

Automatic detection and segmentation of colorectal tumor based on multimodality magnetic resonance imaging fusion and co-predictive learning

Yunhao Ge ^a, Bin Li ^a, Weixin Yan ^{a*}

^aRobotics Institute of Shanghai Jiao Tong University, Shanghai, China

ABSTRACT

Colorectal tumor detection and segmentation by a deep learning algorithm in magnetic resonance imaging (MRI) is a difficult task due to many mutually-affected challenges, including the complicated anatomical environments, limited information extracted from radiology images and inadequate training samples. Regarding the diagnosis algorithm, some end-to-end deep convolutional neural networks (CNNs) have achieved remarkable success in some medical image detection and segmentation tasks. However, due to the limitation of the features extracted from single modality MRI with a single algorithm, which are basics of diagnosis, the precision of detection and segmentation of colorectal tumor diagnosis are hard to improve. Instead, experienced doctors make their final decision by combining radiological data, pathology data, morphological data, and clinical information. In this paper, a co-predictive learning method, which consists of two parallel prediction algorithms, is proposed to imitate the diagnostic process of doctors by fusing multimodality MRI ((T1 weighted (T1-w), T2 weighted (T2-w) and diffusion-weighted imaging (DWI)) information and combining different detection algorithms (CNNs and machine learning methods). The T2-w MRI, which performs best in single modality detection after the experiment, is used to perform main-stream detection and segmentation. While the other MRI modalities (T1w, DCE), which have different complementary information on the classification task, can extract both high level features and artificial features with different algorithms and predict the possibility whether the target image contains tumors by the Xgboost algorithm. The image-level predictions are used to make main-stream co-prediction to improve the detection precision. Meanwhile, with the help of specific diagnostic characteristics of colorectal tumor and specific treatment skills of the medical image, the hardness of segmentation can be reduced to improve the quality of segmentation. Experimental results demonstrate that our co-predictive method can achieve better performance both qualitatively and quantitatively for colorectal tumor automatic detection and segmentation, compared with the state-of-the-art methods.

INTRODUCTION

Colorectal cancer is one of the most common malignant tumors of the digestive system, ranked the third most common cancer worldwide making up about 10% of all cases. Due the inapparent symptoms, most patients are late stage when diagnosed, missing the best treatment time [1]. Therefore, it is critical to diagnose and assess colorectal cancer timely. Magnetic Resonance Imaging (MRI), which create detailed images of the organs and tissues within the body, provides a noninvasive and more accurate way for colorectal cancer detection [2-3]. MRI can also help to predict clinically relevant endpoints, by predicting whether the patient will show a positive response to treatment can achieve significant clinical meaning. Because organ-preserving treatment strategies can be considered for them as an alternative to standard surgical resection [1]. Encouraging results have been achieved for volumetric measurements by multimodality MRI (diffusion-weighted imaging (DWI), apparent diffusion coefficient (ADC)). However, most of results are calculated from regions of interest (ROI) of the tumor that are typically obtained after manual tumor segmentation by experienced readers, which are highly time-consuming thus greatly limits their usage [4]. In the field of biomedical imaging in particular, deep learning has been largely utilized for automatic detection and segmentation purpose [5][6]. To the best of our knowledge, few studies have been conducted on the automatic localization and segmentation of colorectal cancer. Because colorectal detection is a difficult task due to many mutually-affected challenges, including the complicated anatomical environments, limited information extracted from radiology images and inadequate training samples. Atlas-based Auto-segmentation of rectal tumors can help improve efficiencies in contouring in the clinical practice setting [7], while the method (shortage)an automatic segmentation procedure use paired dynamic contrast-enhanced MRI (DCE-MRI) achieves a successful result. However, the requirement of large scales of image pair data greatly limits their usage as in some situation DCE-MRI pair data is infeasible to obtain [8]. A deep learning method is proposed to automatic segment the rectal tumor from real patient [8]. However, in the actual application for medical image diagnosis, judging if the target patient contains a colorectal tumor or not is significantly more important and difficult than segment the tumor in positive patient. Moreover, most deep learning methods in medical image diagnose are end-to-end network, which complete the detection or segmentation using a single algorithm depending on only the information from a single modality image. Compared with the experienced radiologist, who prefers making diagnosis after combining the information from multimodality MRI and other pathology or

dynamic information to improve the precision of detection and segmentation. To solve the above problems, in this paper, a co-predictive neural network has been used to automatically localize and segment colorectal rectum tumors.

Every patient had nearly 30 slides MRI in each modality, but for the patients with colorectal cancer, not all the 30 slides contained the tumors, which increased the hardness of detection and segmentation. Since T1 weighted (T1-w) MRI and T2 weighted (T2-w) MRI can extract more features by CNNs [5] about the shape and texture, which can help in the location and segmentation, while DCE give a much clearer and less noisy signal to make image-level prediction about whether the target image contains a tumor. Different modalities of MRI may be more sensible for machine learning algorithm, especially in limited dataset circumstance. Besides, doctors can also make segmentation by specific characteristics of colorectal cancer, like by limiting the output region (only around rectal tissues) and limiting the output of the location, which performs better than a non-maximum suppression algorithm. These are much more complex for end-to-end algorithm to finish.

To combine the advantages of different features extracted by multimodality MRI, the proposed co-predictive neural network includes two parallel prediction algorithms, which is called main-stream and side-stream. The main-stream algorithm achieves detection and segmentation on the T2-w MRI, because T2 performs best among different single modality MRI detection in our CNNs experiments. The other MRI modalities (T1-w, DCE), have different complementary information on classification tasks, can be used to extract high level features and perform other statistics, texture or gray features by different methods, which are used to predict the possibility whether the image contains tumors by the Xgboost algorithm [9]. The image-level predictions are used to make co-prediction in main-stream CNNs to improve the detection precision. This adjustment is useful to improve the precision of image-level detection. To take advantage of the specific characteristics of colorectal tumor in diagnosis and improve the detection and segmentation precision, crop methods are applied before the detection and segmentation using a simple CNN. In addition, to reduce the computational cost, in the main-stream neural network, the detection and segmentation task share the weight in the CNN part, which is used to extract the crucial high-level features.

MATERIALS and METHODS

Dataset

This retrospective study was approved by the Institutional Review Board (IRB) of the Sir Run Run Shaw Hospital, College of Medicine, Zhejiang University (Zhejiang, China); and the informed consent requirement was waived from the IRB. The colorectal tumor database (CTD) comprised of 18,720 images from 208 patients, who were diagnosed between June 2009 and November 2014, in accordance with the inclusion and exclusion criteria. The inclusion criteria for patients in this study were: (i): rectal cancer confirmed by biopsy; (ii) locally advanced disease determined by pre-treatment MRI images (stage T3 or T4); (iii) pre-treatment MRI scan performed within 4 weeks before treatment. The exclusion criteria included: (i): received complete neoadjuvant chemoradiotherapy before MRI scans; (ii) suffering from synchronous distant metastasis; (iii) poor MRI image quality due to motion artifacts.

