# Abdominal Multi-organ Segmentation with Organ-Attention Networks and Statistical Fusion

Yan Wang[a,1], Yuyin Zhou[a,1], Wei Shen[b,a], Seyoun Park[c,*], Elliot K. Fishman[c], Alan L. Yuille[d,a],

[a] Department of Computer Science, Johns Hopkins University, USA
[b] Key Laboratory of Specialty Fiber Optics and Optical Access Networks, Shanghai University, China
[c] Department of Radiology and Radiological Science, Johns Hopkins University, USA
[d] Department of Cognitive Science, Johns Hopkins University, USA

*spark139@jhmi.edu

## Abstract

Accurate and robust segmentation of abdominal organs on CT is essential for many clinical applications such as computer-aided diagnosis and computer-aided surgery. But this task is challenging due to the weak boundaries of organs, the complexity of the background, and the variable sizes of different organs. To address these challenges, we introduce a novel framework for multi-organ segmentation of abdominal regions by using organ-attention networks with reverse connections (OAN-RCs) which are applied to 2D views, of the 3D CT volume, and output estimates which are combined by statistical fusion exploiting structural similarity. More specifically, OAN is a two-stage deep convolutional network, where deep network features from the first stage are combined with the original image, in a second stage, to reduce the complex background and enhance the discriminative information for the target organs. Intuitively, OAN reduces the effect of the complex background by focusing attention so that each organ only needs to be discriminated from its local background. RCs are added to the first stage to give the lower layers more semantic information thereby enabling them to adapt to the sizes of different organs. Our networks are trained on 2D views (slices) enabling us to use holistic information and allowing efficient computation (compared to using 3D patches). To compensate for the limited cross-sectional information of the original 3D volumetric CT, e.g., the connectivity between neighbor slices, multi-sectional images are reconstructed from the three different 2D view directions. Then we combine the segmentation results from the different views using statistical fusion, with a novel term relating the structural similarity of the 2D views to the original 3D structure. To train the network and evaluate results, 13 structures were manually annotated by four human raters and confirmed by a senior expert on 236 normal cases. We tested our algorithm by 4-fold cross-validation and computed Dice-Sørensen similarity coefficients (DSC) and surface distances for evaluating our estimates of the 13 structures. Our experiments show that the proposed approach gives strong results and outperforms 2D- and 3D-patch based state-of-the-art methods in terms of DSC and mean surface distances.

## 1 Introduction

Segmentation of the internal structures, like body organs, in medical images is an essential task for many clinical applications such as computer-aided diagnosis (CAD), computer-aided surgery (CAS) and radiation therapy (RT). However, despite intensive studies of automatic or semi-automatic segmentation methods, there remain challenges which need to be overcome before these methods can be applied to clinical environments. In particular, detailed abdominal organ segmentation on CT is a challenging task both for

---

[1]The first two authors equally contributed to the work.

manual human annotation and for automatic segmentation algorithms for various reasons including the morphological complexity of the structures, the large variations between inter- and intra-subjects, and image characteristics such as low contrast of soft tissues.

Early studies of abdominal organ segmentation focused on specific single organs, for example relatively large isolated structures such as the liver [12, 20, 23] or critical structures such as blood vessels [17, 19]. However, most of the algorithms were based on specific features of the target organ, and so extensibility to the simultaneous segmentation of multiple organs was limited. For multi-organ segmentation, atlas-based approaches were adopted for many applications [2, 7, 13, 15, 16, 37, 40]. The general framework of atlas-based segmentations is to deformably register selected atlas images with segmented structures to the target image. Critical issues for this approach, which affect performance accuracy, include proper atlas selection, accurate deformable image registration, and label fusion. In particular, for the abdominal region, inter-subject variations are relatively large compared with other parts of the body (e.g., the brain) so the segmentation results are dependent on deformable registration between inter-subjects from the limited set of atlases, which is a challenging problem that critically affects the final accuracies. In addition, computational time is strongly dependent on the number of atlases. Therefore, selection of the proper number and types of atlases is a critical factor for both of the accuracy and efficiency.

Recently, learning-based approaches exploiting large datasets have been applied to the segmentation of medical images [4, 8, 9, 11, 14, 24–26, 33, 41]. In particular, deep convolutional neural networks (CNN) have been very successful [4, 8, 9, 11, 14, 24, 26, 28, 29, 33]. Targets include regions in the brain [4, 11, 14], chest [33], and abdomen [9, 28, 29]. The performance results of CNNs for organs (and even tumors) reach, or outperform, alternative state-of-the-art methods. Unlike multi-atlas-based approaches, deep networks do not require selecting a specific atlas or require deformable registration from training sets to a target image. In this study, we apply deep network approaches to abdominal organ segmentation.

Most studies based on deep networks, however, focused on a single structure segmentation, particularly for abdominal regions, and there are few studies of multi-organ segmentation partly due to technical challenges discussed later. We note that fully convolutional networks (FCNs) [21] have been generally accepted for organ segmentations on CT scans [8, 30, 39] partly because they give state-of-the-art performance for semantic segmentation of natural images [5, 21]. But there are three major characteristics of abdominal CT which we must address in order to obtain strong performance on multi-organ segmentation.

Firstly, many abdominal organs have weak boundaries between spatially adjacent structures on CT, e.g. between the head of the pancreas and the duodenum. In addition, the entire CT volume includes a large variety of different complex structures. Morphological and topological complexity includes anatomically connected structures such as the gastrointestinal (GI) track (stomach, duodenum, small bowel and colon) and vascular structures. The correct anatomical borders between connected structures may not be always visible in CT, especially in sectional images (i.e., 2D slices), and may be indicated only by subtle texture and shape change, which causes uncertainty even for human experts. This makes it hard for deep networks to distinguish the target organs from the complex background.

Secondly, there are large variations in the relative sizes of different target organs, e.g. the liver compared to the gallbladder. This causes problems when applying deep networks to multi-organ segmentation because lower layers typically lack semantic information when segmenting small structures. The same problem has been observed in semantic segmentation of natural images where the segmentation performance on small regions is typically much worse than on large regions, motivating the need to introduce mechanisms which attend to the scale [6].

Thirdly, although CT scans are high-resolution three-dimensional volumes, most current deep network methods were designed for 2D images. To overcome the limitations of using 2D CNNs for 3D images, Setio *et al.* [33] used multiple 2D patches reconstructed from 9 different directions around the target region for the task of pulmonary nodule detection. Zhuang *et al.* [40] used 2D axial, coronal, and sagittal slices for pancreas detection at the coarse level and also for segmentation at the finer level. More recently, there are studies which use 3D deep networks [8, 14, 24, 27, 30]. These, however, are not networks that act on the entire 3D CT volume but instead are local patch-based approaches (due to complex challenges of 3D deep networks discussed later in this paragraph). To address the problems caused by restricting to image patches, [14, 30] used a hierarchical approach with multi-resolutions, which reduces the dimension

**Figure 1.** The overall framework.

of the whole volume for initial detection and focuses on smaller regions at the finer resolution. But this strategy is best suited to a single target structure. Roth *et al.* [27] applied a bigger patch size to deal with the whole dense pancreatic volume, but this was also for single pancreas segmentation and hard to extend to the whole abdominal region. In general, 3D deep networks face far greater complex challenges than 2D deep networks. Both approaches rely heavily on graphics processing units (GPUs) but these GPUs have limited memory size which makes it difficult when dealing with full 3D CT volumes compared to 2D CT slices (which require much less memory). In addition, 3D deep networks typically require many more parameters than 2D deep networks and hence require much more training data, unless they are restricted to patches. But there is limited training data for abdominal CT images, because annotating them is challenging and requires expert human radiologists, which makes it particularly difficult to apply 3D deep networks to abdominal multi-organ segmentation. We have, however, implemented a 3D patch based approach for comparison.

To deal with the technical difficulties for abdominal multi-organ segmentation on CT, we introduce a novel framework of an organ-attention 2D deep networks with reverse connections (OAN-RC) followed by statistical fusion to combine the information from the three different views exploiting structural similarity using local isotropic 3D patches. OAN is a two-stage deep network, which computes an organ-attention map (OAM) from typical probability map of labels for input images in the first stage and combines OAM to the original input image for the second stage. This two-stage strategy effectively reduces the complexity of the background while enhancing the discriminative information of target structures (by concentrating attention close to the target structures). By training OAM with additional deep network, uncertainties and errors from the first stage are adjusted and the fidelity of the final probability map is improved. In this procedure, we apply reverse connections [18] to the first stage so that we can localize organ information at different scales by assisting the lower layers with semantic information.

