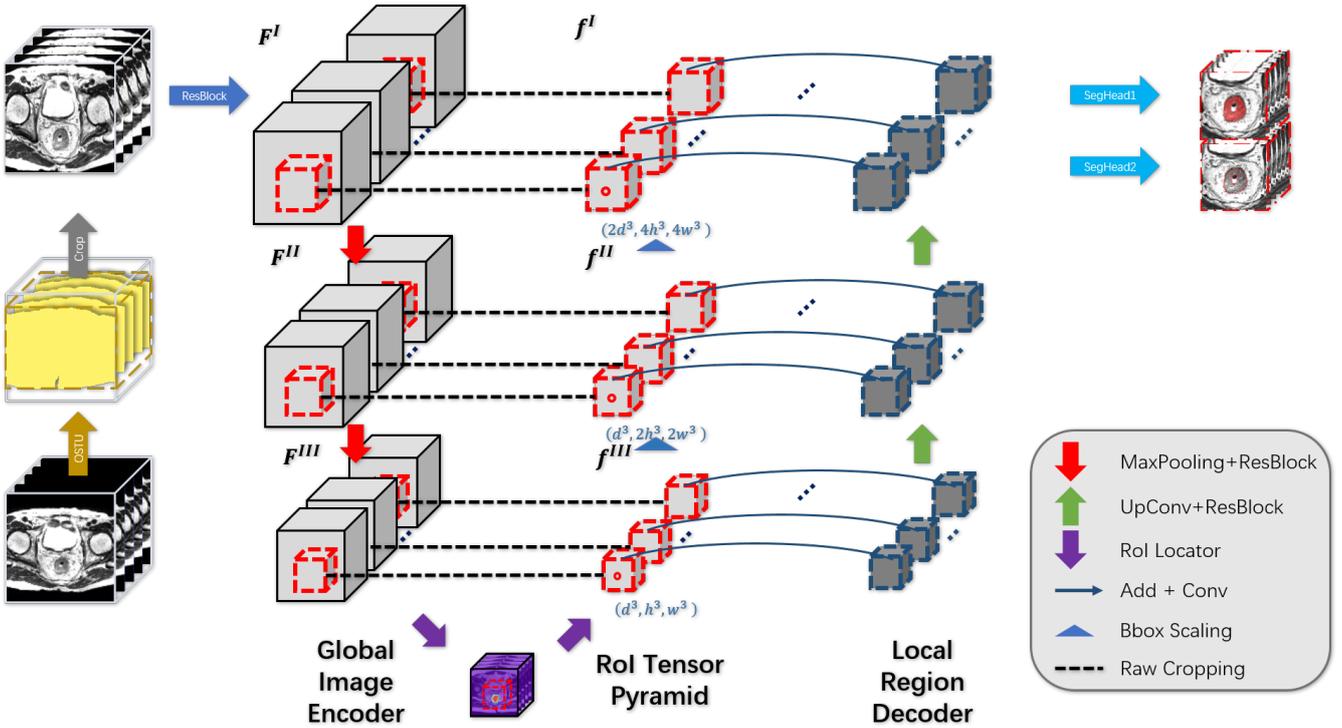


Figure 2: The illustration of 3D RU-Net. The network consists of the Global Image Encoder, the RoI Tensor Pyramid and the Local Region Decoder. A bounding box is predicted using feature maps  $\mathbf{F}^{III}$  and are extended as a Bounding Box Pyramid, then the corresponding RoI Tensor Pyramid ( $f^I, f^{II}, f^{III}$ ) is extracted from ( $F^I, F^{II}, F^{III}$ ) and memory-efficient multilevel feature fusion for in-region segmentation is performed in the decoder stage.

mentation methods. To effectively train the model, we design a hybrid loss function to help the network both handle small objects in big volumes and focus on accurately recognizing ambient borders in local RoIs, and additionally adopt low-cost multi-receptive field ensemble strategy for better robustness. Experiments conducted on 64 acquired scans demonstrated the efficacy of our method and ablation studies validate the contribution gain of each component from our framework.

Our main contributions are summarized as follows:

1. We propose a 3D joint RoI localization-segmentation framework with a shared Global Image Encoder for global-understanding based RoI localization, and a Local Region Decoder working on pyramid-designed in-region features for RoI segmentation. This design enables fast and memory efficient detail-preserving whole volume segmentation with full use of large receptive fields compared to its competing counterparts.
2. Considering automatic class rebalancing and better boundary discrimination, we propose a Dice formulated global-to-local multi-task hybrid loss (MHL) function to further improve the accuracy. Additionally, the accelerated framework encourages us to employ a multiple receptive field model ensemble strategy to suppress the false positives and refine the boundary details at an acceptable speed cost.
3. Extensive experiments on the acquired dataset

proved the efficacy of our proposed framework. Furthermore, our method is inherently general and can be applied in other similar applications.

The remainder of this paper is organized as follows. We describe our method in Section II and report the experimental results in Section III. Section IV further discusses some insights as well as issues of the proposed method. The conclusions are drawn in Section V.

## 2 Methodology

In this section, to address slow prediction and limited effective receptive field issues of non-joint models along with detail-losing and bad bounding box issues of joint models discussed in Section 1, we propose a framework to effectively localize and segment colorectal tumors from whole volume 3D images.

### 2.1 Construction of 3D RU-Net

The proposed 3D RU-Net architecture is illustrated in Fig. 2. We input whole image volumes to Global Image Encoder for multi-level feature encoding, employ an encoder-only RoI locator for RoI localization, crop in-region feature tensors from multi-scale feature maps using RoI Pyramid Layer, and design a Local Region Decoder sub-network to perform multi-level feature fusion for high-resolution cancerous tissue segmentation.

### 2.1.1 Global Image Encoder

Due to limited GPU memory of commonly used devices and dramatically increased parameters of 3D convolution kernels, it's essential to carefully design the 3D backbone feature extractor to avoid GPU memory overflow and overfitting.

Instead of constructing a complete 3D version of encoder-decoder architecture like 3D FPN, or directly extending popular backbones [49–51] to 3D, a compact encoder-only network named the Global Image Encoder is constructed to process whole volume images rather than dealing with context-limited small parts as common practices do. Specifically, the encoder employs a stack of ResBlocks [50] and MaxPooling layers to encode whole volume images. Each Residual Block has three convolutional layers, three Instance Normalization Layers [52], three ReLU layers and a Skip Connection for better gradient flowing. The Instance Normalization Layer is used for better robustness given  $batchsize = 1$  in 3D segmentation tasks.

### 2.1.2 RoI Locator

The RoI Locator is a template where any method that employs encoder-only backbones for target detection can be employed. Due to aspect ratio diversity of number-limited training samples, learning accurate bounding box regression can be difficult. For this specific 3D semantic segmentation task, we recommend taking full advantage of available voxel-level masks as is discussed below for simplicity and more robust bounding box prediction.

Specifically, we avoid degrading voxel-wise labels to object-wise labels to learn anchor fitting. Instead, the locator is designed as a module taking feature map  $F^{III}$  as input, consisting of a convolutional layer with kernel size 1 and *Sigmoid* activation function. This module is trained to predict down-sampled segmentation masks from global images. To tackle the extremely imbalanced foreground-to-background ratio, instead of partial sampling, i.e. sampling a fixed proportion of foreground and background or employing OHEM [53], the locator is trained towards Dice loss, which will be introduced in subsection 2.2. Then we perform a fast 3D connectivity analysis to compute desired bounding boxes formulated as  $Bbox^{III} = (z^3, y^3, x^3, d^3, h^3, w^3)$  where  $(z^3, y^3, x^3)$  denotes the starting coordinates and  $(d^3, h^3, w^3)$  denotes depth, height and width of  $Bbox^{III}$  in feature map  $F^{III}$ .