All patients underwent MRI scans including T1-weighted, T2-weighted and dynamic contrast-enhanced (DCE). All the MRI scans were performed at our institution by using a 3.0 Tesla MRI scanner (Signa HDxt, GE Medical Systems, Milwaukee, WI, USA) with a phased-array body coil. No special bowel preparation was performed during the scan. A quality assurance check was performed on the MRI machine to ensure the consistency of the image quality. The T2-weighted images were acquired by using T2-weighted fast spin echo sequence. The acquisition parameters were as follows: repetition time (TR): 2840ms; echo time(TE): 131ms; image resolution 0.49×0.49×4 mm; matrix: 512×512. The array spatial sensitivity encoding technique (ASSET) was used with an acceleration factor of 2. The T1-weighted sequences were obtained using a spoiled gradient echo sequence. The T1-weighted sequences parameters were: TR: 4.4 ms; TE: 1.9 ms; flip angle: 12°; bandwidth: 325.5 kHz; image resolution: 0.7×0.7×2 mm; matrix: 512×512. For the DCE scan, contrast injection and data acquisition were performed simultaneously. All patients were injected with 0.1 mmol/kg body-weight gadolinium-diethylenetriamine penta-acetic acid (Gd-DTPA) at a speed of 2.5 mL/s. Four phases data were acquired, before the injection of the contrast agent (Time-1), 15 seconds after the injection (Time-2), 60 seconds after the injection (Time-3), and 120 seconds after the injection (Time-4). All images were reviewed on MIM® (MIM software Inc., Cleveland, OH, USA) by an experienced rectal MRI radiologist. Then the three-dimensional (3D) tumor volumes, excluding the intestinal lumen, were manually segmented as ground truth. The volumes in the T1-weighted and T2-weighted images were segmented separately. The segmentation results were also reviewed by another anorectal clinician.

The CTD was divided in two independent datasets, called the training dataset and validation dataset. The training dataset involving 188 patients was used to train the co-predictive model. The validation dataset involving 20 patients was used to test the performance of the model developed.

The detailed information of the RCD was summarized in Table 1.

Table1. Information of the colorectal tumor database (CTD)

Dataset	Amount	T1-weighted image	T2-weighted image	DCE	Ground Truth	Total(images)
Train and validation	188 patients	5640	5640	5640	5640	16920
test	20 patients	600	600	600	600	1800

Proposed Method

Given the multiparametric MRI (mp-MRI) data of a patient, including T1-w, T2-w and DWI image sequences, our goal is to (1) automatically classify whether each slice of the mp-MRI contains tumor or not, and (2) for slices classified as positive for colorectal tumor, localize the position of the cancerous tissues and segment the area to obtain the pixel level classification information, the ROI. The framework of our automated co-predictive colorectal tumor detection and segmentation system, which is depicted in Fig. 1, consists of 2 main steps. First, after recording of the multimodality MRI, an automated method is applied to detect the rectal region using a squared bounding box (bbox) on every T2-w slice. Then, the bbox is used to aligned all T1-w and DWI slices and the content within each bbox is cropped and the intensities are normalized. The ground truth in training is extracted from the manually colorectal tumor constructed CTD. Second, the cropped T2-w image is input into the main-stream co-predictive algorithm and 2 items are output: (1) the bbox regression as detection output, which contains the rectal tumors (2) the output of the binary mask of the ROI, where each pixel is predicted as a colorectal tumor region. Meanwhile, the T1w and DWI modality images go through the shared structure CNNs in the main-stream algorithm to output the high-level features. Other artificial features are also obtained by the machine learning algorithm, which contain statistical features, gray-level co-occurrence matrices (GLCM) features [10] and local features. Then, the extracted features are input into a Xgboost classifier and output the possibility of the image-level prediction on whether this cropped rectal square contains a tumor or not, the prediction is used to make the co-prediction in the classification part of the main-stream algorithm. Our co-predictive CNNs has two main advantages: (1) high accuracy of image-level detection with the novel fusion solution by different modalities of MRI. By using different algorithms to extract high-level and artificial features from the aligned T1w and DWI, the Xgboost prediction depending on these combined features can improve the detection accuracy on the image-level (most of slice in each patient may contain no tumor regain), as well as reduce the difficulty of the segmentation of the ROI. (2) The bbox output and binary mask depend on the same CNN structure and the same heatmap feature, which significantly reduce the computational cost in main-stream. Additionally, the main-stream and side-stream also have shared CNNs structures to avoid the repetitive computation.

The framework of a fully automatic co-predictive method for colorectal tumor detection and segmentation includes five key components, as shown in Fig.1.