More specifically, we apply OAN-RC to each sectional slice, which is an extreme form of anisotropic local patches but include the whole semantic (i.e. volume) information from one viewing direction. This yields segmentation information from separate sets of multi-sectional images (axial, coronal, and sagittal planes in this study similarly to most of medical image platforms for 2D visualization). We statistically fuse the three sources of information using local isotropic 3D patches based on direction-dependent local structural similarity. The basic fusion framework uses expectation-maximization (EM) similar to [2, 36]. But, unlike typical statistical fusion methods used for atlas-based segmentation, the input volumes and the target volumes for segmentation in our problem are the same. But different structures and texture patterns, from different viewing directions, will often generate nonidentical segmentations in 3D. Our

strategy is to exploit structural similarity by computing a direction-dependent local property at each voxel. This models the structural similarity from the 2D images to the original 3D structure (in the 3D volume) by local weights. This structural statistical fusion improves our overall performance by combining the information from the three different views in a principled manner and also imposing local structure.

Figure 1 describes the graphical concept of our framework. Our proposed algorithm was tested on 236 abdominal CT scans of normal cases collected as a part of FELIX project for pancreatic cancer research [22]. By experiments, our method showed robust and high fidelities to the ground-truth for all target structures with smooth boundaries. It outperformed 3D patch-based algorithms as well as 2D-based in terms of DICE-similarity coefficient and average surface distance with memory and computational efficiency.

# 2 Organ-Attention Networks with Reverse Connections

Given a 3D volume of interest (VOI) of a scanned CT image $V \subset \mathbb{R}^3$, our goal is to find the label of each voxel $v \in V$. The target structures (i.e., the labeled structures) are restricted to be organs which do not overlap with each other, so every voxel $v$ should be assigned to a label in a finite set $\mathcal{L}$. In this section we introduce our proposed organ-attention networks with Reverse connections (annotated as OAN-RC) which is run separately on three different views, and then in the next section we describe our novel structural similarity statistical fusion method which combines the segmentation results obtained from the OAN-RCs on the three different views.

## 2.1 Two-stage Organ Attention Network

We first introduce the OAN, which is composed of two jointly optimized stages. The first stage (stage-I) transforms the organ segmentation probability map to provide spatial attention to the second stage (stage-II), so that the segmentation network trained in stage-II is more discriminative for segmenting organs (because it only has to deal with local context). To assist the lower layers in stage-I with more semantic information, we employ reverse connections (Sec. 2.2), which pass semantic information down from high layers to low layers. The OAN is trained in an end-to-end fashion to enhance the learning ability of all stages.

The input images to our OAN are reconstructed 2D slices from axial, sagittal and coronal directions. Based on the normal vector directions of the sagittal ($X$), coronal ($Y$) and axial ($Z$) planes, we denote the 2D images by $\mathbf{I}_i^X$, $\mathbf{I}_j^Y$ and $\mathbf{I}_k^Z$ respectively, where $i = 1, \ldots, n_x$, $j = 1, \ldots, n_y$, $k = 1, \ldots, n_z$ and $n_x$, $n_y$, $n_z$ are the numbers of slices for the three directions, respectively, and $\bigcup_i \mathbf{I}_i^X = \bigcup_j \mathbf{I}_j^Y = \bigcup_k \mathbf{I}_k^Z = V$. Following the work of [39], we train an individual OAN for each direction.

Fig. 2 illustrated our organ-attention-network architecture. The network consists of two stages, where each stage is a segmentation network. For notational simplicity, we denote an input 2D slice by $\mathbf{I} \subset \mathbb{R}^{H \times W}$ and its corresponding label map by $\mathbf{T} = \{t_i\}_{i=1,\ldots,H \times W}$. Stage-I outputs a probability map $\mathbf{P}^{(1)} = f(\mathbf{I}; \mathbf{\Theta}^{(1)}) \subset \mathbb{R}^{H \times W \times |\mathcal{L}|}$ for each label at every pixel, where the probability density function $f(\cdot; \mathbf{\Theta}^{(1)})$ is a segmentation network parameterized by $\mathbf{\Theta}^{(1)}$. We use FCN [21] with reverse connections, which is explained in Sec. 2.2, as $\mathbf{\Theta}^{(1)}$. FCN is the backbone network throughout the paper. Each element $p_{i,l}^{(1)} \in \mathbf{P}^{(1)}$ is the probability that the $i$-th pixel in the input slice belongs to label $l$, where $l = 0$ is the background, and $l = 1, ..., |\mathcal{L}|$ are target organs. We define $p_{i,l}^{(1)} = \sigma(a_{i,l}^{(1)}) = \frac{\exp(a_{i,l}^{(1)})}{\sum_{t=0}^{|\mathcal{L}|} \exp(a_{i,t}^{(1)})}$, where $a_{i,l}^{(1)}$ is the activation value of the $i$-th pixel on the $l$-th channel dimension. Let $\mathbf{A}^{(1)} = \{a_{i,l}^{(1)}\}_{i=1,\ldots,H \times W, l=0,\ldots,|\mathcal{L}|}$ be the activation map. The objective function to minimize for $\mathbf{\Theta}^{(1)}$ is given by

$$\mathcal{J}^{(1)}(\mathbf{\Theta}^{(1)}) = -\frac{1}{H \times W} \left[ \sum_{i=1}^{H \times W} \sum_{j=0}^{|\mathcal{L}|} \mathbf{1}\left(t_i = l\right) \log p_{i,l}^{(1)} \right], \tag{1}$$

**Figure 2.** The architecture of our two-stage organ-attention network with reverse connections. The organ-attention network (OAN) is composed of two jointly optimized stages, where the first stage (stage-I) transforms the organ segmentation probability map by spatial attention to the second stage (stage-II). Hence the organ segmentation map generated in the organ-attention module guides the latter computation. The reverse connections, described in Sec. 2.2, modify the first stage of OAN as shown by dashed lines.

where $\mathbf{1}(\cdot)$ is an indicator function.

Using a preliminary organ segmentation map to guide the computation of a better organ segmentation can be thought as employing an attentional mechanism. Towards this end, we propose an organ-attention module by

$$\mathbf{Q} = \mathbf{W} * \mathbf{P}^{(1)} + \mathbf{b}, \tag{2}$$

where $*$ denotes the convolution operator, $\mathbf{W}$ indicates the convolutional filters, and $\mathbf{b}$ is the bias. (2) embeds cross-organ information into a single organ-attention map, $\mathbf{Q}$, which learns discriminative spatial attention for different organs automatically. By combining $\mathbf{Q}$ with the original input $\mathbf{I}$, we get an image which emphasizes each organ by

$$\mathbf{I}^{(2)} = \mathbf{I} \star \mathbf{Q}, \tag{3}$$

where $\star$ is the element-wise product operator. We apply $\mathbf{I}^{(2)}$ to the input of stage-II, and the probability of stage-II then becomes $\mathbf{P}^{(2)} = f(\mathbf{I}^{(2)}; \mathbf{\Theta}^{(2)})$.

In order to drive stage-II to focus on organ regions without needing to deal with complicated non-local background, we define a selection function, $\mathbf{1}(\mathbf{P}_0^{(1)} \leqslant \rho)$ where $\mathbf{P}_0^{(1)} = \{p_{i,0}^{(1)}\}_{i=1,\ldots,H \times W}$ is the probability map provided by stage-I. In stage-II, we only accept the region if $p_{i,0}^{(1)} > \rho$ and do not back-propagate it to stage-I. The loss function for stage-II is formulated as

$$\mathcal{J}^{(2)}(\mathbf{\Theta}^{(2)}, \mathbf{W}, \mathbf{b}) = -\frac{1}{H \times W} \left[ \sum_{i=1}^{H \times W} \sum_{j=0}^{|\mathcal{L}|} \mathbf{1}\left(p_{i,0}^{(1)} \leqslant \rho\right) \cdot \mathbf{1}\left(t_i = l\right) \log p_{i,l}^{(2)} \right]. \tag{4}$$

To jointly optimize stage-I and stage-II, we define a loss function aiming at estimating parameters $\mathbf{\Theta}^{(1)}$, $\mathbf{\Theta}^{(2)}$, $\mathbf{W}$, and $\mathbf{b}$ by optimizing

$$\mathcal{J} = h^{(1)} \mathcal{J}^{(1)}(\mathbf{\Theta}^{(1)}) + h^{(2)} \mathcal{J}^{(2)}(\mathbf{\Theta}^{(2)}, \mathbf{W}, \mathbf{b}), \tag{5}$$

where $h^{(1)}$ and $h^{(2)}$ are the fusion weights.