### 2.1.3 RoI Pyramid Layer

As is illustrated in Fig. 2, we propose a novel layer named RoI Pyramid Layer. Instead of bin-fitting the RoI tensor cropped from a manually selected single-scale feature map, in this paper, we propose to extract a group of raw multi-level feature tensors from each feature scale named as RoI Tensor Pyramid for full utilization of multi-level features and better mask details.

To extract an RoI Tensor Pyramid for a detected target, we first construct a Bounding Box Pyramid ( $Bbox^I, Bbox^{II}, Bbox^{III}$ ) from a given bounding box

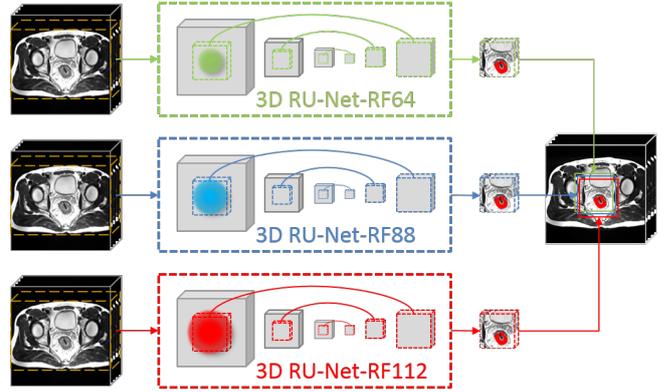


Figure 3: 3D RU-Net-RF64, 3D RU-Net-RF88 and 3D RU-Net-RF112 are of different dilation rates. The green, blue and red spheres indicate receptive fields of  $26 \times 64 \times 64$ ,  $26 \times 88 \times 88$  and  $26 \times 112 \times 112$ , respectively. In the output end, their predictions are averaged.

$Bbox^{III} = (z^3, y^3, x^3, d^3, h^3, w^3)$ . Specifically, Bounding Box Pyramid is computed iteratively following Bbox Scaling criterion listed below:

$$Bbox^{i-1} = (z^i \times s_z^i, y^i \times s_y^i, x^i \times s_x^i, d^i \times s_z^i, h^i \times s_y^i, w^i \times s_x^i) \quad (1)$$

where  $(s_z^i, s_y^i, s_x^i)$  denotes the stride configuration of  $MaxPooling^i$  layer. Given the Bounding Box Pyramid ( $Bbox^I, Bbox^{II}, Bbox^{III}$ ), we crop raw RoI Tensor Pyramid ( $f^I, f^{II}, f^{III}$ ) from whole volume feature maps  $F^I, F^{II}$  and  $F^{III}$  without applying any bin-fitting operation and form an RoI Tensor Pyramid for posterior Local Region Decoder branch.

### 2.1.4 Local Region Decoder

Given a RoI Tensor Pyramid, we construct a sub-network for in-region segmentation named as Local Region Decoder by applying successful multilevel feature fusion mechanism. The construction of the decoder is more or less symmetrical to the encoder part with skip connections to fuse feature maps of corresponding scales, while the beneficial difference lies on much smaller sizes of the decoder branch's feature tensors. Since no shape distortion or scale normalization is included in the RoI Pyramid Layer, this module restores the spatial dimension of the RoI region without losing details. The same set of decoder weights is used to iteratively process different RoIs if multiple RoIs are localized.

## 2.2 Dice-based Multi-task Hybrid Loss Function

In multi-task learning practices, each task faces different challenges. In our case, the Global Image Encoder mainly suffers from class imbalance issue, while the Local Region Decoder has to focus on the exact boundaries of the target regions. Thus we propose a Dice-based multi-task loss (MHL) function to effectively learn these tasks.

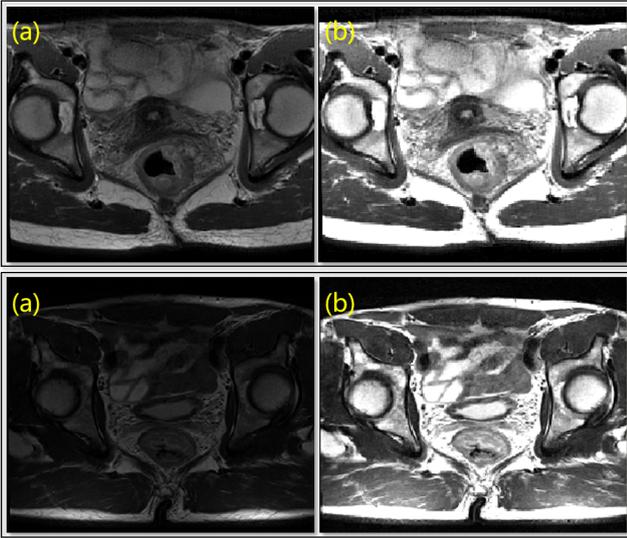


Figure 4: Examples of: (a) Original Images (b) Normalized Images. The intensity of homogeneous tissues from images acquired under different imaging configurations are normalized to identical ranges.

### 2.2.1 Dice Loss Formulation

Inspired by the success of [9], we apply Dice loss function to formulate the optimization objective, since it serves as an effective hyper-parameter free class balancer to help the network learn objects of small size and weak saliency. The Dice loss is defined as:

$$L_d(P, G) = 1 - 2 \times \frac{\sum_{i=1}^N p_i g_i + \epsilon}{\sum_{i=1}^N p_i + \sum_{i=1}^N g_i + \epsilon} \quad (2)$$

where the sums are computed over the  $N$  voxels of the predicted volume  $p_i \in P$  and the ground truth volume  $g_i \in G$ .  $\epsilon$  is a minimal smoothness term that avoids division by 0 and is set as  $10^{-4}$ . In the optimization stage, the Dice loss is minimized by gradient descent using the following derivate:

$$\frac{\partial L_d(P, G)}{\partial p_k} = -2 \times \frac{\sum_{i=1}^N p_i g_i - g_k \sum_{i=1}^N (p_i + g_i)}{[\sum_{i=1}^N (p_i + g_i)]^2} \quad (3)$$

### 2.2.2 Dice Loss for Global Localization

To tackle the class imbalance issue of the global image RoI localization task, we employ the aforementioned Dice loss:

$$L_{global} = L_d(P_{global}, G_{global}) \quad (4)$$

where  $P_{global}$  and  $G_{global}$  denotes predictions of the localization top and down-sampled annotations.

### 2.2.3 Dice-based Contour-aware Loss for Local Segmentation

Compared to the localization task, the in segmentation branch needs multiple constraints to acquire better boundary-sensitive segmentation results. In semantic segmentation practices, the ambiguous borders are the most

difficult to learn but learned with insufficient attention. Borrowing the insight of previous exploration of adding an auxiliary contour-aware side task [43], we further formulate the side task using Dice loss to help it tackle the extreme sparsity of contour labels in 3D space. Practically we add an extra  $1 \times 1 \times 1$  convolutional layer activated by *Sigmoid* function at the output terminal of the segmentation branch to predict the contour voxels, trained in parallel with the region segmentation task. Taking the side task into account, the loss function of the segmentation branch  $L_{local}$  is denoted as following by summarizing the weighted losses:

$$L_{local} = L_d(P_{region}, G_{region}) + \lambda_c L_d(P_{contour}, G_{contour}) \quad (5)$$

where  $\lambda_c = 0.5$ , denoting the auxiliary task weight to ensure that the region segmentation task dominates while other tasks take effects.