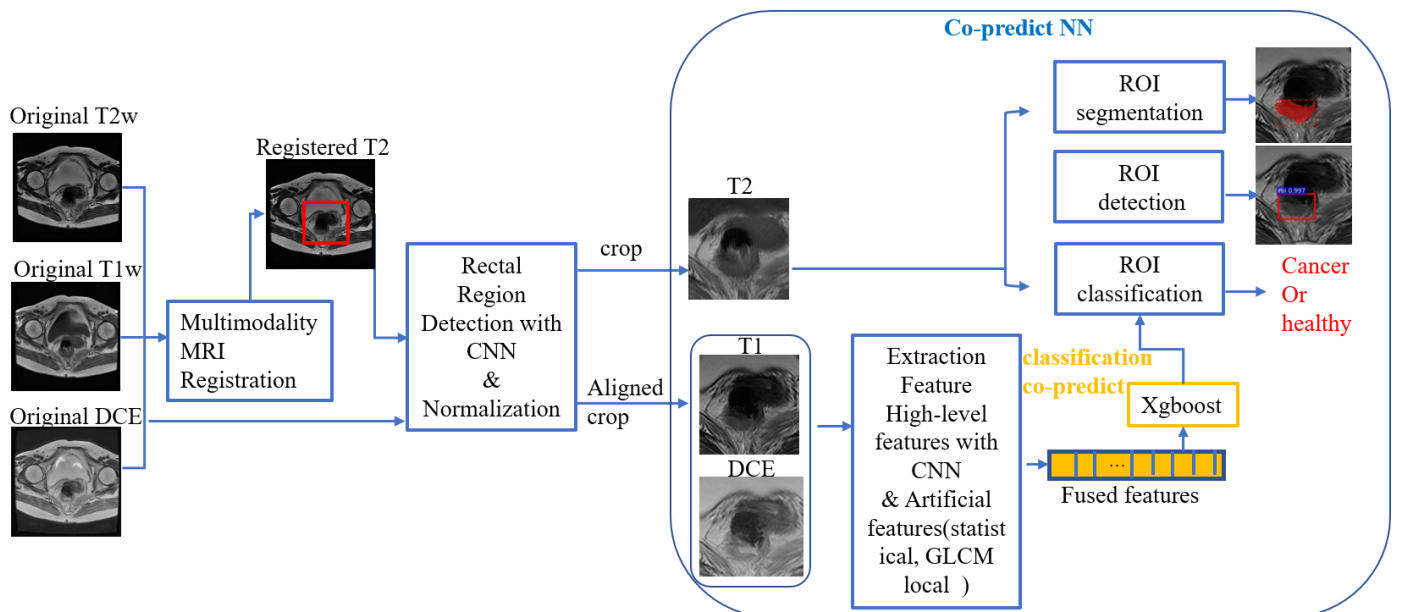


Fig. 1. The framework of the automated rectal localization system, which includes five key components: (1) multimodality (T1-w, T2-w and DWI) image registration; (2) rectal detection with CNN and normalization; (3) main-stream of co-predictive algorithm with ROI detection and segmentation; (4) tumor classification via Xgboost; (5) combination of the main-stream and side-stream to make the co-prediction.

Image preprocessing

The image preprocessing was performed including multi-modality image registration and intensity normalization. For the multi-modality MRI images, the DWI images were selected as the reference image, and the T1-weighted and T2-weighted images were deformed to match them by using the elastix toolbox [11]. The images were normalized, and the formula was as follows:

$$G_{norm} = \frac{G - G_{min}}{G_{max} - G_{min}} \quad (1)$$

where G was grey intensity for each pixel in the images, G_{norm} was normalized grey value, G_{min} was the minimum grey value of the image and G_{max} was the maximum grey value of the image.

Automated rectal detection and normalization

Once we achieve the multi-modality image registration and intensity normalization, we used a simple CNNs regression to automatically crop a square region including the whole rectal region for every T2-w slice along the transverse direction. The architecture of our CNNs model is depicted in Fig. 2. The CNN model is trained by a set of original T2-w slices using the manually labeled squared bbox to note the rectal regions. There are three output parameters from the CNN regression: the (x, y) coordinate position of the center of the square bbox, and the length l of the side of the square bbox. Note that these three parameters are normalized by the origin T2-w size, the range of the x, y coordinate position is $[-1, 1]$ and the range of length l is $[0, 1]$. We adopted the activation of the tanh function as the final output layer of the x, y position, and the sigmoid function as the final output layer of the length l . and the corresponding loss function is defined as:

$$loss = \frac{1}{3} (|\tanh(p_1 - x_t)| + |\tanh(p_2 - y_t)| + |sigmoid(p_3 - l_t)|) \quad (2)$$

where p_1, p_2 , and p_3 are the three final outputs of the CNN as shown in Fig. 2, and (x_t, y_t, l_t) are the normalized center coordinates and the length of the rectal region.

$$\begin{aligned} x_t &= \frac{2x - w}{w} \\ y_t &= \frac{2y - w}{w} \\ l_t &= \frac{l}{w} \end{aligned} \quad (3)$$

where (x, y, l) indicate the non-normalized coordinates and length ranging from 0 to w , and w is the width and height of the input square image with the size of 256x256. The CNN parameters are updated for the rectal region detection by minimizing the loss function.

Main-stream Network Architecture

The two main functions of the co-predictive neural network are implemented in main-stream, ROI detection and Segmentation. Regarding the object only task, the framework of Faster R-CNN, which consists of two stages. The first, the Region Proposal Network (RPN) [12], proposes candidate object bounding boxes. The second, which is similar to Fast R-CNN [12], extracts features using the ROI Pool operator from each candidate box and performs classification and bbox regression. This structure can also share the features used by two parts for faster inference. For adding the segmentation function as well, Mask R-CNN [14] has an excellent CNN framework and a mature pipeline for combining the object detection and segmentation tasks. Mask R-CNN improves the function of Faster R-CNN [12], which has two outputs for each candidate object, a class label and a b-box offset, by adding a third branch that outputs the object mask. Thus, Mask-R-CNN achieves a conceptually simple, flexible, and general framework for object instance segmentation. The design of the main-stream structure is similar to that of Mask-R-CNN, which consists of three main parts. First, as Table.2, we use ResNet101[15] as the backbone architecture, which can extract the high-level features from the input image. The architecture of Resnet 101 is as shown in Table 2. In the main-stream structure, the 91conv layers from conv1 to conv4-x are shared parts of the CNNs, the output of conv4_x is used as the input of the ROI align [14] layers. In addition, the output of conv4_x is also used as input of the second part—the RPN [12], which proposes the candidate object b-boxes that may contain the colorectal tumor. In other words, we share parts of the backbone architecture with the RPN. The proposed region is mapped into the feature map of the ROI align layers in main-stream, the features in the proposed region then go through the conv5_x to output the final high-level feature map of the input image. Ultimately, the third part—head Architectures is used to construct three parallel outputs: classification, detection and segmentation.

Table. 2 Architecture of the ResNet-101

Layer name	Output size	structure
Conv1	128×128	7×7, 64, stride=2
Conv2_x	64×64	3×3 max pool, stride=2
		$\begin{bmatrix} 1 \times 1, & 64 \\ 3 \times 3, & 64 \\ 1 \times 1, & 256 \end{bmatrix} \times 3$
Conv3_x	32×32	$\begin{bmatrix} 1 \times 1, & 128 \\ 3 \times 3, & 128 \\ 1 \times 1, & 512 \end{bmatrix} \times 4$
Conv4_x	16×16	$\begin{bmatrix} 1 \times 1, & 256 \\ 3 \times 3, & 256 \\ 1 \times 1, & 1024 \end{bmatrix} \times 23$
Conv5_x	8×8	$\begin{bmatrix} 1 \times 1, & 512 \\ 3 \times 3, & 512 \\ 1 \times 1, & 2048 \end{bmatrix} \times 3$