## 2.2 Reverse Connections

FCNs [21] have shown good segmentation results in recent studies, especially for single organ segmentation. However, for multi-organ segmentation, lower layers typically lack semantic information, which may

**Figure 3.** The reverse connections architecture of OAN stage-I. The network has reverse connections to the output of convolutional layers. In the training step, both backbone network and reverse connection side-outputs are supervised by the ground-truth. Finally, all reverse connection side-outputs and the output of backbone network are fused and made to approach ground-truth.

lead to inaccurate segmentation particularly for smaller structures. Therefore, we propose reverse connections which feed coarse-scale (high) layer information backward to fine-scale (low) layer for semantic segmentation of multi-scale structures, inspired by [18]. This enables us to connect abstract high-level semantic information to the more detailed lower layers so that all the target organs have similar levels of details and abstract information at the same layer. The reverse connections framework for stage-I is shown in Fig. 3. Fig. 4 illustrates a reverse connection block. Let $\mathbf{R}_n$ denote the reverse connection map of the $n$-th convolutional layer in the backbone network, i.e. FCN in this study, where $\mathbf{C}_n$ is the output of the $n$-th convolutional layer. A convolutional layer (with 512 channels by $3 \times 3$ kernels) is added after $\mathbf{C}_n$, and a deconvolutional layer (with 512 channels by $4 \times 4$ kernels) is applied after $\mathbf{R}_{n+1}$. $\mathbf{R}_n$ is then obtained via an element-wise summation of these two maps. $\mathbf{R}_7$ is the output of a convolutional layer (with 512 channels by $2 \times 2$ kernels) grafted onto $\mathbf{C}_7$. Let $\mathbf{w}^n$ denote the corresponding weights for obtaining $\mathbf{R}_n$. Following [18], we add reverse connections from $\mathbf{C}_4$ to $\mathbf{C}_7$.

With these learnable reverse connections, the semantic information of the lower layers can be enriched. In order to drive learned reverse connection maps to produce segmentation results approaching the ground-truth, we make each reverse connection map associate with a classifier. As the side-output layers proposed in [18] are designed for detection purposes, they are not suitable for our task. Instead we follow the side-outputs used in [38]. More specifically, a convolutional layer (with $|\mathcal{L}|$ channels by $1 \times 1$ kernels) is added on top of $\mathbf{R}_n$, whose output is denoted as $\mathbf{V}_n$, and followed by a deconvolutional layer (with $|\mathcal{L}|$ channels). We denote the weights of the $n$-th side-output layer by $\boldsymbol{\theta}^n$. The loss function for side-output

**Figure 4.** A reverse connection block.

layers $\mathcal{J}^{(s,1)}$ is defined as

$$\mathcal{J}^{(s,1)}(\mathbf{\Theta}^{(1)}, \mathbf{w}, \boldsymbol{\theta}) = \sum_{n=4}^{7} h_n^{(s,1)} \ell_n^{(s,1)}\left(\mathbf{\Theta}^{(1)}, \mathbf{w}^n, \boldsymbol{\theta}^n\right), \tag{6}$$

where $\ell_n^{(s,1)} = -\frac{1}{H \times W}\left[\sum_{i=1}^{H \times W}\sum_{j=0}^{|\mathcal{L}|} \mathbf{1}\left(t_i = l\right) \log p_{i,l}^{(s,1)}\right]$ and $p_{i,l}^{(s,1)}$ is the probability output of the $n$-th side-output layer.



**Figure 5.** Feature fusion strategy. A deep-to-shallow refinement is adopted for multi-scale side-output features. The final activation map $(\mathbf{A}^{(f,1)})$ for stage-I is an element-wise addition of the side-output activation map $(\mathbf{A}^{(r,1)})$ and the backbone network activation map $(\mathbf{A}^{(b,1)})$.

In order to combine the learned reverse connection maps of fine layers and coarse layers, we add up the predictions (i.e., $\mathbf{V}_n$) of the reverse connection maps from high layer to low layer gradually. First, $\mathbf{V}_6$ is fused with a $2\times$ upsampling of $\mathbf{V}_7$ by an element-wisely addition. Then we follow the same strategy and gradually merge $\mathbf{V}_5$ and $\mathbf{V}_4$, as shown in Fig. 5. To obtain a fused activation map $\mathbf{A}^{(f,1)} = \{a_{i,l}^{(f,1)}\}_{i=1,...,H \times W, l=0,...,|\mathcal{L}|}$ from the activation map of both side-outputs (i.e., $\mathbf{A}^{(r,1)}$) and convolutional layers in the backbone network (i.e., $\mathbf{A}^{(b,1)}$), a scale function is adopted followed by an element-wise addition by

$$\mathbf{A}_l^{(f,1)} = h_l^{(r,1)}\mathbf{A}_l^{(r,1)} + h_l^{(b,1)}\mathbf{A}_l^{(b,1)}, \qquad l = 0,...,|\mathcal{L}| \tag{7}$$

where $\mathbf{A}_l$ indicates the $l$-th channel of the activation map. $h_l^{(r,1)}$ and $h_l^{(b,1)}$ are fusion weights. Then the fused probability map, $\mathbf{P}^{(f,1)} = \{p_{i,l}^{(f,1)}\}_{i=1,\ldots,H\times W, l=0,\ldots,|\mathcal{L}|}$, can be obtained by $p_{i,l}^{(f,1)} = \sigma(a_{i,l}^{(f,1)})$. The final objective function for stage-I is defined by

$$
\begin{aligned}
\mathcal{J}^{(1)}(\boldsymbol{\Theta}^{(1)}, \mathbf{w}, \boldsymbol{\theta}) &= h^{(b,1)} \mathcal{J}^{(b,1)}(\boldsymbol{\Theta}^{(1)}) \\
&+ h^{(s,1)} \mathcal{J}^{(s,1)}(\boldsymbol{\Theta}^{(1)}, \mathbf{w}, \boldsymbol{\theta}) + h^{(f,1)} \mathcal{J}^{(f,1)}(\boldsymbol{\Theta}^{(1)}, \mathbf{w}, \boldsymbol{\theta}),
\end{aligned}
\tag{8}
$$

where $h^{(b,1)}$, $h^{(s,1)}$ and $h^{(f,1)}$ are fusion weights, and

$$
\mathcal{J}^{(f,1)}(\boldsymbol{\Theta}^{(1)}, \mathbf{w}, \boldsymbol{\theta}) = -\frac{1}{H \times W} \left[ \sum_{i=1}^{H \times W} \sum_{j=0}^{|\mathcal{L}|} \mathbf{1}\,(t_i = l) \log p_{i,l}^{(f,1)} \right].
\tag{9}
$$

Note that in our full system with the two-stage organ-attention network and reverse connections, all the parameters are optimized simultaneously by standard back-propagation

$$
\begin{aligned}
(\hat{\boldsymbol{\Theta}}^{(1)}, \hat{\mathbf{w}}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Theta}}^{(2)}, \hat{\mathbf{W}}, \hat{\mathbf{b}}) \\
= \arg\min \{ \mathcal{J}^{(1)}(\boldsymbol{\Theta}^{(1)}, \mathbf{w}, \boldsymbol{\theta}) + h^{(2)} \mathcal{J}^{(2)}(\boldsymbol{\Theta}^{(2)}, \mathbf{W}, \mathbf{b}) \}.
\end{aligned}
\tag{10}
$$

## 2.3 Testing Phase

In the testing stage, given a slice $\mathbf{I}$, we obtain the stage-I and stage-II probability map by

$$
\begin{aligned}
\mathbf{P}^{(1)} &= f(\mathbf{I}; \hat{\boldsymbol{\Theta}}^{(1)}, \hat{\mathbf{w}}, \hat{\boldsymbol{\theta}}) \\
\mathbf{P}^{(2)} &= f(\mathbf{I}; \hat{\boldsymbol{\Theta}}^{(2)}, \hat{\mathbf{W}}, \hat{\mathbf{b}}),
\end{aligned}
\tag{11}
$$

where $f(\cdot, \cdot)$ is the network functions defined in Sec. 2.1. A fused probability map of $\mathbf{P}^{(1)}$ and $\mathbf{P}^{(2)}$ is then given by

$$
\mathbf{P} = \mathbf{P}^{(1)} \circ \mathbf{1}(\mathbf{P}_0^{(1)} > \rho) + \mathbf{P}^{(2)} \circ \mathbf{1}(\mathbf{P}_0^{(1)} \leqslant \rho).
\tag{12}
$$

The final label map $\mathbf{S} = \{s_i\}_{i=1,\ldots,H\times W}$ is determined by $s_i = \arg\min_{l \in \mathcal{L}} p_{i,l}$.

