

PROCEEDINGS OF SPIE

SPIEDigitalLibrary.org/conference-proceedings-of-spie

Super-mask-based object localization for auto-contouring in head and neck radiation therapy planning

Yubing Tong, Jayaram K. Udupa, Drew A. Torigian

Yubing Tong, Jayaram K. Udupa, Drew A. Torigian, "Super-mask-based object localization for auto-contouring in head and neck radiation therapy planning," Proc. SPIE 10951, Medical Imaging 2019: Image-Guided Procedures, Robotic Interventions, and Modeling, 109512L (8 March 2019); doi: 10.1117/12.2511973

SPIE.

Event: SPIE Medical Imaging, 2019, San Diego, California, United States

Super-mask-based object localization for auto-contouring in head and neck radiation therapy planning

Yubing Tong¹, Jayaram K. Udupa^{1*}, Drew A. Torigian¹

¹Medical Image Processing Group, 602 Goddard building, 3710 Hamilton Walk, Department of Radiology, University of Pennsylvania, Philadelphia, PA 19104, United States.

*Corresponding author.

Abstract

We have presented a variety of methods for object recognition based on the Automatic Anatomy Recognition (AAR) framework at previous SPIE conferences, including AAR recognition via optimal threshold on intensity, AAR recognition via composite information from intensity and texture, and AAR recognition with the optimal hierarchical structure design, and via neural networks to learn object relationships. The purpose of this paper is to introduce new features for the AAR-based recognition procedure and improve the performance of object localization for auto-contouring in head and neck (H&N) radiation therapy planning, specifically for some of the most challenging objects. The proposed super-mask technique first registers images used for model building among themselves optimally by using a minimal spanning tree in the complete graph formed with images as nodes to determine the order of registering images. Subsequently, we build a super-mask by combining the similarly registered binary images corresponding to each object by taking (**S1**) union of all binary images, (**S2**) intersection among all binary images, or (**S3**) the voting-based fuzzy mask created by adding the binary images. The super-mask is then used to confine search for optimum localization of the object in the given image. A large-scale H&N computed tomography (CT) data set with 216 subjects and over 2000 3D object samples were utilized in this study. The super-mask-based object localization approach within the AAR framework improved the recognition accuracy by 25-45% compared with the previous AAR strategy, especially for the most challenging H&N objects. On low quality images, the new method achieves recognition accuracy within 2 voxels on 50-60% of the cases.

Key words: Automatic anatomy recognition (AAR), optimal spanning tree, texture, image quality, head and neck cancer, radiation therapy, computed tomography (CT)

1. Introduction

During the radiation therapy (RT) treatment planning process, segmentation of organs at risk (OARs) is a key step. Furthermore, almost all quantitative medical image analysis depends on accurate object segmentation. Automatic Anatomy Recognition (AAR) is a body-wide multiple object segmentation approach [1] for which segmentation is designed as two dichotomous steps: object recognition (or localization) and object delineation. Recognition is the high-level process of determining the whereabouts of an object, and delineation is the meticulous low-level process of precisely indicating the space occupied by an object. Object segmentation can be improved by separately improving recognition and delineation processes. This study focuses on improving object recognition.

A variety of methods for object recognition based on the AAR framework have been presented at previous SPIE conferences, including AAR recognition via optimal threshold on intensity, texture, optimal hierarchical structure design [2 - 5], and AAR recognition with neural networks to learn object relationships [6]. We presented methods to recognize not only solid objects, but also lymph node zones which do not show intensity boundaries in images [7, 8]. The purpose of this paper is to introduce new features for the AAR-based recognition procedure and improve the performance of object localization for auto-contouring in head and neck (H&N) RT planning. Our efforts to improve AAR recognition (AAR-R) performance include three parts as follows:

- 1) Combining texture and intensity information in AAR recognition and significantly improve recognition performance.
- 2) Exploring strategies of integrating different image-based information including intensity, texture, and super-mask for object recognition.

- 3) Evaluating AAR–R performance in the context of image quality on a very large number of studies and 3D object samples (216 and over 2000, respectively). This is by far the largest study on real-world clinical H&N CT data sets that demonstrates object localization accuracy as a function of object/image quality. Instead of only a minority of images being used for testing, most of the images (185/216 = 86%) are used for testing in this study. The size of our data sets is significantly larger than the size involved in even some deep-learning approaches.