Side-stream Network Architecture with Multimodality MRI Fusion

The input of the side-stream includes two modalities of MRI, namely T1-w and DWI, these two modalities do not perform as well as the T2-w in the single modality ROI detection and Segmentation based on our previous experiment. Thus, we utilize the information in T1-w and DWI with the side-stream to improve the classification decision in the main-stream. Two methods are simultaneously used to extract the features in the ROI image. First, the shared part of the CNNs in main-stream (from conv1 to conv4-x in Table.2) was used to extract a set of complex high-level convolutional features. The global average pooling (GAP) operation generates the final output by averaging all scores within feature maps and the dimension of the final output is 512-d. Second, the artificial features composed of GLCM features, location features, and artificial features are obtained. Before the extraction of artificial features, the ROI image needs to be converted to the gray size. The artificial features consist of the statistical parameters, GLCM features and location features. Below is a list of twelve features extracted from the mask area of the original image; they were selected to be combined with the output of the final convolutional layer to train and test our method. Statistical parameters: four statistical parameters of the ROI region were extracted: mean, variance, skewness and kurtosis. For each unique pixel value f_k $p_f(f_k)$ is the probability

of this unique pixel value in the whole ROI. The four parameters are calculated as follows (Eq.4-7):

$$mean : \mu = \sum_{k=1}^N f_k p_f(f_k) \quad (4)$$

$$variance : \sigma^2 = \sum_{k=1}^N (f_k - \mu)^2 p_f(f_k) \quad (5)$$

$$skewness : ske = \sum_{k=1}^N [(f_k - \mu)^3 p_f(f_k)] / \sigma^3 \quad (6)$$

$$kurtosis : kur = \sum_{k=1}^N [(f_k - \mu)^4 p_f(f_k)] / \sigma^4 \quad (7)$$

The GLCM features: the GLCM features consists of the sum entropy (SE), sum average (SA), difference variance (DV) and difference entropy (DE). The *SE* is a logarithmic function of ROI in consideration; *SA* is calculated from the ROI and the size of the gray scale; *DV* is a variance measure between the ROI intensities calculated as a function of the *SE* previously calculated; *DE* is an entropy measure which provides a measure of no uniformity, while taking into consideration a measure of the difference obtained from the original image. These four parameters are calculated as follows:

$$SE = - \sum_{i=2}^{2N_g} p_{x+y}(i) \log\{p_{x+y}(i)\} \quad (8)$$

$$SA = \sum_{i=2}^{2N_g} i p_{x+y}(i) \quad (9)$$

$$DV = \sum_{i=2}^{2N_g} (i - SE)^2 p_{x-y}(i) \quad (10)$$

$$DE = - \sum_{i=2}^{2N_g} p_{x-y}(i) \log\{p_{x-y}(i)\} \quad (11)$$

Location features: four parameters about the ROI location and shape were extracted: convexity (C(S)), compactness (C), aspect ratio (AR), area ratio (R_Area). The four parameters are calculated as follows (Eq.12-15):

$$C(S) = \frac{A}{Area(CH(S))} \quad (12)$$

$$C = \frac{P^2}{4\pi A} \quad (13)$$

$$AR = \frac{D_y}{D_x} \quad (14)$$

$$R_Area = \frac{Area_ROI(in_pixels)}{Area_window(in_pixels)} \quad (15)$$

where S is a ROI, $CH(S)$ is its convex hull and A is the ROI's area, P is the ROI's perimeter, and $Area_window = D_x * D_y$, D_x is the width's ROI and D_y is the height's ROI.

Ultimately, The output of the fourth layer—including 16 channels of 1x1 array as well as the 12 features extracted from the original image—is

fully connected to Xgboost and outputs the category of the input clinical image. The image-level decision of the Xgboost can adjust the classification prediction in main-stream, the predicted result will combine both the main-stream and side-stream classification results. As shown in fig 4, the classification result q in the side-stream based on T1w and DWI contains the weight Wq in the final prediction, while the main-stream classification results contain the weight Wp . If we have a high probability, which means that the image may contain colorectal tumor, the final prediction will be improve by the weight Wp . In contrast, if we have low probability, which means that the image may be that of a healthy rectum, the final prediction will be reduced to reluctant the output of the detection box.

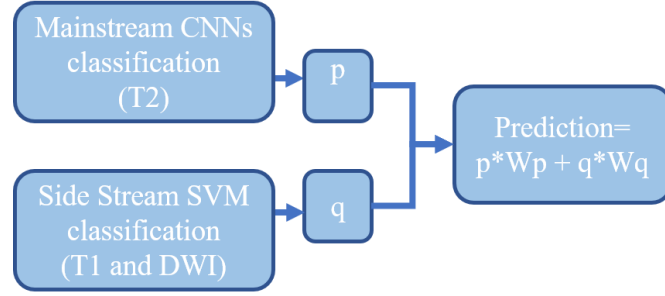


Fig. 4 Framework of the co-predictive network in the ROI classification part

Co-predictive Network Training

Loss function

Formally, during training, we define a multi-task loss function on each sampled ROI as $L = L_{cls} + L_{box} + L_{mask}$. Different from the Mask-RCNN, which contains a multi-class task, our co-predictive network only needs a two-task class. Thus, the classification loss L_{cls} and b-box loss L_{box} are as follows (Eq.16-18):

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*) \quad (16)$$

$$L_{cls}(p_i, p_i^*) = -\log[p_i^* p_i + (1 - p_i^*)(1 - p_i)] \quad (17)$$

$$L_{reg}(t_i, t_i^*) = \begin{cases} 0.5(t_i - t_i^*)^2 & \text{if } |t_i - t_i^*| < 1 \\ |t_i - t_i^*| - 0.5 & \text{otherwise} \end{cases} \quad (18)$$

where N_{cls} and N_{reg} are the mini-batch size and the number of anchor locations respectively, i is the index of an anchor in a mini-batch

and p_i is the predicted probability of the anchor i being colorectal tumor. The ground-truth label p_i^* is 1 if the anchor is positive for colorectal

tumor, and it is 0 if the anchor is negative; t_i is a vector representing the four parameterized coordinates of the predicted bounding box[12]; t_i^* is that of the ground-truth box associated with a positive anchor.

The classification loss L_{cls} is a log loss over two classes (object vs. not object) as in Eq.17[12]. For the regression loss, we use the Eq.18 robust

loss function (smooth L1) as previously defined [12]. The outputs of the cls and reg layers consist of $\{p_i\}$ and $\{t_i\}$, respectively. Different from Mask-RCNN, our definition of L_{mask} allows the network to generate masks for only 1 class, we typically used a per-pixel softmax and defined L_{mask} as the average binary cross-entropy loss.