# 3 Statistical Label Fusion Based on Local Structural Similarity

As described in Sec. 1, our OAN-RC is based on 2D images which is an extreme case of 3D anisotropic patches. In this section, we propose to fuse anisotropic information obtained from different viewing directions using isotropic 3D local patches to estimate the final segmentation. Let us denote the segmentation results by $\mathbf{S}^j, (j = 1, \ldots, M = 3)$, which are obtained as described in Sec. 2.3 from the axial (Z), sagittal (X), and coronal (Y) OAN-RCs. Depending on the viewing directions, sectional images contain different structures and may have different texture patterns in the same organs. These differences can cause nonidentical segmentations by the deep network as shown in Fig. 6 in 3D. In addition, there is no guarantee of connectivity between neighbor slices by independent use of slices for training and testing. Possible naïve approaches for determining the final segmentation in 3D from the OAN-RC results can be boolean operations such as union or intersection. Majority voting (MV) is another candidate for efficient fusion, however, theses approaches assume the same global weights of OAN-RC results. From the observations that the performance level of segmentation, e.g. sensitivity, can be different from viewing directions for each organ, we set the performance level to be an unknown variable when computing the probability of labeling. This concept is similar to the label fusion algorithms using expectation-maximization (EM) framework such as STAPLE (simultaneous truth and performance level estimation) and its extensions [1, 2, 36].

Let us denote the true label of the $V$ by $\mathbf{T}$, which is unknown, and the unknown performance level parameter of segmentation by $\theta$. The segmentations from the deep networks $\mathbf{S} = \{\mathbf{S}^j | j = 1, ..., M\}$ are

**Figure 6.** An example of multi-planar reconstruction view of OAN-RC estimations

observed values. Under this condition, the basic EM framework is performed by following two steps in an iterative manner: 1) to compute $Q^0(\theta|\theta^{(k)}) = E_{\mathbf{T}}\left[\ln L(\theta|\mathbf{S}, \mathbf{T})|\mathbf{S}, \theta^{(k)}\right]$ which is the expected value of the log likelihood, $\ln L(\theta|\mathbf{S}, \mathbf{T}) = \ln P(\mathbf{S}, \mathbf{T}|\theta)$, under the current estimate of the parameters $\theta^{(k)}$ at $k^{th}$ iteration, and 2) to find the parameter $\theta^{(k+1)}$ which maximizes $Q^0(\theta|\theta^{(k)})$.

The maximization step can be written as

$$
\begin{aligned}
\theta^{(k+1)} &= \arg\max_{\theta} E_{\mathbf{T}}\left[\ln P(\mathbf{S}, \mathbf{T}|\theta)|\mathbf{S}, \theta^{(k)}\right] \\
&= \arg\max_{\theta} E_{\mathbf{T}}\left[\ln P(\mathbf{S}|\mathbf{T}, \theta)P(\mathbf{T})|\mathbf{S}, \theta^{(k)}\right] \\
&= \arg\max_{\theta} \sum_{\mathbf{T}} \ln\left\{P(\mathbf{S}|\mathbf{T}, \theta)P(\mathbf{T})\right\} P(\mathbf{T}|\mathbf{S}, \theta^{(k)}) \\
&= \arg\max_{\theta} \sum_{\mathbf{T}} \left\{\ln P(\mathbf{S}|\mathbf{T}, \theta) + \ln P(\mathbf{T})\right\} P(\mathbf{T}|\mathbf{S}, \theta^{(k)}).
\end{aligned}
\tag{13}
$$

By assuming independence between $\mathbf{T}$ and $\theta$ in our problem, the second term $\sum_{\mathbf{T}} \ln P(\mathbf{T})P(\mathbf{T}|\mathbf{S}, \theta^{(k)})$ in (13) becomes free of $\theta$ and the maximization step can be written as

$$
\begin{aligned}
\theta^{(k+1)} &= \arg\max_{\theta} \sum_{\mathbf{T}} \ln P(\mathbf{S}|\mathbf{T}, \theta)P(\mathbf{T}|\mathbf{S}, \theta^{(k)}) \\
&= \arg\max_{\theta} E_{\mathbf{T}}\left[\ln P(\mathbf{S}|\mathbf{T}, \theta)|\mathbf{S}, \theta^{(k)}\right].
\end{aligned}
\tag{14}
$$

Therefore, we redefine $Q^0(\theta|\theta^{(k)})$ as $Q(\theta|\theta^{(k)}) = E_{\mathbf{T}}\left[\ln P(\mathbf{S}|\mathbf{T}, \theta)|\mathbf{S}, \theta^{(k)}\right]$.

The performance level parameter in this framework is a global property representing the overall confidence of deep network segmentation for the whole volume. However, it can also vary according to the voxel spatial locations via the local and neighbor structures as we use 2D slices for the initial segmentation.

Therefore, we propose to combine local structural similarity shown from a specific viewing direction to the original 3D volume and the global performance level, conceptually similar to local weighted voting [31]. We compute the probability of correspondence between 2D images and the 3D volume by structural similarity (SSIM) [35] by

$$
\begin{aligned}
\alpha_i^j &= P\left(\ell_2(I_i^j)|\ell_3(V_i)\right) \equiv SSIM\left(\ell_2(I_i^j), \ell_3(V_i)\right) \\
&= \frac{\left(2\mu_{\ell_2(I_i^j)}\mu_{\ell_3(V_i)} + c_1\right)\left(2\sigma_{\ell_2(I_i^j)\ell_3(V_i)} + c_2\right)}{\left(\mu_{\ell_2(I_i^j)}^2 + \mu_{\ell_3(V_i)}^2 + c_1\right)\left(\sigma_{\ell_2(I_i^j)}^2 + \sigma_{\ell_3(V_i)}^2 + c_2\right)},
\end{aligned}
\tag{15}
$$

where $\alpha_i^j$ is the SSIM from the $j^{th}$ viewing direction at the $i^{th}$ voxel. $c_1$ and $c_2$ are user-defined constants, and $\ell_2(I_i)$ and $\ell_3(V_i)$ represent local 2D and 3D patches centered at the $i^{th}$ voxel, respectively. $\mu_\ell$ and $\sigma_\ell$ are the average and standard deviation of the patch $\ell$, respectively, and $\sigma_{\ell_2(I_i)\ell_3(V_i)}$ is the covariance of $\ell_2(I_i)$ and $\ell_3(V_i)$. Fig. 7 shows an example of the structural similarity computed on different viewing directions as a color map.

Considering the local image properties, the expectation of log likelihood function in our problem becomes

$$
\begin{aligned}
Q\left(\theta|\theta^{(k)}\right) &= E\left[\ln P\left(\mathbf{S}, I|\mathbf{T}, V, \theta\right)|\mathbf{S}, I, V, \theta^{(k)}\right] \\
&= \sum_{\mathbf{T}} \ln P\left(\mathbf{S}, I|\mathbf{T}, V, \theta\right) P\left(\mathbf{T}|\mathbf{S}, I, V, \theta^{(k)}\right).
\end{aligned}
\tag{16}
$$

The global underlying performance level parameters of the deep network segmentations is defined as

$$
\theta_{js's} \equiv P\left(\mathbf{S}_i^j = s'|\mathbf{T}_i = s, \theta_{js's}^{(k)}\right),
\tag{17}
$$

where $\theta_{js's}$ is the probability of the voxel labeled as $s'$ from the $j^{th}$ deep network with the current estimated performance value $\theta_{js's}^{(k)}$, when the true label is $s$.