Finally, the overall loss function is:

$$L = L_{global} + L_{local} + \beta \|W\|_2^2 \quad (6)$$

where  $\beta = 10^{-4}$  denotes the balance of weight decay term and  $W$  denotes the parameters of the whole network.

## 2.3 Multiple Receptive Field Model Ensemble

Due to the limited accuracy of single models, ensemble of multiple models is considered as an effective practice to perform robust inference, and is widely employed in practical cases, at a cost of computational expensiveness.

Encouraged by the dramatically accelerated framework, in this paper, we propose to employ multiple receptive field model ensemble strategy by fusing models of identical structure but with different receptive field settings. This is a generalization to the multi-resolution strategy proposed in [48] that applies identical receptive field to images with different spatial resolutions, which is actually formulating different spatial receptive fields. Such generalization gets rid of detail-losing down-sampling and allows each model contribute to boundary details equally.

In detail, as is illustrated in TABLE 1, we first construct an original 3D R-U-Net of receptive field  $26 \times 64 \times 64$ , named 3D RU-Net-RF64. Next, we tune the dilation rate of *ResBlock3* as 2, enlarging the receptive field to  $26 \times 88 \times 88$  and formulate 3D RU-Net-RF88; We further tune the dilation rates of *ResBlock2*, *ResBlock3* and *ResBlock4* as 2 and construct a 3D R-U-Net of receptive field  $26 \times 112 \times 112$  named 3D RU-Net-RF112.

In the inference stage, as is shown in Fig. 3, three networks' outputs are averaged to generate the final prediction. Major voting produces similar scores and is therefore not discussed.

Part Name	Input Layer	Module Name	Kernel	Out Channels	Receptive Field 1	Receptive Field 2	Receptive Field 3
Encoder	Image	ResBlock1	$1 \times 3 \times 3$	48	$1 \times 7 \times 7$	$1 \times 7 \times 7$	$1 \times 7 \times 7$
	ResBlock1	MaxPooling1	$1 \times 1/2 \times 1/2$	48	-	-	-
	MaxPooling1	ResBlock2	$3 \times 3 \times 3$	96	$7 \times 20 \times 20$	$7 \times 20 \times 20$	$7 \times 34 \times 34$
	ResBlock2	MaxPooling2	$1/2 \times 1/2 \times 1/2$	96	-	-	-
	MaxPooling2	ResBlock3	$3 \times 3 \times 3$	192	$20 \times 46 \times 46$	$20 \times 70 \times 70$	$20 \times 46 \times 46$
	ResBlock3	Locator (sigmoid)	$1 \times 1 \times 1$	1	$20 \times 46 \times 46$	$20 \times 70 \times 70$	$20 \times 82 \times 82$
RoI Pyramid Layer	Locator,ResBlock1	RoI Tensor I	-	48	$1 \times 7 \times 7$	$1 \times 7 \times 7$	$1 \times 7 \times 7$
	Locator,ResBlock2	RoI Tensor II	-	96	$7 \times 20 \times 20$	$7 \times 20 \times 20$	$7 \times 34 \times 34$
	Locator,ResBlock3	RoI Tensor III	-	192	$20 \times 46 \times 46$	$20 \times 70 \times 70$	$20 \times 82 \times 82$
Decoder	RoI Tensor III	UpConv1	$2 \times 2 \times 2$	96	-	-	-
	RoI Tensor II,UpConv1	Add1	-	96	-	-	-
	Add1	ResBlock4	$3 \times 3 \times 3$	96	$26 \times 58 \times 58$	$26 \times 82 \times 82$	$26 \times 106 \times 106$
	ResBlock4	UpConv2	$1 \times 2 \times 2$	48	-	-	-
	RoI Tensor I,UpConv	Add2	-	48	-	-	-
	Add2	ResBlock5	$1 \times 3 \times 3$	48	$26 \times 64 \times 64$	$26 \times 88 \times 88$	$26 \times 112 \times 112$
	ResBlock5	SegHead1 (sigmoid)	$1 \times 1 \times 1$	1	$26 \times 64 \times 64$	$26 \times 88 \times 88$	$26 \times 112 \times 112$
	ResBlock5	SegHead2 (sigmoid)	$1 \times 1 \times 1$	1	$26 \times 64 \times 64$	$26 \times 88 \times 88$	$26 \times 112 \times 112$

Table 1: Parameters and connectivity of the network.

## 3 Experiments

### 3.1 Dataset and Preprocessing

#### 3.1.1 Dataset

The dataset contains a total of 64 MR images of the pelvic cavity of T2 modality whose ZYX spacings range from  $3.6 \times 0.31 \times 0.31$  mm to  $4.0 \times 1.0 \times 1.0$  mm. Target areas were labeled voxel-wisely by experienced radiologists, and contour labels were automatically generated from the region labels of one-voxel thickness using erosion and subtraction operations. An 3D image has mostly one and up to two RoIs containing cancerous tissues.

#### 3.1.2 Preprocessing

Different spacing rates are normalized to  $4.0 \times 1.0 \times 1.0$  as the HighRes set. Some part-based methods listed in TABLE 2 employ down-sampled image sets, namely LowRes set of  $4.0 \times 2.0 \times 2.0$  mm spacing and MidRes set of  $4.0 \times 1.5 \times 1.5$  spacing. To normalize the intensities of input images acquired under different imaging configurations and field of views, we perform in-body intensity normalization to exclude the affect of inconsistent body-to-background ratios. By OTSU [54] thresholding, connectivity analysis and closing operation, body masks are extracted as foreground and other voxels are set as background. The mean intensity and standard deviation are computed within the body mask according to following formulas:

$$Mean(X) = \frac{1}{N_{mask}} \sum_{i \in mask} x_i \quad (7)$$

$$std(X) = \sqrt{\frac{1}{N_{mask}} \sum_{i \in mask} (x_i - Mean(X))^2} \quad (8)$$

where  $x_i \in X$  denotes the intensity of a voxel and  $N_{mask}$  denotes the count of mask voxels. Then the image is normalized according to standard normalization criterion.

A few examples of the comparison between original images and intensity-normalized images are illustrated in Fig. 4.

Before feeding the images to the network, we crop the input images according to minimum bounding boxes of the body masks to further reduce the GPU memory footprint. Additionally, in the training stage, we performed on-the-fly data augmentation when feeding training samples. Applied random operations include 0.9X to 1.1X scaling, flipping w.r.t. the X axis, 0.9X to 1.1X intensity jittering, and RoI translation that shifts the RoI center by -50% to 50% width long each axis.

### 3.2 Implementation Details

Our implementation is publicly available at <https://github.com/huangyjhurst/3D-RU-Net>.

#### 3.2.1 Hyper-Parameters

The network’s detailed connectivity and kernel configuration are illustrated in Table 1. Specifically, to fit the anisotropic spacing of the acquired dataset which has larger spacing along Z axis, flat kernels of  $1 \times 3 \times 3$ , pooling rate of  $1 \times 1/2 \times 1/2$  and up-sampling rate of  $1 \times 2 \times 2$  are employed by the input and output blocks, *i.e.* ResBlock1, MaxPooling1, UpConv2, ResBlock5. Initial experiments demonstrate that adding MaxPoolings, ResBlocks or channels does not improve the performance, hence we tune receptive field setting by applying dilated convolution rather than adding layers.