Implement details

The ground-truth definition is like the Faster-R-CNN, the ROI is considered positive if it has IOU with a ground-truth box of at least 0.5 and negative otherwise. The mask target is the intersection between the ROI and its associated ground-truth mask. We trained the main-stream and side-stream separately. We followed the “image-centric” sampling strategy [12] to train this network. Each mini-batch arises from a single image that contains many positive and negative example anchors. The main-stream is trained first, and transfer learning is used due to the size of our training dataset: we randomly initialized all head layers by drawing weights from a zero-mean Gaussian distribution with standard deviation of 0.01. All other layers in the backbone are ResNet101 initialized by pretraining a model for ImageNet classification [12], as it is standard practice [12]. The training tricks are used as the training pipeline [12] in Faster R-CNN to share convolutional layers between the RPN and the backbone. We first trained the head part only, to avoid breaking the function for extracting high-level features, we froze all layers in the backbone: conv1_x to conv4_x. Then, we fine-tuned all the other layers. The side-stream is trained after the main-stream, as we used the shared structure to acquire the features of T1w and DWI, the image level ground truth and extracted features are used to train the Xgboost. While the shared structure in side-stream do not need training again. When predicting, the co-predictive network is implemented with the fusion of the multimodalities of MRI.

EXPERIMENTS and RESULTS

Each mini-batch has 4 images per graphics processing unit (GPU) and each image has 32 sampled ROIs, with a ratio of positive to negatives of 1:3 [12]. We trained on 2 GPUs (Titan X with 16G memory per GPU), thus the batch size is 8. There were 100 steps per training epoch, and we trained the ‘head’ part with 200 epochs and ‘all layers’ with another 200 epochs. For the training we used a learning rate of 0.001. We also used a weight decay of 0.0001 and a momentum of 0.9. In the testing dataset, there were 30 image slices for each patient (each slice has 3 modalities of MRI, T1-w, T2-w and DCE), and 30 ground-truth images, which contain the information healthy. if the image has an ROI, the ground-truth will be positive 1 and contain the contour line of the ROI, otherwise if the image is healthy, without the ROI of colorectal tumor, the ground-truth will be negative 0. There were 122 positive image slices and 478 negative image slices in the test dataset.

Colorectal tumor detection and crop

In the part of the automated rectal detection and crop, our proposed method can achieve a 100% precision with our dataset. Some of the results are shown in Fig. 5, we enlarged the region of view and fixed it to 256*256 to encompass some of the surrounding information which can improve the performance in later detection and segmentation.

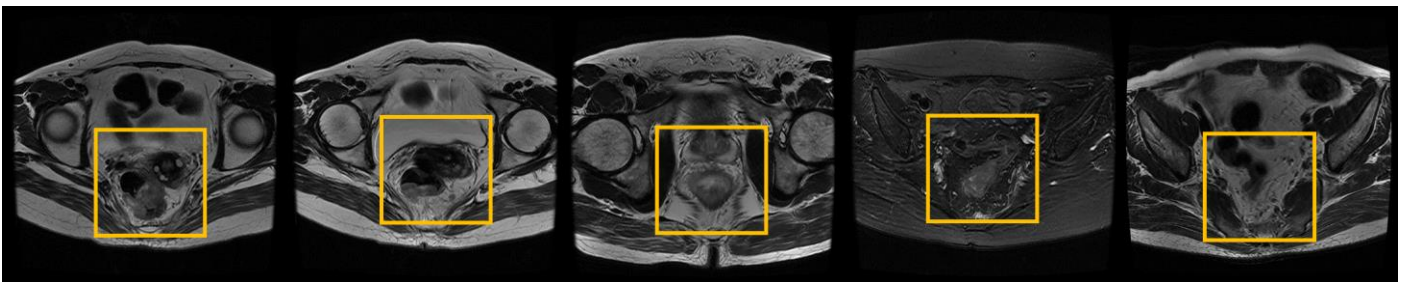


Fig.5 Colorectal tumor region detection and crop based on CNN

Image-level Tumor Detection

We calculated the information of the image-level sensibility (Recall), and used the average precision (AP) to evaluate the performance on the image-level colorectal tumor detection, In the actual application for medical image diagnosis, judging if the target MRI contains a rectal tumor or not is significantly more important and difficult than the other tasks. Before the decision of using the T2-w image as input in mainstream, we try different single-modalities as mainstream input to compare the performance on the image-level detection of colorectal tumor, the

performance on validation dataset is as follows:

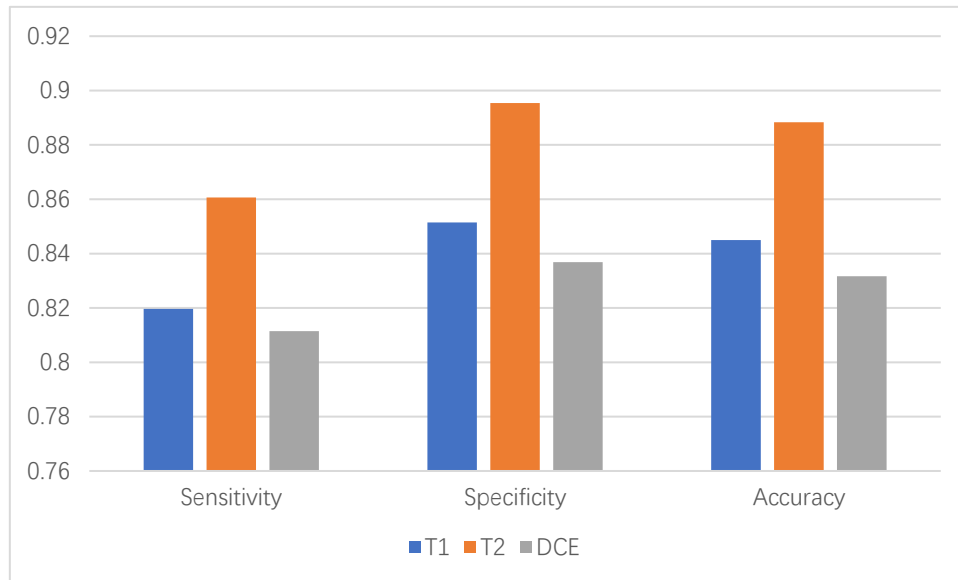


Fig.6 Image level detection performance using single modality in main-stream

Thus, we chose the T2-w image to be the input of the main-stream.

To choose the best machine learning algorithm in the final prediction step in side-stream, we used the validation dataset to evaluate the performance of different algorithms on the image level detection. We use three machine learning algorithms, Xgboost, SVM and MLP, for the classification model. For Xgboost, the depth of tree is 8 and the number of tree is 110. For SVM, we use the Gaussian RBF kernel and the sigma value of RBF kernel is 0.1. For MLP, we use three hidden layers MLP with 50, 25 and 18 hidden units in each hidden layer. The experiment results are as follows.