To make the problem simple, we assume conditional independence between labeling and the original volume intensities. The labeling probability with the target image intensity then becomes

$$
\begin{aligned}
&P\left(\mathbf{S}_i^j = s', \ell_2(I_i^j)|\mathbf{T}_i = s, \ell_3(V_i), \theta_{js's}^{(k)}\right) \\
&= P\left(\mathbf{S}_i^j = s'|\mathbf{T}_i = s, \theta_{js's}^{(k)}\right) P\left(\ell_2(I_i^j)|\ell_3(V_i)\right) \\
&= \theta_{js's}\alpha_i^j.
\end{aligned}
\tag{18}
$$

## 3.1 E-step

In the expectation step (E-step), we estimate the probability of voxelwise labels. Let us denote the probability that the true label of $i^{th}$ voxel is $s \in \mathcal{L}$ at the $k^{th}$ iteration by $\omega_{si}^{(k)}$. When the deep network segmentations $\mathbf{S}$ and performance level parameters at the $k^{th}$ iteration $\theta^{(k)}$ are given, $\omega_{si}^{(k)}$ can be then described as

$$
P\left(\mathbf{T}_i = s|\mathbf{S}, I, V, \theta^{(k)}\right) \equiv \omega_{si}^{(k)},
\tag{19}
$$

where $\theta \in \mathbb{R}^{N \times |\mathcal{L}| \times |\mathcal{L}|}$ is the vector of all $(\theta_{js's})^T$. From the independence between $\mathbf{S}^X$, $\mathbf{S}^Y$, and $\mathbf{S}^Z$, we apply Bayesian theorem to (19).

$$
\omega_{si}^{(k)} = \frac{P(\mathbf{T}_i = s)\prod_j P\left(\mathbf{S}_i^j = s', \ell_2(I_i^j)|\mathbf{T}_i = s, \ell_3(V_i), \theta_j^{(k)}\right)}{\sum_n P(\mathbf{T}_i = n)\prod_j P\left(\mathbf{S}_i^j = s', \ell_2(I_i^j)|\mathbf{T}_i = n, \ell_3(V_i), \theta_j^{(k)}\right)},
\tag{20}
$$

where $P(T_i = s)$ is a *priori* of the $i^{\text{th}}$ voxel. By applying (18) to (20), we then obtain the probability of voxelwise labeling as

$$\omega_{si}^{(k)} = \frac{P(\mathbf{T}_i = s) \prod_j \theta_{js's}^{(k)} \alpha_i^j}{\sum_n P(\mathbf{T}_i = n) \prod_j \theta_{js'n}^{(k)} \alpha_i^j}. \tag{21}$$

## 3.2   M-step

In the maximization step (*M-step*), the goal is to find the performance parameters, $\theta$, which maximize (16) with the current given parameters. Considering each $\mathbf{S}^j$ and $\theta_j$ independently, the expectation of log likelihood function in (16) can be expressed with the estimated voxelwise probability in *E-step*. Then the performance parameter of each segmentation can be formulated to find the solution which maximizes the summation of voxelwise probability as

$$\theta_j^{(k+1)} = \arg\max_{\theta_j} Q\left(\theta_j | \theta_j^{(k)}\right) = \arg\max_{\theta_j} \sum_i Q_i\left(\theta_j | \theta_j^{(k)}\right), \tag{22}$$

where $Q_i = E[\ln P(\mathbf{S}_i, \ell_2(I_i) | \mathbf{T}_i, \ell_3(V_i), \theta^{(k)}) | \mathbf{S}, I, V, \theta^{(k)}]$ at $i^{th}$ voxel. By applying (19) and (18), (22) becomes

$$
\begin{aligned}
\theta_j^{(k+1)} &= \arg\max_{\theta_j} \sum_i \sum_s P(\mathbf{T}_i = s | \mathbf{S}, I, V, \theta^{(k)}) \times \ln P\left(\mathbf{S}_i^j, \ell_2(I_i^j) | \mathbf{T}_i = s, \ell_3(V_i), \theta_j^{(k)}\right) \\
&= \arg\max_{\theta_j} \sum_i \sum_s \omega_{si}^{(k)} \ln P\left(\mathbf{S}_i^j, \ell_2(I_i^j) | \mathbf{T}_i = s, \ell_3(V_i), \theta_j^{(k)}\right) \\
&= \arg\max_{\theta_j} \sum_{s'} \sum_{i:\mathbf{S}_i^j = s'} \sum_s \omega_{si}^{(k)} \times lnP\left(\mathbf{S}_i^j = s', \ell_2(I_i^j) | \mathbf{T}_i = s, \ell_3(V_i), \theta_j^{(k)}\right) \\
&= \arg\max_{\theta_j} \sum_{s'} \sum_{i:\mathbf{S}_i^j = s'} \sum_s \omega_{si}^{(k)} \ln \theta_{js's} \alpha_i^j.
\end{aligned}
\tag{23}
$$

From the definition of $\theta$ in (17), the summation of probability mass function, $\sum_{s'} \theta_{js's}^{(k)}$, must be 1, and (22) becomes a constrained optimization problem which can be solved by introducing a Lagrange multiplier, $\lambda$. We then obtain the optimal solution by making the first gradient zero as

$$0 = \frac{\partial}{\partial \theta_{js's}} \left[ Q\left(\theta_j | \theta_j^{(k)}\right) + \lambda \sum_{s'} \theta_{js's} \right]. \tag{24}$$

By applying the derivation of $Q$ in (16), (22) and (23), (24) becomes

$$
\begin{aligned}
0 &= \frac{\sum_{i:\mathbf{S}_i^j = s'} \omega_{si}^{(k)} \alpha_i^j}{\theta_{js's}} + \lambda \\
\theta_{js's}^{(k+1)} &= \frac{\sum_{i:\mathbf{S}_i^j = s'} \omega_{si}^{(k)} \alpha_i^j}{-\lambda}.
\end{aligned}
\tag{25}
$$

By substituting the constraint of $\sum_{s'} \theta_{js's}^{(k)} = 1$, we can obtain the final optimal solution as

$$\theta_{js's}^{(k+1)} = \frac{\sum_{i:\mathbf{S}_i^j = s'} \alpha_i^j \omega_{si}^{(k)}}{\sum_i \omega_{si}^{(k)}}. \tag{26}$$

The two steps, (21) and (26), are then computed alternatively in the EM iterations until they converge. From the final values of (21), the final segmentation can be computed by graph-based approaches such as [3].

**Figure 7.** The local structural similarity map between 2D slices and the 3D volume. Each row is captured from the same similarity map computed on one viewing direction. Each column shows the captures images at the same location computed from different viewing directions.

## 3.3 Parallel computing using GPUs

The fusion step can be efficiently computed in a parallel way on a GPU. The local structural similarity $\alpha_i^j$ of $i$-th voxel in $j$th deep network and *priori* $P(T_i)$ can be computed for each voxel and saved as a pre-processing step. In the EM iterations, as shown in (21), the probability can be computed and updated for each structure at each voxel. In our implementation, a GPU thread is logically allocated for each voxel. However, to reduce the used memory and computation cost, the target volume of interest (VOI) for each structure $s$ is computed in an extended region as $\delta = 4$ voxels for each direction from $V(\bigcup_j \mathbf{S}^j = s)$ in our implementation. For parallel computing, one CPU thread is allocated to a structure and launches a kernel of one GPU to compute EM iteration for each structure.

## 4  Experimental Results

We evaluated our methods on 236 abdominal CT images of normal cases under an IRB (Institutional Review Board) approved protocol in Johns Hopkins Hospital as a part of the FELIX project for pancreatic cancer research [22]. CT images were obtained by Siemens Healthineers (Erlangen,Germany) SOMATOM Sensation and Definition CT scanners. CT scans are composed of $(319 - 1051)$ slices of $(512 \times 512)$ images, and have voxel spatial resolution of $([0.523 - 0.977] \times [0.523 - 0.977] \times 0.5)\, mm^3$. All CT scans are contrast enhanced images and obtained in the portal venous phase.

A total of 13 structures for each case were segmented by four human annotators/raters, one case by one person, and confirmed by an independent senior expert. The structures include the aorta, colon, duodenum, gallbladder, interior vena cava (IVC), kidney (left, right), liver, pancreas, small bowel, spleen, stomach, and large veins. Vascular structures were segmented only outside of the organs in order to make the structures exclusive to each other (i.e. no overlaps).