#### 3.2.2 Training Process

The backbone network were initialized using criterion proposed in [55], then pre-trained using our previous work’s patch-wise HL-FCN [48]. We used Adam [56] optimizer at a learning rate of  $10^{-4}$ . The weights of convolution kernels were penalized with  $10^{-4}$  L2 norm for better generalization capability. Then, we first train the RoI locator until evaluation loss no longer decrease, then jointly train the RoI locator and the segmentation branch. In each joint training iteration, we accumulate the losses of the RoI Locator, SegHead1 and SegHead2.

### 3.3 Evaluation Metrics

#### 3.3.1 Dice Similarity Coefficient (DSC)

The Dice similarity coefficient (DSC) measures a general overlap rate that equally assigns significance to recall rate and false positive rate. DSC is denoted as:

$$DSC(P, G) = \frac{2|P \cap G|}{|P| + |G|} \quad (9)$$

where the metric is scored in  $[0,1]$ . Better prediction generates a score closer to 1.0. Since this network is trained towards this metric, DSC is not enough to evaluate the performance.

#### 3.3.2 Voxel-wise Recall Rate

We also employ voxel-wise recall rate to evaluate the recall capability of different methods.

$$Recall = \frac{|P \cap G|}{|G|} \quad (10)$$

#### 3.3.3 Average Symmetric Surface Distance (ASD)

We define the shortest distance of an arbitrary voxel of one volume’s surface to another volume’s surface as:

$$d(a_k, B) = \min_{b_i \in S(B), a_k \in S(A)} \|a_k - b_i\| \quad (11)$$

where  $a_k$  denotes  $k$ th voxel from extracted surface  $S(A)$  of volume  $A$ ,  $b_i$  denotes  $i$ th voxel from extracted surface  $S(B)$  of volume  $B$ , and  $\|\cdot\|$  denotes Euclidean distance. Then the evaluation value is defined as:

$$ASD = \frac{\sum_{p_k \in S(P)} d(p_k, G) + \sum_{g_k \in S(G)} d(g_k, P)}{|S(P)| + |S(G)|} \quad (12)$$

where  $|S(P)|$  and  $|S(G)|$  denote the number of surface voxels.

Specifically, this metric is sensitive to failures such as debris outliers predicted far away from the colon region or complete failure to recall an object. The long distance makes up for the small size of the debris and produce large error penalty. If a failure segmentation has 0 recall rate, its surface distance is set as 50 mm, which is big enough to be a strong penalty.

#### 3.3.4 Average Inference Time

We include average inference time to evaluate speed in the inference stage. Since this metric is decided by the size of the input volume, the standard deviation is not evaluated. The tested methods are all performed on a workstation platform with 2x Xeon E5 CPU (8C16T) @ 2.4 Ghz, 128GB RAM and an NVIDIA Titan Xp GPU with 12GB GPU memory. The code is implemented with PyTorch and the inference speed is evaluated under volatile mode.

#### 3.3.5 Typical GPU Memory Footprint

By analyzing this metric, we describe the GPU memory efficiency of the proposed methods by tracking the total GPU memory footprint given an input volume of typical size  $40 \times 180 \times 320$  voxels.

### 3.4 Results

For evaluation, four-fold cross-validation was conducted on 64 scans and their mean scores are reported in TABLE. 2. Comparison of predicted masks between different methods is illustrated in Fig. 5; Eight volume predictions are illustrated in Fig. 6.

#### 3.4.1 Ablation Studies

Firstly, we conduct a full ablation study to evaluate the contribution of each proposed component, listed in the upper section of TABLE. 2.

Compared to the part based 3D U-Net [8] built with the ResBlocks described in 2.1.1, 3D U-Net-RF64+DL [9] and 3D U-Net-RF64+HL [48] using Dice loss and hybrid loss improved the performance by alleviating the class imbalance problem. Specifically, we acquired  $(d, h, w) = (24, 96, 96)$  patches at a stride of 50% window overlapping for training and predicting and found that despite that details are sacrificed, down-sampling, *i.e.* using MidRes and LowRes image sets, significantly boosts the performance due to enlarged physical receptive fields. As a step further, can we enlarge the receptive field defined by the network’s convolution kernels rather than down-sampling the images for better performance without sacrificing details? Following the criterion stated in 2.3, we tuned dilation rates of 3D U-Net+HL to form 3D U-Net-RF64+HL, 3D U-Net-RF88+HL and 3D U-Net-RF112+HL, with receptive fields of  $26 \times 64 \times 64$ ,  $26 \times 88 \times 88$  and  $26 \times 122 \times 122$ , respectively. The experimental results demonstrate that enlarging receptive field without enlarging patches does not take the performance to the level of down-sampling based methods. These results highlight that the input volume size hindered the receptive field from taking advantage of wide range contexts.

Apparently, the pipeline can be accelerated by employing either of a non-joint or a joint detection-segmentation framework, and most of the false positives can be eliminated as well. To further emphasize the merit of whole volume joint training and cross-module feature sharing enabled by the proposed method named as 3D RU-Net+MHL, a detection-segmentation cascaded model without these properties, namely 3D FCN+3D U-Net, is designed and evaluated. The cascaded model consists of a standalone Global Image Encoder for RoI detection and a full 3D U-Net replacing the Local Region Decoder for in-region segmentation, trained towards Dice loss and Hybrid loss, respectively. We also tuned the 3D U-Net’s receptive field as  $26 \times 64 \times 64$ ,  $26 \times 88 \times 88$  and  $26 \times 122 \times 122$ , and did not notice a significant performance difference ( $< 0.4\%$ ). On the other hand, the proposed method, however, enjoyed a higher Dice score (from

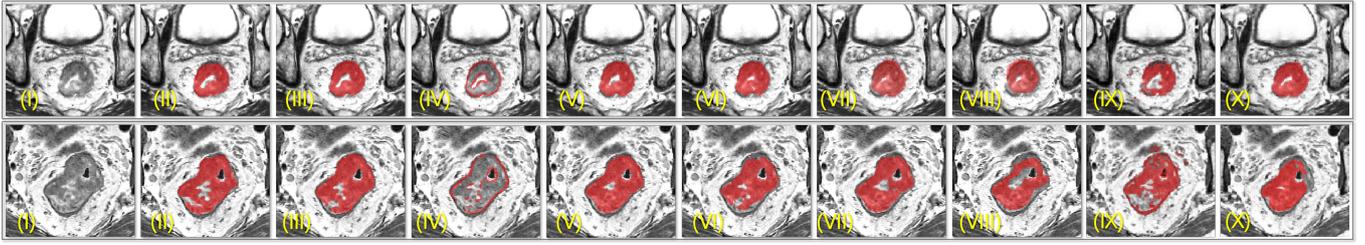


Figure 5: (I) Cancerous region, (II) Expert delineation, (III) Proposed method(predicted regions), (IV) Proposed method (predicted contours) (V) 3D U-Net+DL [9] (Ensemble) (VI) 3D U-Net [8] (VII) 3D FCN+3D U-Net (VIII) 3D Mask R-CNN [37](IX) Super-Voxel clustering [4](X) 2D kU-Net+LSTM [16]