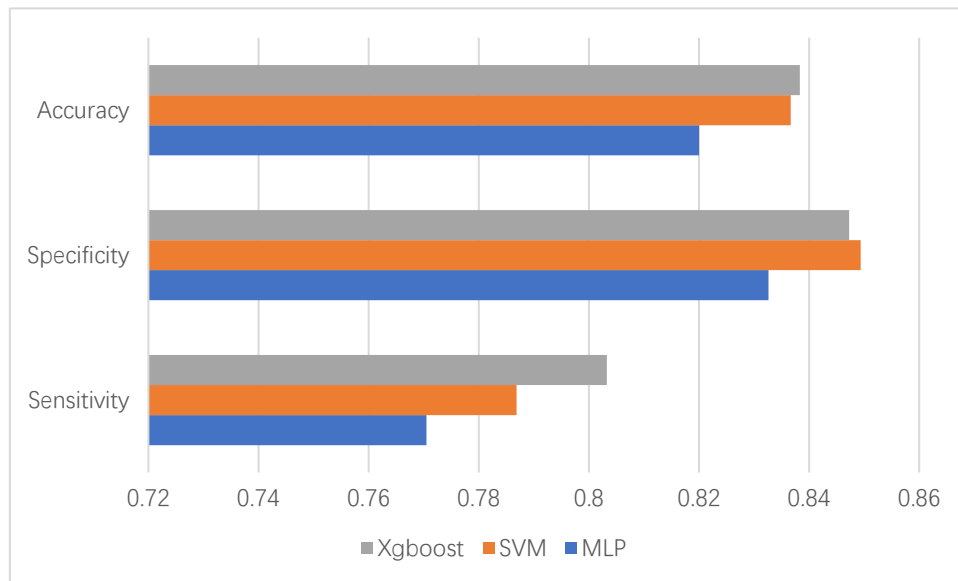


Fig.7 Image level detection performance using single modality in main-stream

Thus, we chose the Xgboost to make final prediction step in side-stream. To optimize the output threshold in main-stream and the combined weight in the co-prediction(ratio between main-stream prediction and side-stream prediction), we manually adjusted the parameters and evaluated the performance on image-level detection using validation dataset. Some of the details are in Table. 3.

Table. 3 Comparison of image level detection performance with different combined weight on CTD test-dev.

Weight ratio (main-stream : side-stream)	1:0.3	1:0.5	1:0.7	1:1	1:1.2
AP	90.7	91.9	90.1	89.3	88.5

We also performed an experiment to choose the best backbone in the main-stream neural network. Some of the details are in Table. 4:

Table. 4 Comparison of image level detection performance with different backbone in main-stream.

Base architecture	ZF	Vgg-16	Resnet-50	Resnet-101	Resnet-152
AP	88.4	89.0	89.6	91.8	91.5

We used the training dataset to train both the Mask-R-CNN network and our proposed co-predictive network, and used the test dataset to evaluate them. The evaluation metrics and results are in Table. 5.

Table. 5 Comparison of image level detection performance with different methods

Algorithm	AP(Test Dataset)	AP(Train Dataset)
Faster-rcnn(ZF)	85.66%	89.29%
Faster-rcnn(Vgg-16)	84.34%	89.21%
Mask-rcnn (ResNet50-FPN)	85.98%	92.60%
Mask-rcnn (ResNet101-FPN)	88.71%	93.46%
Mask-rcnn (ResNet152-FPN)	89.59%	93.38%
Co-predict	92.03%	96.40%

Shortage and future task: After the analysis in experiment, the result shows most of the false positive instance in test dataset are start or end part of colorectal tumor, where the tissue of tumor are fragmented and mixed with healthy texture. Thus the 3D convolution[18] may solve this problem by involving the information cross the adjacent tumor slice, while requiring larger computing space of device.

Rectal tumor Segmentation

Impact of automated colorectal detection and crop: Fig.8 shows the contribution of using the automated colorectal detection and crop to reduce false positives and improve the segmentation results. Fig.8(a) shows the original MR image with ground truth which highlighted in red, and Fig.8(b) shows the detection and segmentation results without colorectal detection and crop, which put the whole axial view original MR image as input. It can be observed that the cyan output is a false positive. Fig.8(c) shows the results using the colorectal detection and crop, and it can be observed that after colorectal detection and crop, the input image change from original whole image as yellow cropped image, which can make the detection and segmentation concentrate most on the correct field of view (FOV), and reduce the false positive output.

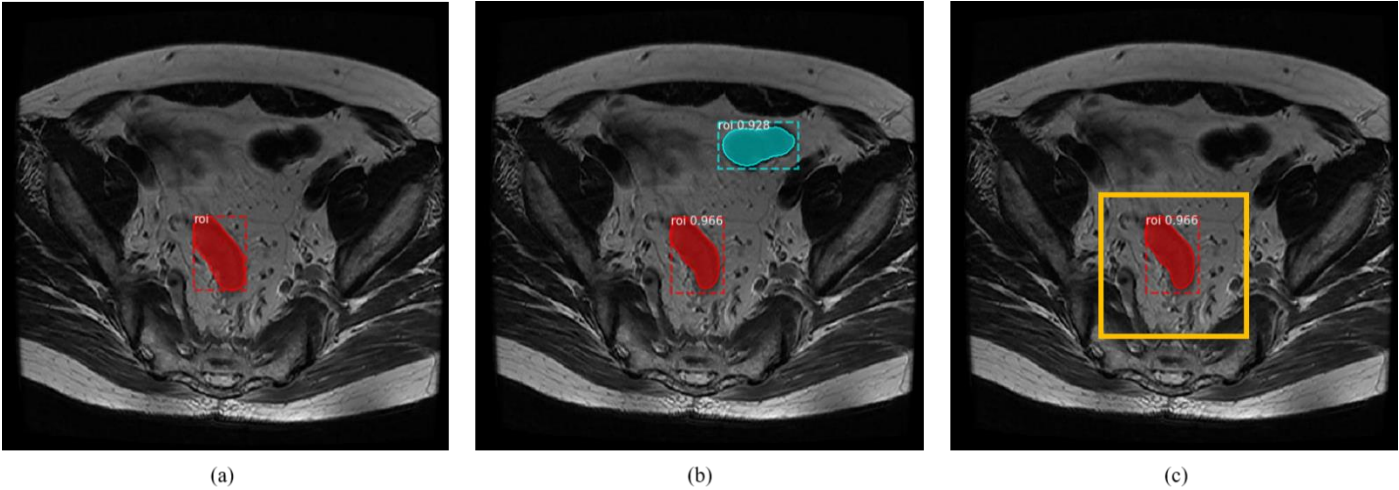


Fig.8 Advantages of using automated colorectal detection and crop to detect and segment the tumor ROI