As explained in Sec. 2, we used OAN-RCs for multi-organ segmentation whose backbone FCNs had been pre-trained by $PascalVOC$ dataset [10]. From the possible variants of FCNs (e.g., FCN-32s, FCN-16s, and FCN-8s), which depend on how they combine the fine detailed predictions [34], we selected FCN-8s in this study because it captures very fine details in the $3^{rd}$ and $4^{th}$ pooling layer, and keeps high-level semantic contextual information from the final layer. Our algorithm was implemented and tested on a workstation with Intel i7-6850K CPU, NVidia TITAN X (PASCAL) GPU. With 236 cases, the initial segmentations using OAN-RCs were tested by four-fold cross-validation. All the input images of OAN-RCs are 1.5 times enlarged by upsampling, which lead to improved performance in our experiments.

In the fusion step, the average probability of $\mathbf{S}^X, \mathbf{S}^Y, \mathbf{S}^Z$ are taken as a *priors* in (21) and the initial performance levels $\theta_{js's}^{(0)}$ were computed by randomly selecting 5 cases and by comparing them to the ground-truth. To compute the local patch-based structural similarity in (15), patches of $(4.5 \times 4.5 \times 4.5)mm^3$ size cubes were used for 3D volume. Since CT voxels are not always isotropic and spatial resolutions can be different between scan volumes, we re-sampled the 3D patch with $0.5mm$ length cubic voxels so that the same size of $(9 \times 9 \times 9)$ 3D patches and $(9 \times 9)$ 2D patches from all directions can be used for all cases in our experiments.

The final segmentation results using OAN-RC with local structural similarity-based statistical fusion (LSSF) were compared with the 3D-patch based state-of-the-art approaches, 3D Unet [8] and hierarchical 3D FCN (HFCN) [30] as well as 2D-based FCN, OAN and OAN-RC with majority voting (MV). For a quantitative comparison, we computed the well-known Dice-Sørensen similarity coefficient (DSC) and the surface distances based on the manual annotations as ground-truth. For a structure $s$, DSC is computed as $\frac{2V(\mathbf{S}=s \bigcap \mathbf{T}=s)}{V(\mathbf{S}=s)+V(\mathbf{T}=s)}$ where $\mathbf{S}$ is the estimated segmentation and $\mathbf{T}$ is the ground-truth, i.e. manual annotations in this study. The surface distance was computed from each vertex of the ground-truth and to the estimates of our algorithms. Fig. 8 shows comparison results by box plots, while Tables 1 and 2 represent the mean and standard deviations for all the 236 cases.

As shown in Fig. 8, the basic OAN-RC outperforms other state-of-the-art approaches and our local structural similarity-based fusion improves the results even more. We note that although DSC shows the relative overall volume similarity, it does not quantify the boundary smoothness or the boundary noise of the results. But evaluating the surface distances, see below, shows that our method works effectively for

**Figure 8.** Box plots of the Dice-Sørensen similarity coefficients of 13 structures to compare performance. As in typical box plots, the box represents the first quartile, median, and the third quartile from the lower border, middle and the upper boarder, respectively, and the lower and the upper whiskers show the minimum and the maximum values. (LSSF: Local Similarity-based Statistical Fusion.)

**Table 1.** DICE-Sørensen similarity coefficient (DSC, %) of thirteen segmented organs. (mean ± standard deviation of 236 cases)

| Structure | 3D U-net | HFCN | FCN MV | OAN MV | OAN-RC MV | OAN-RC LSSF |
|---|---|---|---|---|---|---|
| Aorta | 87.0±12.3 | 88.3± 8.8 | 85.0±4.2 | 85.5± 4.2 | 85.3± 4.1 | **91.8**± 3.5 |
| Colon | 77.0±11.0 | 79.3± 9.2 | 80.3± 9.1 | 81.5± 9.4 | 82.0± 8.8 | **83.0**± 7.4 |
| Duodenum | 66.8±12.8 | 70.3± 10.4 | 70.2±11.3 | 72.6±11.4 | 73.4±11.1 | **75.4**± 9.1 |
| Gallbladder | 85.4±10.3 | 87.9± 7.5 | 87.8± 8.3 | 88.9± 6.2 | 89.4± 6.1 | **90.5**± 5.3 |
| IVC | 80.8±10.2 | 84.7± 5.9 | 84.0± 6.0 | 85.6± 5.8 | 86.0± 5.5 | **87.0**± 4.2 |
| Kidney(L) | 83.9±22.4 | 95.2± 2.6 | 96.1± 2.0 | 96.2± 2.2 | 95.9± 2.3 | **96.8**± 1.9 |
| Kidney(R) | 88.0±14.4 | 95.6± 4.5 | 95.8± 4.9 | 95.9± 4.9 | 96.0± 2.5 | **98.4**± 2.1 |
| Liver | 91.4± 9.9 | 95.7± 1.8 | 96.8± 0.8 | 97.0± 0.9 | 97.0± 0.8 | **98.0**± 0.7 |
| Pancreas | 79.3±11.7 | 81.4±10.8 | 84.3± 4.9 | 86.2± 4.5 | 86.6± 4.3 | **87.8**± 3.1 |
| Small bowel | 69.9±17.3 | 71.1±15.0 | 76.9±14.0 | 78.0±13.8 | 79.0±13.4 | **80.1**±10.2 |
| Spleen | 89.6±9.5 | 93.1± 2.1 | 96.3± 1.9 | 96.4± 1.9 | 96.4± 1.7 | **97.1**± 1.5 |
| Stomach | 90.1± 7.2 | 93.2± 5.4 | 93.9± 3.2 | 94.2± 2.9 | 94.2± 3.0 | **95.2**± 2.6 |
| Veins | 60.7±23.7 | 74.5±10.5 | 74.8±10.7 | 76.8±11.2 | 77.4±12.1 | **80.7**± 9.3 |

**Table 2.** Average surface distances of thirteen segmented organs for all 236 cases. (mean ± standard deviation of average surface distances in $mm$)

| Structure | 3D U-net | HFCN | FCN MV | OAN MV | OAN-RC MV | OAN-RC LSSF |
|---|---|---|---|---|---|---|
| Aorta | 0.44 ±1.01 | 0.42±0.58 | 0.56±0.47 | 0.47±0.42 | 0.44±0.28 | **0.39**±0.21 |
| Colon | 6.75±9.01 | 6,35±8.12 | 6.27±7.44 | 5.65±7.25 | 4.07±5.72 | **3.59**±4.17 |
| Duodenum | 2.01±2.46 | 1.70±2.18 | 1.71±2.25 | 1.49±1.87 | 1.54±1.43 | **1.36**±1.31 |
| Gallbladder | 1.31±0.76 | 1.21±0.50 | 1.22±0.52 | 1.12±0.50 | 1.05±0.41 | **0.95**±0.37 |
| IVC | 1.57±1.53 | 1.15±1.05 | 1.26±1.08 | 1.16±1.38 | 1.12±1.24 | **1.08**±1.03 |
| Kidney(L) | 0.77±1.04 | 0.41±0.42 | 0.36±0.47 | 0.34±0.47 | 0.30±0.33 | **0.30**±0.30 |
| Kidney(R) | 1.39±2.01 | 1.03±1.68 | 1.05±1.74 | 0.74±1.32 | 0.54±1.09 | **0.45**±0.89 |
| Liver | 1.89±3.21 | 1.60± 0 | 1.61±2.98 | 1.39±2.64 | 1.32±1.74 | **1.23**±1.52 |
| Pancreas | 1.78±1.05 | 1.51±0.80 | 1.41±0.88 | 1.19±0.82 | 1.17±0.72 | **1.05**±0.65 |
| Small bowel | 4.21±5.78 | 4.01±6.01 | 3.91±6.05 | 3.20±4.05 | 3.37±5.48 | **3.01**±3.35 |
| Spleen | 0.98±0.56 | 0.59±0.37 | 0.60±0.36 | 0.56±0.40 | 0.47 ±0.27 | **0.42**±0.25 |
| Stomach | 2.78±5.89 | 2.50±5.02 | 2.51±5.13 | 2.36±5.65 | 1.88±1.64 | **1.68**±1.55 |
| Veins | 2.31±4.51 | 1.75±3.51 | 1.69±3.61 | 1.92±6.48 | 1.40±3.61 | **1.21**±3.05 |

both the whole volumes and the boundaries of the organs.