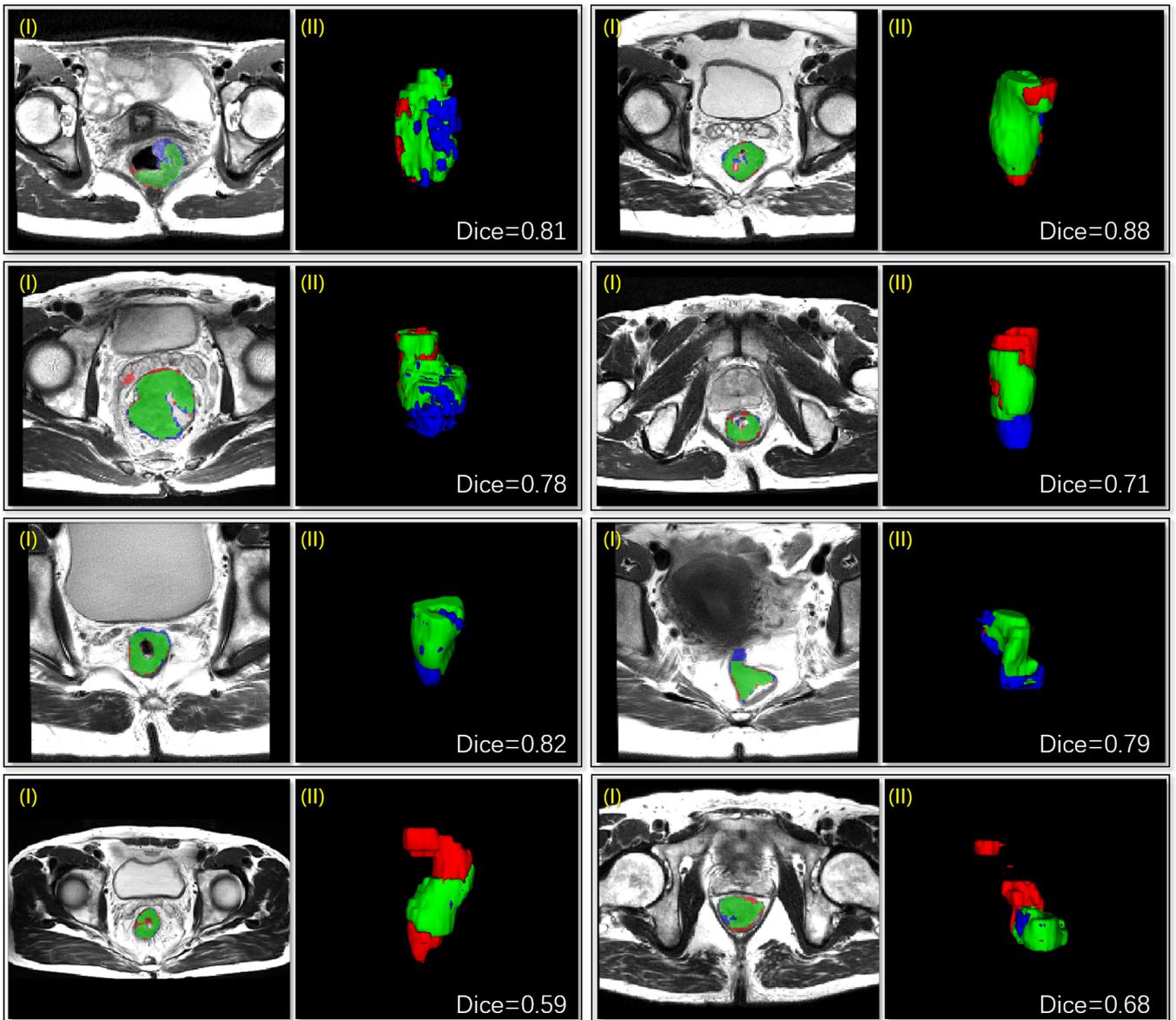


Figure 6: (I) selected 2D slices (II) 3d segmentation masks. Green indicates true positives; Red indicates false positives; Blue indicates false negatives.

Table 2: Ablation studies and comparisons.

Method	DSC[%]	Recall[%]	ASD[mm]	$AvgN_{RoIs}$	Loc/Seg GPU Time[s]	Loc/Seg CPU Time[s]
<b>3D RU-Net+MHL+Ensemble</b>	<b>75.5±10.7</b>	<b>77.8±14.8</b>	<b>2.45±3.26</b>	-	0.61	38.10
3D RU-Net-RF112+MHL (HighRes)	74.2±10.6	78.6±13.9	3.02±3.55	1.8	0.22/0.01	15.30/0.79
3D RU-Net-RF88+MHL (HighRes)	73.7±10.5	75.4±14.8	2.83±3.57	2.2	0.17/0.01	10.44/0.45
3D RU-Net-RF64+MHL (HighRes)	72.7±12.5	76.2±17.2	2.62±3.05	2.7	0.15/0.01	8.98/0.34
3D FCN+3D U-Net+HL+Ensemble	73.4±11.6	78.5±15.3	3.10±3.71	-	0.72	44.45
3D FCN+3D U-Net-RF122+HL (HighRes)	72.0±12.1	78.2±15.5	3.41±4.14	2.2	0.22/0.03	15.30/1.77
3D FCN+3D U-Net-RF88+HL (HighRes)	71.6±12.5	75.9±16.6	3.42±3.73	2.6	0.17/0.02	10.44/1.12
3D FCN+3D U-Net-RF64+HL (HighRes)	71.7±11.9	79.1±14.9	3.56±4.07	3.1	0.15/0.02	8.98/0.96
3D U-Net-RF64+HL+Ensemble [48]	72.1±13.9	72.2±17.2	3.83±4.95	-	18.11	616.70
3D U-Net-RF64+HL [48] (LowRes)	69.9±12.5	70.2±14.9	3.90±4.43	-	2.25	88.90
3D U-Net-RF64+HL [48] (MidRes)	70.0±14.5	72.1±17.3	5.48±7.06	-	5.60	180.62
3D U-Net-RF64+HL [48] (HighRes)	67.7±18.4	69.2±21.3	10.24±14.59	-	10.26	346.72
3D U-Net-RF112+HL [48] (HighRes)	67.8±13.8	70.5±14.8	12.44±13.70	-	14.62	827.57
3D U-Net-RF88+HL [48] (HighRes)	66.9±17.4	71.8±15.7	14.16±11.54	-	12.47	358.88
3D 3D U-Net-RF64+DL+Ensemble [9]	69.9±13.7	72.4±18.0	4.18±5.89	-	18.11	616.70
3D U-Net-RF64+DL [9] (LowRes)	68.5±13.8	68.5±19.5	4.19±5.75	-	2.25	88.90
3D U-Net-RF64+DL [9] (MidRes)	67.3±15.3	70.2±17.7	5.70±7.31	-	5.60	180.62
3D U-Net-RF64+DL [9] (HighRes)	66.0±18.2	70.9±22.0	10.32±12.11	-	10.26	346.40
3D U-Net-RF64 [8] (HighRes)	61.7±19.2	57.3±23.9	4.26±4.35	-	10.26	346.40
2D U-Net+3D U-Net+Ensemble [34]	72.0±13.6	76.1±18.3	3.86±5.46	-	1.021	79.70
2D U-Net+3D U-Net [34] (LowRes)	70.2±12.4	74.5±15.8	4.11±5.03	5.3	0.15/0.02	7.68/0.53
2D U-Net+3D U-Net [34] (MidRes)	69.1±17.7	73.7±21.0	6.05±9.53	6.1	0.18/0.02	16.95/0.87
2D U-Net+3D U-Net [34] (HighRes)	69.4±14.1	76.2±18.2	6.23±8.77	7.1	0.25/0.03	38.29/1.22
2D kU-Net+BDC-LSTM [16] (HighRes)	69.3±13.1	79.1±16.7	7.81±6.88	-	0.51	39.22
Super-Voxel Clustering [4] (HighRes)	62.6±14.9	60.2±18.2	6.54±5.96	-	-	15.13
3D Mask R-CNN [37] (HighRes)	56.4±19.0	58.5±25.6	7.93±10.33	-	0.55	35.88
3D Mask R-CNN [37] (MidRes)	54.6±17.3	61.0±24.5	9.05±8.53	-	0.32	18.07
3D Mask R-CNN [37] (LowRes)	52.0±16.8	55.0±24.1	7.02±7.24	-	0.24	11.61