Impact of applying diagnostic characteristics of colorectal tumor: Fig.9 shows the contribution of using diagnostic characteristics of colorectal tumor as priori knowledge into co-predictive method. Different with normal detection task, there is only one ROI at most in one axial colorectal slice. Thus, we only output one ROI which has the highest predicted value. Fig.9(a) shows the original MR image with ground truth which highlighted in red, and Fig.9(b) shows the detection and segmentation results applying colorectal detection and crop. It can be observed that the red and cyan ROI output are both in the field of cropped colorectal, which means only using the colorectal crop can not exclude the false positive. Fig.9(c) shows the results adding the diagnostic characteristics of colorectal tumor as priori knowledge into co-predictive method, and it can be observed that by restricting the number of output ROI with predicted value, the false positive in cropped colorectal FOV can be

reduced.

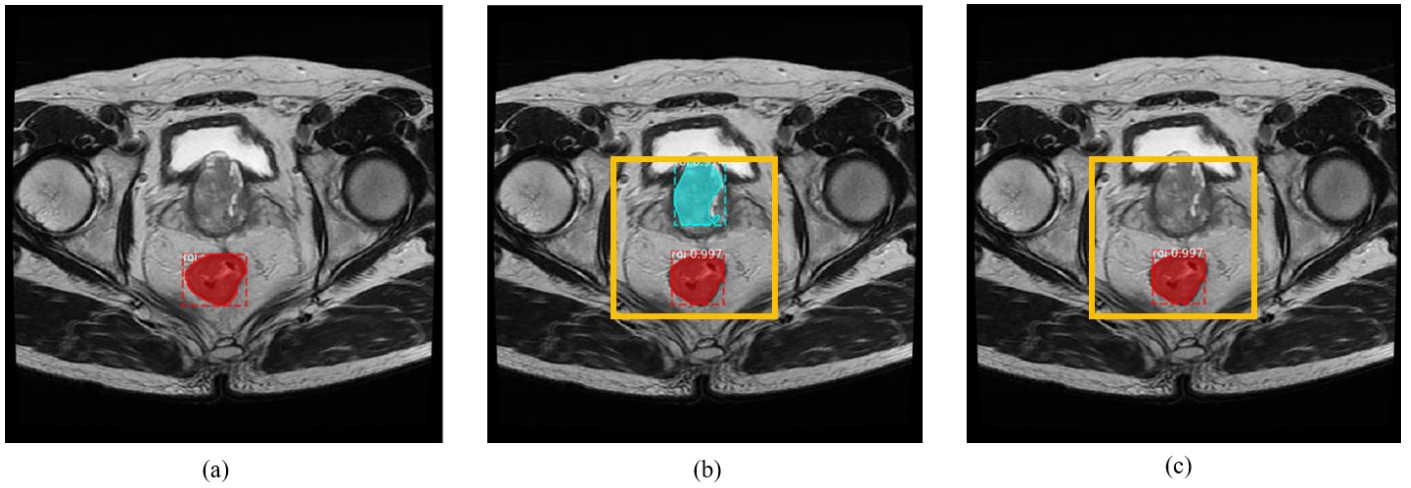


Fig.9 Advantages of using diagnostic characteristics of colorectal tumor to detect and segment the tumor ROI

The results are measured in terms of pixel intersection-over-union (IOU) to quantitatively evaluate the performance in pixel-level segmentation. Table.6 shows the quantitative evaluation results of comparison of segmentation pixel-level IOU on the CTD testing dataset with different methods including Mask-rnn, Unet [16] Vnet[17] and our co-predicted learning method. It can be observed that by automated colorectal crop, and co-predictive learning, the result can be further improved.

Table. 6 Comparison of segmentation pixel-level IOU on the CTD testing dataset with different methods.

Algorithm	IOU
Mask-rnn (ResNet101-FPN)	62.7%
U-net	53.5%
V-net	60.2%
Co-predict	70.1%