Tables 1 and 2 represent the mean and standard deviations of performance measures for 13 critical organs. Similar to the box plots, they show that our OAN-RCs with statistical fusion improves the overall mean performance and also reduces the standard deviations significantly.

The OAN-RC training and testing can be computed in parallel for each view direction. In our experiments, the training took 40 hours for $120,000$ iterations for 177 training cases and the average testing time for each volume was 76.73 seconds. The fusion time depended on the volume of the target structure, and the average computation time for 13 organs was 6.87 seconds.

# 5 Discussion

Multi-organ segmentation using OAN-RCs alone, without the statistical fusion, gave similar or better performance compared with the state-of-the-art approaches summarized in [16]. In the specific case of the pancreas, state-of-the-art methods showed (mean ± standard deviations) segmentation accuracies as $74.4 \pm 20.2(\%)$ on 140 cases [32], $78.5 \pm 14.0(\%)$ on 150 cases [16], $78.0 \pm 8.2(\%)$ on 82 cases [29] and $75.74 \pm 10.47(\%)$ (on the whole slice) versus $82.4 \pm 5.7(\%)$ (reduced region of interest) on 82 cases [39] in terms of DSC. We cannot make a direct comparison because in these datasets CT images and manual segmentations (i.e. annotation) for the ground-truth are different from each other. But our OAN-RCs segmentations on our larger dataset shows similar or better performances in terms of DSC. Among target organs, our performance on structures such as gallbladder and pancreas, whose sizes are relatively small and have particularly weak boundaries improves significantly from using basic FCNs or using OANs without reverse connections.

**Figure 9.** 3D photo-realistic rendering of the ground-truth (left) and the results from OAN-RC with statistical fusion (right). The aorta, duodenum, IVC, liver, kidneys, pancreas, duodenum, spleen, and stomach are rendered. The difference between our results and the ground-truth are almost visually indistinguishable. To differentiate adjacent organs and from manual segmentation, different color setting were applied to the our methods results.

Moreover, as shown in Sec. 4, our statistical fusion based on local structural similarity improves the overall segmentation accuracies in terms of both DSC and average surface distances. In particular, there are significant performance improvements for the minimum values as shown in Fig. 8, which helps explain the robustness of the algorithm. The differences can be depicted more clearly by visualizing the 3D surfaces as shown in Figs 10 - 11. The noise of the deep network segmentations is distributed over large regions, without much connectivity, and occasionally they show significantly different patterns. But our fusion step exploits structural similarity which outputs clean and smooth boundaries by effectively combining different information based on the local structure of the original 3D volume.

When applying our proposed method and interpreting the evaluation results, we must address several considerations.

As shown in our experiments, our proposed algorithm also outperforms 3D patch based approaches. But 3D (isotropic) patch-based approaches have several issues which make it hard to apply to this problem. To make bigger patch size, they require more parameters and hence require more training data or, if this is not available, significant data augmentation (.e.g, by scaling, rotation, and elastic deformation). In addition, there can be practical memory limitation on GPUs which restricts the expandable patch size. The limited patch size means that the deep networks receptive field sizes contains only limited local information which is problematic for multi-organ segmentation and the discontinuities between the patches also raises problems. It is possible that solutions to these three problems may make 3D patch based methods work better in the future. Unlike 3D approaches, the local structure-similarity used in our fusion method effectively combine the information from anisotropic patches to 3D at each voxel. Fig. 9 shows an example generated by our proposed algorithm, which is visually indistinguishable from manual segmentation for almost all target structures.

The ground-truth used in this study for training and evaluation was specified using manual annotations by human observers. It is well known that there can be significant inter-/intra-observer variations in manual segmentation. But, as explained before, the ground-truth was created by four human observers and checked by experts in a visual way, and we randomly divided testing groups in our 4-fold cross-validation to avoid biased comparison. However, it is still possible that inaccuracies due to human variability may affect the evaluation as well as the training. This can be further intensively explored as separate experiments.

Another possible consideration when applying the proposed approach is the image quality which

(a) Liver (upper) and Pancreas (lower)



(b) Pancreas (Upper) and Duodenum (lower)

**Figure 10.** Effects of local structural similarity-based statistical fusion (LSSF) for estimating 3D surfaces. From left to right, the manual segmentation (ground-truth), initial segmentations from OAN-RCs with X, Y, Z slices, and the results of our proposed algorithm with statistical fusion. (a) When $\mathbf{S}^X$, $\mathbf{S}^Y$, and $\mathbf{S}^Z$ show similar result, statistical fusion produces smoother and less-noisy boundaries. (b) Surface estimation examples when initial OAN-RCs give differing results. But our approach effectively fuses the information, exploiting the local structural similarity.

**Figure 11.** Examples of FCN, OAN, OAN-RC, and OAN-RC. The manual segmentation (ground-truth), FCN MV, OAN MV, OAN-RC MV, OAN-RC LSFF (from left to right). (a) Pancreas: DSC(%) and surface distances (mean± standard deviation in $mm$) to the ground-truth are 72.5 and $2.13 \pm 1.74$ (FCN MV), 77.2 and $1.90 \pm 1.77$ (OAN MV), 82.4 and $1.33 \pm 1.31$ (OAN-RC MV), and 85.5 and $0.71 \pm 0.81$ (OAN-RC LSSF), respectively. (b) Stomach: DSC(%) and surface distances (mean± standard deviation in $mm$) to the ground-truth are 92.5 and $2.44 \pm 1.27$ (FCN MV), 93.6 and $1.63 \pm 1.14$ (OAN MV), 94.9 and $2.25 \pm 1.30$ (OAN-RC MV), and 97.1 and $1.26 \pm 0.88$ (OAN-RC LSSF), respectively.

can affect both of manual annotations and deep network segmentation results. Various factors such as spatial resolution, level of artifacts and reconstruction kernels should be considered. The dataset used in this study has been collected between 2005 to 2009 in the same institute with control over the scanning parameters. As explained in Sec. 4, the CT protocol is the portal venous phase and the spatial resolution is almost isotropic. But different scanning parameters and artifacts may affect our algorithms performance when applied to other datasets.

The same issues about manual segmentations and image qualities can be raised in general segmentation and evaluations. Specifically for our proposed approach, especially in the fusion step, the way of computing *priori*, $P(T)$, used in (21) can in practice affect the final segmentation. But considering that the deep network segmentation results from different viewing-directions are independently obtained, the mean can be accepted in general. However, if the deep network segmentations show clear tendencies towards over-estimation or under-estimation, then different types of models for *priors* may need to be used in order to improve the final result for practical applications.

One of the main advantages of our algorithm is the efficient computation time. The segmentation of 13 organs of the whole volume takes similar to or less than 1 minute with better performance reported than the state-of-the-art methods [16]. Hence our approach can be practically useful in clinical environments.

# 6 Conclusion

In this paper, we proposed a novel framework for multi-organ segmentation using OAN-RCs with statistical fusion exploiting structural similarity. Our two-stage organ-attention network reduces uncertainties at weak boundaries, focuses attention on organ regions with simple context, and adjusts FCN error by training the combination of original images and OAMs. Reverse connections deliver abstract level semantic information to lower layers so that hidden layers can be assisted to contain more semantic information and give good results even for small organs. The results are improved by the statistical fusion, based on local structural similarity, which smooths our noise and removes biases leading to better overall

segmentation performance in terms of DSC and surface distances. We showed that our performance is better than previous state of the art algorithms. Our framework is not specific to any particular body region, but gives high quality and robust results for abdominal CTs, which are typically challenging regions due to their low contrast, large intra-/inter-variations, and different scales. In addition, the efficient computational time of our algorithm makes our approach practical for clinical environments such as CAD, CAS or RT.

# References

1. A. J. Asman and B. A. Landman. Formulating spatially varying performance in the statistical fusion framework. *IEEE Transactions on Medical Imaging*, 31(6):1326–1336, 2012.

2. A. J. Asman and B. A. Landman. Non-local statistical label fusion for multi-atlas segmentation. *Medical Image Analysis*, 17(2):194–208, 2013.

3. Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(11):1222–1239, 2001.