Table 3: GPU memory footprint tracking given an input volume of size  $40 \times 180 \times 320$  and RoI size of  $24 \times 96 \times 96$ 

Part Name	Layer Name	Size	GPU Memory Footprint	Part GPU Memory Footprint
Encoder	ResBlock1	9 nodes $\times$ $40 \times 180 \times 320 \times 48$ channels	3796.88 MBytes	6302.047MBytes
	MaxPooling1	1 node $\times$ $40 \times 90 \times 160 \times 48$ channels	105.47 MBytes	
	ResBlock2	9 nodes $\times$ $40 \times 90 \times 160 \times 96$ channels	1898.45 MBytes	
	MaxPooling2	1 node $\times$ $20 \times 45 \times 80 \times 96$ channels	26.37 MBytes	
	ResBlock3	9 nodes $\times$ $20 \times 45 \times 80 \times 192$ channels	474.60 MBytes	
	Locator (sigmoid)	1 node $\times$ $20 \times 45 \times 80 \times 1$ channel	0.27MBytes	
RoI Tensor Pyramid	RoI Tensor1	1 node $\times$ $24 \times 96 \times 96 \times 48$ channels	40.50 MBytes	65.82 MBytes
	RoI Tensor2	1 node $\times$ $24 \times 48 \times 48 \times 96$ channels	20.25 MBytes	
	RoI Tensor3	1 node $\times$ $12 \times 24 \times 24 \times 192$ channels	5.06 MBytes	
Local Region Decoder	UpConv1	1 node $\times$ $24 \times 48 \times 48 \times 96$ channels	20.25 MBytes	669.93 MBytes
	Add1	1 node $\times$ $24 \times 48 \times 48 \times 96$ channels	20.25 MBytes	
	ResBlock4	9 nodes $\times$ $24 \times 48 \times 48 \times 96$ channels	182.25 MBytes	
	UpConv2	1 node $\times$ $24 \times 96 \times 96 \times 48$ channels	40.50 MBytes	
	Add2	1 node $\times$ $24 \times 96 \times 96 \times 48$ channels	40.50 MBytes	
	ResBlock5	9 nodes $\times$ $24 \times 96 \times 96 \times 48$ channels	364.50 MBytes	
	SegHead1 (sigmoid)	1 node $\times$ $24 \times 96 \times 96 \times 1$ channels	0.84 MBytes	
	SegHead2 (sigmoid)	1 node $\times$ $24 \times 96 \times 96 \times 1$ channels	0.84 MBytes	
Standard Decoder	UpConv1	1 node $\times$ $40 \times 90 \times 160 \times 96$ channels	210.94 MBytes	6978.55 MBytes
	Add1	1 node $\times$ $40 \times 90 \times 160 \times 96$ channels	210.94 MBytes	
	ResBlock4	9 nodes $\times$ $40 \times 90 \times 160 \times 96$ channels	1898.45 MBytes	
	UpConv2	1 node $\times$ $40 \times 180 \times 320 \times 48$ channels	421.88 MBytes	
	Add2	1 node $\times$ $40 \times 180 \times 320 \times 48$ channels	421.88 MBytes	
	ResBlock5	9 nodes $\times$ $40 \times 180 \times 320 \times 48$ channels	3796.88 MBytes	
	SegHead1 (sigmoid)	1 node $\times$ $40 \times 180 \times 320 \times 1$ channels	8.79 MBytes	
	SegHead2 (sigmoid)	1 node $\times$ $40 \times 180 \times 320 \times 1$ channels	8.79 MBytes	

72.7% to 74.2%) by enlarging the receptive field, which is a significant performance boost considering that no extra parameter or module is included. Another merit of feature sharing lies on the observation that a joint trained model’s Locator module and Local Region Decoder module share a highly consistent behavior pattern except for detail richness. On the contrary, the cascaded model suffers from more in-region false positives due to the non-joint training scheme and the patch size limit of the receptive field, therefore produced higher recall rates along with more false positives (larger  $AvgN_{RoIs}$ ) and scored lower mean Dice and larger mean ASD.

Nevertheless, the proposed method is significantly faster than part-based methods, which is almost impossible for CPU-only deployment, while our method costs 15 seconds to score 74.2% mean Dice and less than 40 seconds to achieve 75.5% mean Dice, providing faster and more accurate predictions. Furthermore, the proposed Local Region Decoder is 2X faster compared to the cascaded model’s segmentation branch, which makes significant difference when an image carries multiple detected RoIs.

Aforementioned performance gains and speedups are enabled by improved memory efficiency: compared to vanilla 3D U-Net, the largest volume size trainable using a device with 12GB GPU memory is increased from  $48 \times 168 \times 168$  to  $48 \times 288 \times 288$ , under aforementioned fixed parameter setting.

### 3.4.2 Cross-Methodology Comparison

Next, we conducted cross-methodology evaluation by comparing the proposed method to other third-party methodologies.

Firstly, 2D U-Net+3D U-Net proposed by [34] is another version of model cascading. Compared to 3D FCN+3D U-Net and 3D RU-Net, a 2D U-Net serving as an RoI locator produces significantly more false positive candidates (larger  $AvgN_{RoIs}$ ) with larger length along the Z-axis, which degrades the performance and is more time costly.

Next, a 2D U-Net+BDC-LSTM [16] is evaluated, whose kU-Net is employed for intra-slice feature extraction and a bidirectional convolutional LSTM is used to explore intra-slice features. Since patch size no longer limits the effective receptive field, we evaluated this method only using the HighRes dataset with a large designed receptive field as is proposed in [16]. It scored similarly compared to a 3D U-Net+HL, highlighting the effectiveness of intra-slice LTSMs. However, it only partially resolved the problem and got higher recall rate along with larger ASD since that full 3D context utilization is still limited and more false positives are produced along the Z axis, and its execution time is significantly longer compared to the proposed method.

Additionally, a 3D-FPN based Mask R-CNN is evaluated. As the scores and figures illustrates, the limitation of 3D Mask R-CNN is two-fold: bad-shaped bounding boxes’ cutting off some parts of the objects, and low-resolution masks generated by the coarse-resolution fea-

ture maps of the FPN backbone.

Finally, we also set a super-voxel clustering based [4] method as the baseline. Without the merit of discriminative 3D deep features, super-voxels are inevitably over-segmented or under-segmented. In our experiments, one of the 64 targets went completely missing and significantly lowered the Dice score, while some wrong super-voxels were chosen as the output mask.

## 4 Discussion

In this paper, we proposed a method to inherit easy-to-train and detail-preserving merits of volume-to-volume 3D FCNs while acquiring fast RoI localization, target completeness and whole volume global understanding of a joint detection-segmentation framework enabled by its large receptive field shared across different tasks. We combined a whole volume RoI localization model named as Global Image Encoder and in-region segmentation model named as Local Region Decoder as a joint model named 3D RoI-aware U-Net (3D RU-Net). As the result, we could segment colorectal tumors accurately and fast.

We notice a recent trend that researches seek to segment medical objects via detection. But most of these works employ independent modules for different tasks, leaving the benefit of fast feature reusing and wide range context utilization of large receptive field not fully enjoyed. As a refinement, the proposed method utilize the pre-extracted globally encoded features for in-region segmentation, providing better understanding of the whole image in the segmentation branch to discriminate background from false positives, and further saved over 50% computing resource for each in-region segmentation, which significantly accelerated the workflow in circumstances where multiple targets are detected.