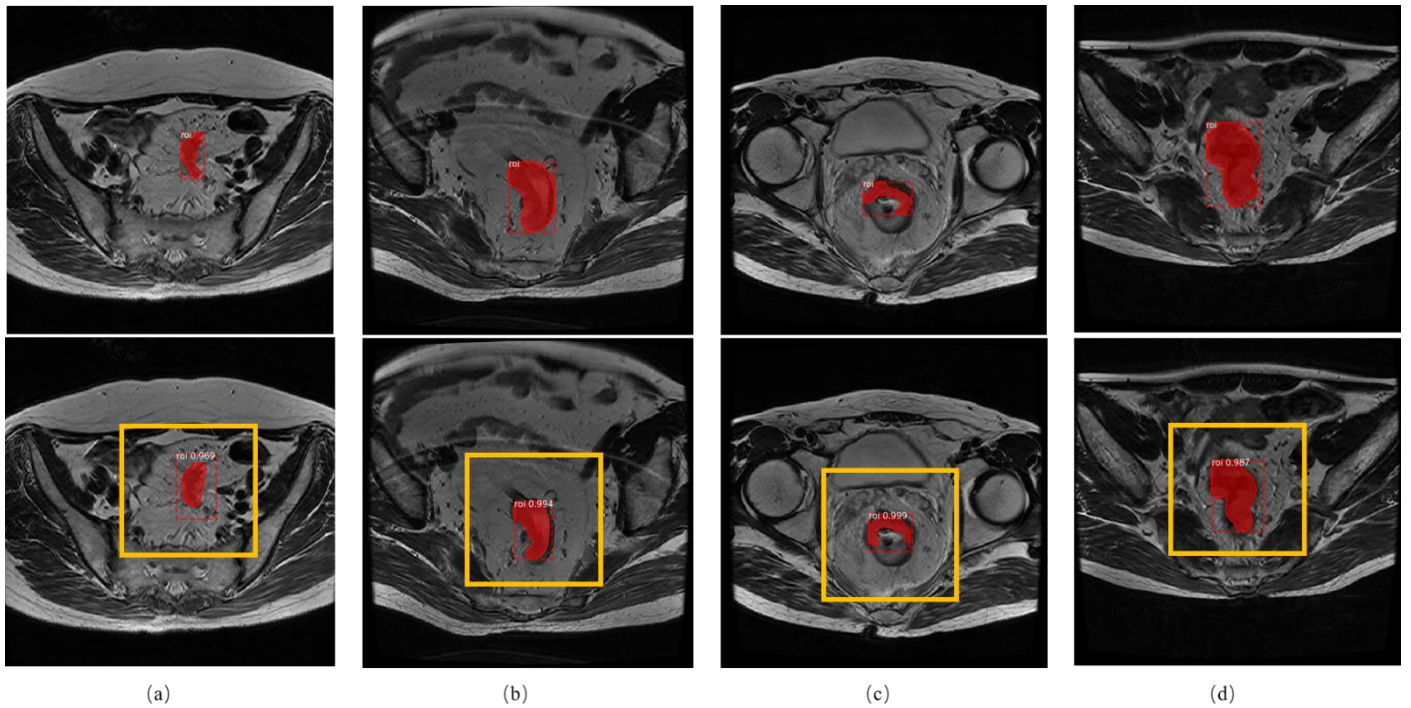


Fig. 10. Samples of segmentation results with co-predictive learning method in the CTD testing dataset. (a)-(d) are four pairs of segmentation results where the top figures are original MRI with ground truth which highlighted with red and the bottom figures are output ROI in red with predicted value. The yellow boxes are FOV after colorectal detection and crop.

Fig.10 shows the Samples of segmentation results with co-predictive learning method in the CTD testing dataset. It can be observed that after colorectal detection and crop, the FOV is reduced which can make the detection concentrate on the correct position and eliminate the influence of

irrelevant region. The performance of segmentation illustrated the advantages of using the co-predictive learning.

Acknowledgment

This work is supported by the National Natural Science Foundation of China under Grant No. 51475305.

CONCLUSIONS

To accurately detect and segment colorectal tumor, a fully automatic co-predictive learning method is proposed, which tries to imitate the diagnostic process of doctors by using multimodal fusion of multiple sequences of MRI and making the method more pertinent combining diagnostic characteristics of colorectal tumor. First, after the registration of multi-modality MRI and preprocessing, a specific CNN structure is designed to localize the position of the colorectal, (x, y) and l for the center and length of the cropped square respectively. Then a proposed co-predictive learning method is used to detect the ROI of the colorectal tumor and perform pixel-to-pixel segmentation. There are two streams in co-predictive learning method. The main-stream neural network detects and segments the ROI of tumor depending on the T2-w modality, while the side-stream, using both shared deep network and combined machine learning algorithm to extract features from T1-w and DWI, can improve the accuracy on the image-level detection in main-stream. In addition, the difficulty of segmentation is reduced by the fusion of multimodality MRI and applying diagnostic characteristics of colorectal tumor as priori knowledge into co-predictive method. Moreover, the computing cost is low due to the shared structure in mainstream, as well as between both the side-stream and main-stream. The colorectal tumor dataset (CTD) is used to quantitatively evaluate the efficiency and generalization capability of our co-predictive learning method. The training of main-stream and Xhboost in side-stream are separated using different loss functions and data structure. The detection results show that the colorectal detection neural network is able to accurately crop all the colorectal area. The tumor detection and segmentation performance with 92.03% image-level detection accuracy and 70.1% pixel-level IOU indicates that the proposed system is superior compared with the state-of-the-art methods.

References

- [1] Lambregts, D. M. J. et al. MRI and Diffusion-weighted MRI Volumetry for Identification of Complete Tumor Responders After Preoperative Chemoradiotherapy in Patients With Rectal Cancer: A Bi-institutional Validation Study. *Ann. Surg.* (2015) 262, 1034–9.
- [2] Martens, M. H. et al. Prospective, multicenter validation study of magnetic resonance volumetry for response assessment after preoperative chemoradiation in rectal cancer: Can the results in the literature be reproduced? *Int. J. Radiat. Oncol. Biol. Phys.* (2015) 93,1005–1014.
- [3] Seierstad, Therese, et al. "MRI volumetry for prediction of tumour response to neoadjuvant chemotherapy followed by chemoradiotherapy in locally advanced rectal cancer." *The British journal of radiology* 88.1051 (2015): 20150097.
- [4] Van Heeswijk, M. M. et al. Automated and semiautomated segmentation of rectal tumor volumes on diffusion-weighted MRI: Can it replace manual volumetry? *Int. J. Radiat. Oncol. Biol. Phys.* (2016) 94, 824–831.
- [5] LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* (2015)521, 436–444
- [6] Greenspan, H., Ginneken, B. van & Summers, R. M. Guest Editorial Deep Learning in Medical Imaging: Overview and Future Promise of an Exciting New Technique. *IEEE Trans. Med. Imaging* (2016)35, 1153–1159.
- [7] Gambacorta, Maria Antonietta, et al. "Clinical validation of atlas-based auto-segmentation of pelvic volumes and normal tissue in rectal tumors using auto-segmentation computed system." *Acta Oncologica* (2013) 52.8, 1676-1681.
- [8] Trebeschi, Stefano, et al. "Deep learning for fully-automated localization and segmentation of rectal cancer on multiparametric MR." *Scientific reports* (2017)7.1, 5301.
- [9] Chen, Tianqi, and Carlos Guestrin. "Xgboost: A scalable tree boosting system." *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining. ACM*, (2016).
- [10] Hall-Beyer, Mryka. "GLCM texture: a tutorial." *National Council on Geographic Information and Analysis Remote Sensing Core Curriculum* (2000).
- [11] Klein, S., Staring, M., Murphy, K., Viergever, M. A. & Pluim, J. P. W. Elastix: a toolbox for intensity-based medical image registration. *IEEE Trans. Med. Imaging* (2010)29, 196–205.

- [12] Ren, S., He, K., Girshick, R., & Sun, J.. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*. (2015) 91-99.
- [13] Nair, Vinod and Hinton, Geoffrey E. Rectified linear units improve restricted Boltzmann machines. In *ICML* (2010) 807–814.
- [14] He, K., Gkioxari, G., Dollár, P., & Girshick, R. Mask r-cnn. In *Computer Vision (ICCV), IEEE International Conference on* (2017) 2980-2988.
- [15] He, K., Zhang, X., Ren, S., & Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016) 770-778.
- [16] Olaf Ronneberger, Philipp Fischer, Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. (2015) arXiv preprint arXiv:1505.04597
- [17] Fausto Milletari, Nassir Navab, Seyed-Ahmad Ahmadi. V-Net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation. (2016) arXiv preprint arXiv:1606.04797
- [18] Chen, H., Dou, Q., Wang, X., Qin, J., Cheng, J. C., & Heng, P. A. 3D fully convolutional networks for intervertebral disc localization and segmentation. In *International Conference on Medical Imaging and Virtual Reality*. (2016) 375-382.