4. 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.

5. L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. 2016.

6. L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille. Attention to scale: Scale-aware semantic image segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3640–3649, 2016.

7. C. Chu, M. Oda, T. Kitasaka, K. Misawa, M. Fujiwara, Y. Hayashi, Y. Nimura, D. Rueckert, and K. Mori. Multi-organ segmentation based on spatially-divided probabilistic atlas from 3D abdominal CT images. *Lecture Notes in Computer Science*, 8150 LNCS(PART 2):165–172, 2013.

8. Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger. 3D U-net: Learning dense volumetric segmentation from sparse annotation. *Lecture Notes in Computer Science*, 9901 LNCS:424–432, 2016.

9. Q. Dou, H. Chen, Y. Jin, L. Yu, J. Qin, and P. A. Heng. 3D deeply supervised network for automatic liver segmentation from CT volumes. *Lecture Notes in Computer Science*, 9901 LNCS:149–157, 2016.

10. M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. "the PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results". "http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html".

11. M. Havaei, A. Davy, D. Warde-Farley, A. Biard, A. Courville, Y. Bengio, C. Pal, P.-M. Jodoin, and H. Larochelle. Brain Tumor Segmentation with Deep Neural Networks. *Medical Image Analysis*, 35:18–31, 2017.

12. T. Heimann, B. van Ginneken, M. Styner, Y. Arzhaeva, V. Aurich, C. Bauer, A. Beck, C. Becker, R. Beichel, G. Bekes, F. Bello, G. Binnig, H. Bischof, A. Bornik, P. Cashman, Y. Chi, A. Cordova, B. Dawant, M. Fidrich, J. Furst, D. Furukawa, L. Grenacher, J. Hornegger, D. Kainmüller, R. Kitney, H. Kobatake, H. Lamecker, T. Lange, J. Lee, B. Lennon, R. Li, S. Li, H. Meinzer, G. Nemeth, D. Raicu, A. Rau, E. van Rikxoort, M. Rousson, L. Rusko, K. Saddi, G. Schmidt, D. Seghers, A. Shimizu, P. Slagmolen, E. Sorantin, G. Soza, R. Susomboon, J. Waite, A. Wimmer, and I. Wolf. Comparison and evaluation of methods for liver segmentation from CT datasets. *IEEE Transactions on Medical Imaging*, 28:1251–1265, 2009.

13. J. E. Iglesias and M. R. Sabuncu. Multi-atlas segmentation of biomedical images: A survey. *Medical Image Analysis*, 24(1):205–219, 2015.

14. K. Jamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, and B. Glocker. Efficient multi-scale 3D CNN with fully connected CRF for accurate brain lesion segmentation. *Medical Image Analysis*, 36:61–78, 2017.

15. T. Kada, M. G. Linguraru, M. Hori, R. M. Summers, N. Tomiyama, and Y. Sato. Abdominal multi-organ segmentation from CT images using conditional shape-location and unsupervised intensity priors. *Medical Image Analysis*, 26(1):1–18, 2015.

16. K. Karasawa, M. Oda, T. Kitasaka, K. Misawa, M. Fujiwara, C. Chu, G. Zheng, D. Rueckert, and K. Mori. Multi-atlas pancreas segmentation: Atlas selection based on vessel structure. *Medical Image Analysis*, 39:18–28, 2017.

17. C. Kirbas and F. Quek. A review of vessel extraction techniques and algorithms. *ACM Computing Surveys*, 36:81–121, 2004.

18. T. Kong, F. Sun, A. Yao, H. Liu, M. Lu, and Y. Chen. RON: reverse connection with objectness prior networks for object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

19. D. Lesage, E. D. Angelini, I. Bloch, and G. Funka-Lea. A review of 3d vessel lumen segmentation techniques: Models, features and extraction schemes. *MEdical Image Analysis*, 13:819–845, 2009.

20. G. Li, X. Chen, F. Shi, W. Zhu, J. Tian, and D. Xiang. Automatic liver segmentation based on shape constraints and deformable graph cut in ct images. *IEEE Transactions on Image Processing*, 24:5315–5329, 2015.

21. J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.

22. C. Lugo-Fagundo, B. Vogelstein, A. Yuille, and E. K. Fishman. Deep learning in radiology: Now the real work begins. *Journal of the American College of Radiology*, 15:364–367, 2018.

23. A. M. Mharib, A. R. Ramli, S. Mashohor, and R. B. Mahmood. Survey on liver ct image segmentation methods. *Artificial Intelligence Review*, 37, 2012.

24. F. Milletari, N. Navab, and S. A. Ahmadi. V-Net: Fully convolutional neural networks for volumetric medical image segmentation. *Proceedings - 2016 4th International Conference on 3D Vision, 3DV 2016*, pages 565–571, 2016.

25. J. Nascimento and G. Carneiro. Multi-atlas segmentation using manifold learning with deep belief networks. *Proceedings - International Symposium on Biomedical Imaging*, 2016-June, 2016.

26. M. Rajchl, M. C. Lee, O. Oktay, K. Kamnitsas, J. Passerat-Palmbach, W. Bai, M. Damodaram, M. A. Rutherford, J. V. Hajnal, B. Kainz, and D. Rueckert. DeepCut: Object Segmentation from Bounding Box Annotations Using Convolutional Neural Networks. *IEEE Transactions on Medical Imaging*, 36(2):674–683, 2017.

27. H. Roth, M. Oda, N. Shimizu, H. Oda, Y. Hayashi, T. Kitasaka, M. Fujiwara, K. Misawa, and K. Mori. Towards dense volumetric pancreas segmentation in ct using 3d fully convolutional networks. 2017.

28. H. R. Roth, A. Farag, L. Lu, E. B. Turkbey, and R. M. Summers. Deep convolutional networks for pancreas segmentation in CT imaging. 2015.

29. H. R. Roth, L. Lu, A. Farag, A. Sohn, and R. M. Summers. Spatial aggregation of holistically-nested networks for automated pancreas segmentation. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 9901 LNCS:451–459, 2016.

30. H. R. Roth, H. Oda, Y. Hayashi, M. Oda, N. Shimizu, M. Fujiwara, K. Misawa, and K. Mori. Hierarchical 3d fully convolutional networks for multi-organ segmentation. 2017.

31. M. Sabuncu, B. Yeo, K. van Leemput, B. Fischi, and P. Golland. A generative model for image segmentation based on label fusion. *IEEE Transactions on Medical Imaging*, 29:1714–1729, 2010.

32. A. Saito, S. Nawano, and A. Shimizu. Joint optimization of segmentation and shape prior from level-set-based statistical shape model, and its application to the automated segmentation of abdominal organs. *Medical Image Analysis*, 2105.

33. A. A. A. Setio, F. Ciompi, G. Litjens, P. Gerke, C. Jacobs, S. J. Van Riel, M. M. W. Wille, M. Naqibullah, C. I. Sanchez, and B. Van Ginneken. Pulmonary Nodule Detection in CT Images: False Positive Reduction Using Multi-View Convolutional Networks. *IEEE Transactions on Medical Imaging*, 35(5):1160–1169, 2016.

34. E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):640–651, April 2017.

35. Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transcations on Image Processing*, 13(4):600–612, 2004.

36. S. K. Warfield, K. H. Zou, and W. M. Wells. Simultaneous truth and performance level estimation (STAPLE): An algorithm for the validation of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921, 2004.

37. R. Wolz, C. Chu, K. Misawa, M. Fujiwara, K. Mori, and D. Rueckert. Automated abdominal multi-organ segmentation with subject-specific atlas generation. *IEEE Transactions on Medical Imaging*, 32(9):1723–1730, 2013.

38. S. Xie and Z. Tu. Holistically-nested edge detection. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.

39. Y. Zhou, L. Xie, W. Shen, Y. Wang, E. K. Fishman, and A. L. Yuille. A fixed-point model for pancreas segmentation in abdominal CT scans. In *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 2017.

40. X. Zhuang and J. Shen. Multi-scale patch and multi-modality atlases for whole heart segmentation of MRI. *Medical Image Analysis*, 31:77–87, 2016.

41. C. Zu, Z. Wang, D. Zhang, P. Liang, Y. Shi, D. Shen, and G. Wu. Robust multi-atlas label propagation by deep sparse representation. *Pattern Recognition*, 63:511–517, 2017.