Compared to successful and general Mask R-CNN for natural object instance segmentation, the advantage of the proposed framework over Mask R-CNN mainly lies on its full utilization of voxel-wise labels for target detection, and the lossless segmentation process similar to volume-to-volume 3D FCNs that fully restores the targets’ dimension. Hindered by additional aspect ratios and number-limited training samples, it’s sub-optimal to degrade voxel-wise labels to target-wise labels, hence bad-shaped bounding boxes are frequently predicted and cut off parts of the targets. Plus, the detail lossing in-region segmentation scheme of Mask R-CNN also produced inferior detail richness compared to end-to-end FCNs.

Finally, it’s significant to point out that the speed and performance gains are enabled by the memory efficiency of the proposed method that eliminates the need of conventional 3D U-Net for sliding-stitching workflow and enables one-step whole volume inference. Here we track the memory footprint to evaluate the memory efficiency of the proposed method in the environment where in-place computing is deactivated thus a ResBlock has nine tensor nodes. Given a typical T2 volume of 3D pelvic image of size  $40 \times 320 \times 320$ , by body cropping, the size typically drops to  $40 \times 180 \times 320$ . With this volume as input, the

GPU memory footprint details are listed in TABLE. 3. By constructing the Local Region Decoder, a GPU can assign 90% of its GPU memory to the encoder to process larger volumes and spend only 10% GPU memory on the segmentation stage, while conventional encoder-decoder networks spend 50% GPU memory on each path, as is hypothetically computed in the Standard Decoder section of TABLE. 3. Therefore, the applicable volume size is dramatically enlarged. In addition, while model ensemble strategy is often considered to be computationally expensive, based on the proposed method, we can have the performance gain at a promisingly acceptable cost.

Although our method achieved competitive results, there are some limitations. Firstly, as is illustrated in Fig. 6, the model is often confused about which slice to start or end, thus this significantly affects the score. As is illustrated in TABLE. 2, all competing methods including applying a bidirectional convolutional LSTM [16] did not thoroughly tackle this issue. As an explanation, this difficulty is data-related and decision about starting and ending slice index can be observer-dependent due to weak contrast in the border of cancerous tissues and low resolution along Z axis. Secondly, for this specific task without the need of discriminating different instances of tumors, we did not include instance separation capability in our design. However, it can be addressed since that the RoI Locator is a template that any encoder-only detection method is applicable, yet it can still be beneficial to fully utilize voxel-wise labels for bounding box refinement.

## 5 Conclusion

In this paper, we proposed a joint RoI localization-segmentation-based framework for fully automatic one-step whole volume colorectal cancer segmentation referred to as 3D RoI-aware U-Net (3D RU-Net). We emphasized the importance and effectiveness of integrating RoI localization and in-region segmentation fed with globally encoded features to perform fast and accurate whole volume segmentation. The proposed method enables the merit of enlarging receptive fields originally limited by GPU memory capacity and ensemble models with different receptive field settings. A Dice-formulated multi-task hybrid loss function is present to smoothen the training process. Experimental results demonstrated impressive superiority in terms of accuracy and speed over competing methods. In principle, the proposed framework is scalable enough to be adopted to other medical image segmentation tasks.

## References

- [1] Feb, “World cancer report 2014,” *World Health Organization*, 2015.
- [2] Rebecca L. Siegel, Kimberly D. Miller, and Ahmedin Jemal, “Cancer statistics, 2017,” *CA: A Cancer Journal for Clinicians*, vol. 67, no. 1, pp. 7–30, 2017.
- [3] Dwarikanath Mahapatra, Peter J Schuffler, Jeroen AW Tielbeek, Jesica C Makanyanga, Jaap Stoker, Stuart A Taylor, Franciscus M Vos, and Joachim M Buhmann, “Automatic detection and segmentation of crohn’s disease tissues from abdominal mri,” *IEEE Trans. on Med. Imaging*, vol. 32, no. 12, pp. 2332–2347, 2013.
- [4] Benjamin Irving, Amalia Cifor, Bartłomiej W Papięż, Jamie Franklin, Ewan M Anderson, Michael Brady, and Julia A Schnabel, “Automated colorectal tumour segmentation in dce-mri using supervoxel neighbourhood contrast characteristics,” in *MICCAI*. Springer, 2014, pp. 609–616.
- [5] V Badrinarayanan, A Kendall, and R Cipolla, “Segnet: A deep convolutional encoder-decoder architecture for image segmentation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [6] Hao Chen, Qi Dou, Xi Wang, Jing Qin, Jack C. Y. Cheng, and Pheng Ann Heng, “3d fully convolutional networks for intervertebral disc localization and segmentation,” in *International Conference on Medical Imaging and Virtual Reality*, 2016, pp. 375–382.
- [7] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*. Springer, 2015, pp. 234–241.
- [8] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger, “3d u-net: learning dense volumetric segmentation from sparse annotation,” in *MICCAI*. Springer, 2016, pp. 424–432.
- [9] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi, “V-net: Fully convolutional neural networks for volumetric medical image segmentation,” in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 565–571.
- [10] Lequan Yu, Xin Yang, Hao Chen, Jing Qin, and Pheng-Ann Heng, “Volumetric convnets with mixed residual connections for automated prostate segmentation from 3d mr images,” in *AAAI*, 2017, pp. 66–72.
- [11] H. Chen, Q. Dou, L. Yu, J. Qin, and P. A. Heng, “Voxresnet: Deep voxelwise residual networks for brain segmentation from 3d mr images,” *Neuroimage*, 2017.
- [12] Qi Dou, Lequan Yu, Hao Chen, Yueming Jin, Xin Yang, Jing Qin, and Pheng-Ann Heng, “3d deeply supervised network for automated segmentation of volumetric medical images,” *Medical image analysis*, vol. 41, pp. 40–54, 2017.
- [13] Fisher Yu and Vladlen Koltun, “Multi-scale context aggregation by dilated convolutions,” *arXiv preprint arXiv:1511.07122*, 2015.

- [14] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia, "Pyramid scene parsing network," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2881–2890.
- [15] Jiazhou Wang, Jiayu Lu, Gan Qin, Lijun Shen, Yiqun Sun, Hongmei Ying, Zhen Zhang, and Weigang Hu, "A deep learning based auto segmentation of rectal tumors in mr images," *Medical physics*, 2018.
- [16] Jianxu Chen, Lin Yang, Yizhe Zhang, Mark Alber, and Danny Z Chen, "Combining fully convolutional and recurrent neural networks for 3d biomedical image segmentation," in *Advances in Neural Information Processing Systems*, 2016, pp. 3036–3044.
- [17] Holger R. Roth, Le Lu, Amal Farag, Hoo Chang Shin, Jiamin Liu, Evrim B. Turkbey, and Ronald M. Summers, "Deeporgan: Multi-level deep convolutional networks for automated pancreas segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2015, pp. 556–564.
- [18] Holger R. Roth, Le Lu, Ari Seff, Kevin M. Cherry, Joanne Hoffman, Shijun Wang, Jiamin Liu, Evrim Turkbey, and Ronald M. Summers, "A new 2.5d representation for lymph node detection using random sets of deep convolutional neural network observations," *Med Image Comput Comput Assist Interv*, vol. 17, no. 1, pp. 520–527, 2014.
- [19] Dong Nie, Li Wang, Ehsan Adeli, Cuijin Lao, Weili Lin, and Dinggang Shen, "3-d fully convolutional networks for multimodal isointense infant brain image segmentation," *IEEE Transactions on Cybernetics*, vol. PP, no. 99, pp. 1–14, 2018.
- [20] Qian Yu, Yinhuan Shi, Jinqun Sun, Yang Gao, Yakang Dai, and Jianbing Zhu, "Crossbar-net: A novel convolutional network for kidney tumor segmentation in ct images," *arXiv preprint arXiv:1804.10484*, 2018.
- [21] Torsten Rohlfing, Daniel B Russakoff, and Jr Maurer, Calvin R., "Performance-based classifier combination in atlas-based image segmentation using expectation-maximization parameter estimation," *IEEE Transactions on Medical Imaging*, vol. 23, no. 8, pp. 983, 2004.
- [22] S Klein, Ua Van-Der-Heide, Im Lips, M Van-Vulpen, M Staring, and Jp Pluim, "Automatic segmentation of the prostate in 3d mr images by atlas matching using localized mutual information," *Medical Physics*, vol. 35, no. 4, pp. 1407–1417, 2008.
- [23] Sean Murphy, Brian Mohr, Yasutaka Fushimi, Hitoshi Yamagata, and Ian Poole, "Fast, simple, accurate multi-atlas segmentation of the brain," in *International Workshop on Biomedical Image Registration*. Springer, 2014, pp. 1–10.
- [24] Bharath Hariharan, Pablo Arbellez, Ross Girshick, and Jitendra Malik, "Hypercolumns for object segmentation and fine-grained localization," in *Computer Vision and Pattern Recognition*, 2015, pp. 447–456.
- [25] Jifeng Dai, Kaiming He, and Jian Sun, "Convolutional feature masking for joint object and stuff segmentation," in *Computer Vision and Pattern Recognition*, 2015, pp. 3992–4000.
- [26] Pedro H. O. Pinheiro, Ronan Collobert, and Piotr Dollar, "Learning to segment objects candidates," in *Advances in Neural Information Processing Systems*, 2015.
- [27] Xiaomeng Li, Hao Chen, Xiaojuan Qi, Qi Dou, Chi Wing Fu, and Pheng Ann Heng, "H-denseunet: Hybrid densely connected unet for liver and tumor segmentation from ct volumes," *IEEE Transactions on Medical Imaging*, vol. PP, no. 99, pp. 1–1, 2017.
- [28] Fangzhou Liao, Xi Chen, Xiaolin Hu, and Sen Song, "Estimation of the volume of the left ventricle from mri images using deep neural networks," *IEEE Transactions on Cybernetics*, vol. PP, no. 99, pp. 1–10, 2017.
- [29] Bharath Hariharan, Pablo Arbellez, Ross Girshick, and Jitendra Malik, "Simultaneous detection and segmentation," in *European Conference on Computer Vision*, 2014, pp. 297–312.
- [30] Uijlings, R. R J., Van De Sande, E. A K., Gevers, Smeulders, and W. M A., "Selective search for object recognition," *International Journal of Computer Vision*, vol. 104, no. 2, pp. 154–171, 2013.
- [31] Pablo Arbeláez, Jordi Pont-Tuset, Jonathan T Barron, Ferran Marques, and Jitendra Malik, "Multiscale combinatorial grouping," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 328–335.
- [32] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, "Faster r-cnn: towards real-time object detection with region proposal networks," in *International Conference on Neural Information Processing Systems*, 2015, pp. 91–99.
- [33] Min Tang, Ziehen Zhang, Dana Cobzas, Martin Jagersand, and Jacob L Jaremko, "Segmentation-by-detection: A cascade network for volumetric medical image segmentation," in *Biomedical Imaging (ISBI 2018), 2018 IEEE 15th International Symposium on*. IEEE, 2018, pp. 1356–1359.
- [34] Anjali Balagopal, Samaneh Kazemifar, Nguyen Dan, Mu Han Lin, Raquibul Hannan, Amir Owrangi, and Steve Jiang, "Fully automated organ segmentation in male pelvic ct images," 2018.
- [35] Ross Girshick, "Fast r-cnn," *Computer Science*, 2015.

- [36] Jifeng Dai, Kaiming He, and Jian Sun, "Instance-aware semantic segmentation via multi-task network cascades," in *Computer Vision and Pattern Recognition*, 2016, pp. 3150–3158.
- [37] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick, "Mask r-cnn," in *IEEE International Conference on Computer Vision*, 2017, pp. 2980–2988.
- [38] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017, vol. 1, p. 4.
- [39] Liu Shu, Qi Lu, Haifang Qin, Jianping Shi, and Jiaya Jia, "Path aggregation network for instance segmentation," 2018.
- [40] Lee R. Dice, "Measures of the amount of ecologic association between species," *Ecology*, vol. 26, no. 3, pp. 297–302, 1945.
- [41] Zhanpeng Zhang, Ping Luo, Change Loy Chen, and Xiaoou Tang, "Facial landmark detection by deep multi-task learning," in *European Conference on Computer Vision*, 2014, pp. 94–108.
- [42] Hao Chen, Lingyun Wu, Qi Dou, Jing Qin, Shengli Li, Jie-Zhi Cheng, Dong Ni, and Pheng-Ann Heng, "Ultrasound standard plane detection using a composite neural network framework," *IEEE transactions on cybernetics*, vol. 47, no. 6, pp. 1576–1586, 2017.
- [43] H. Chen, X. Qi, L. Yu, Q. Dou, J. Qin, and P. A. Heng, "Dcan: Deep contour-aware networks for object instance segmentation from histology images.," *Medical Image Analysis*, vol. 36, pp. 135–146, 2017.
- [44] Haocheng Shen, Ruixuan Wang, Jianguo Zhang, and Stephen J. McKenna, "Boundary-aware fully convolutional network for brain tumor segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2017, pp. 433–441.
- [45] Dou Qi, Chen Hao, Lequan Yu, Qin Jing, and Pheng Ann Heng, "Multilevel contextual 3-d cnns for false positive reduction in pulmonary nodule detection," *IEEE Transactions on Biomedical Engineering*, vol. 64, no. 7, pp. 1558–1567, 2017.
- [46] Konstantinos Kamnitsas, Christian Ledig, Virginia F. J. Newcombe, Joanna P. Simpson, Andrew D. Kane, David K. Menon, Daniel Rueckert, and Ben Glocker, "Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation," *Medical Image Analysis*, vol. 36, pp. 61, 2016.
- [47] Botian Xu, Yaqiong Chai, Cristina M Galarza, Chau Q Vu, Benita Tamrazi, Bilwaj Gaonkar, Luke Macyszyn, Thomas D Coates, Natasha Lepore, and John C Wood, "Orchestral fully convolutional networks for small lesion segmentation in brain mri," in *Biomedical Imaging (ISBI 2018), 2018 IEEE 15th International Symposium on*. IEEE, 2018, pp. 889–892.
- [48] Yi-Jie Huang, Qi Dou, Zi-Xian Wang, Li-Zhi Liu, Li-Sheng Wang, Hao Chen, Pheng-Ann Heng, and Rui-Hua Xu, "Hl-fcn: Hybrid loss guided fcn for colorectal cancer segmentation," in *Biomedical Imaging (ISBI 2018), 2018 IEEE 15th International Symposium on*. IEEE, 2018, pp. 195–198.
- [49] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [51] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, vol. 1, p. 3.
- [52] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky, "Instance normalization: The missing ingredient for fast stylization," 2016.
- [53] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick, "Training region-based object detectors with online hard example mining," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 761–769.
- [54] N Otsu, "A threshold selection method from gray-level histogram," *IEEE Trans Smc*, vol. 9, no. 1, pp. 62–66, 1979.
- [55] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [56] Diederik Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *Computer Science*, 2014.