2. Materials and Methods

Image data

This retrospective study was conducted following approval from the Institutional Review Board at the Hospital of the University of Pennsylvania along with a Health Insurance Portability and Accountability Act (HIPAA) waiver. Data sets from the Department of Radiation Oncology, University of Pennsylvania, are utilized in this study, which are planning CT studies from 216 H&N cancer patients from among existing patient cases. The ground truth contour data for the cases were created by dosimetrists in the process of routine RT planning of these patients. The non-serial data sets constitute 54 cases gathered from each of four groups: 40-59-year-old males and females (denoted G_{M1} and G_{F1} , respectively), and 60-79-year-old males and females (denoted G_{M2} and G_{F2} , respectively). We corrected contours in only those cases with gross deviations from our standardized definitions of H&N OARs [9]. The CT images have a scene size of $512 \times 512 \times 110-140$, and a voxel size of $0.93 \times 0.93 \times 1.5-2 \text{ mm}^3$ to $1.6 \times 1.6 \times 3 \text{ mm}^3$. The total number of 3D CT scans involved in this study was thus 216, and the number of 3D object samples was 2199.

Objects and Fuzzy Model building

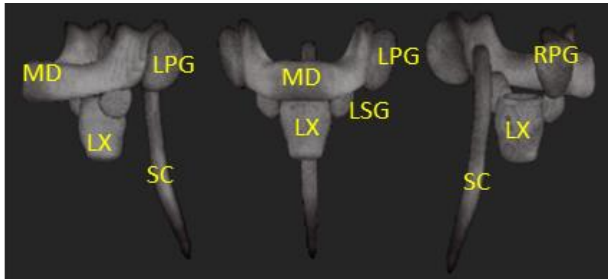


Figure 1. Volume-rendered 3D fuzzy object models for MD, PG (LPG, RPG), LSG, LX, and SC.

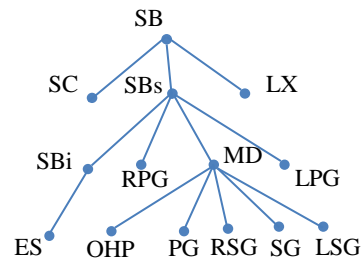


Figure 2. The optimal hierarchy, H_{opt} .

The anatomic objects considered in this study are the same as those in [4, 5], including Skin outer Boundary (SB) which was further sub-divided into an inferior portion below the neck (SBi) and a superior portion (SBs) in the neck, Left and Right Parotid Glands and their union called Parotid Glands (LPG, RPG, PG), Left and Right Submandibular Glands and their union called Submandibular Glands (LSG, RSG, SG), Esophagus (ES), supraglottic/glottic Larynx (LX), Spinal Canal (SC), Mandible (MD), and Oropharynx constrictor muscle (OHP). We illustrate results of all objects although more attention is given to LPG, RPG, PG, LSG, RSG, SG, LX, and OHP since recognition of those objects is still very challenging for the previous approaches.

The *Fuzzy Anatomy Model* of the H&N body region B for a group G , $FAM(B, G) = (H, M, \rho, \lambda, \eta)$, was built from the binary and gray scale images following mostly the methodology in [1]. Note that H denotes a hierarchical arrangement of the objects; M is a set of fuzzy models with one model for each object; ρ represents the parent to offspring relationship in G in the hierarchy; λ is a set of scale ranges one for each object; η includes a host of parameters representing object properties such as the range of variation of image intensity, etc. of each object. Figure 1 shows some samples of models constructed and used in this study, and Figure 2 shows the hierarchical structure used in this study, which is the optimal hierarchy, H_{opt} , and can be derived from the training data sets as detailed in [4, 5].

AAR recognition

The recognition process proceeds hierarchically following the order indicated by H_{opt} . The root object is first recognized following Ref [1]. To recognize any object O_i in the hierarchy, the fuzzy object model $FM(O_i)$ is first placed in the given image I at an initial location derived from the learned parent-child relationship with respect to the

parent object. This strategy is named as one-shot recognition since it uses only prior information. Then, the model is adjusted for its pose p (translation, scaling, and rotation) in I to best match a binary image resulting from thresholding I at a threshold that is optimal for O_i . This procedure of finding the optimal pose p^* can be described as follows.

$$p^* \in \underset{p}{\operatorname{argmin}}(|FM^p(O_i) - J| + |J - FM^p(O_i)|), \quad (1)$$

where J is a binary image resulting from thresholding I at the optimal threshold for O_i , and $| \cdot |$ denotes fuzzy subtraction operation between fuzzy model and J . The optimal threshold is found correspondingly in the image/ texture image [1, 4, 5]. Image and/or texture property that is best suited for each object is determined, and this information is stored in the model $FAM(B, G)$ in the element η and used at recognition. In this study, texture property “maximum probability of occurrence” derived from the co-occurrence matrix is used.

AAR super-mask technique

The super-mask technique makes use of the optimal atlas construction approach [10], which was designed to create atlases in an optimal manner by using a well-established graph theoretic approach based on a minimum spanning tree. The main steps include:

- 1) Create an undirected complete graph for a symmetric arc cost function or a directed graph for asymmetric arc cost function by calculating the cost (based on mean squared difference (MSD)) between each possible pair of images in the training set used for model building.
- 2) Calculate the minimum spanning tree (MST).
- 3) Perform registration among images hierarchically following the above tree.

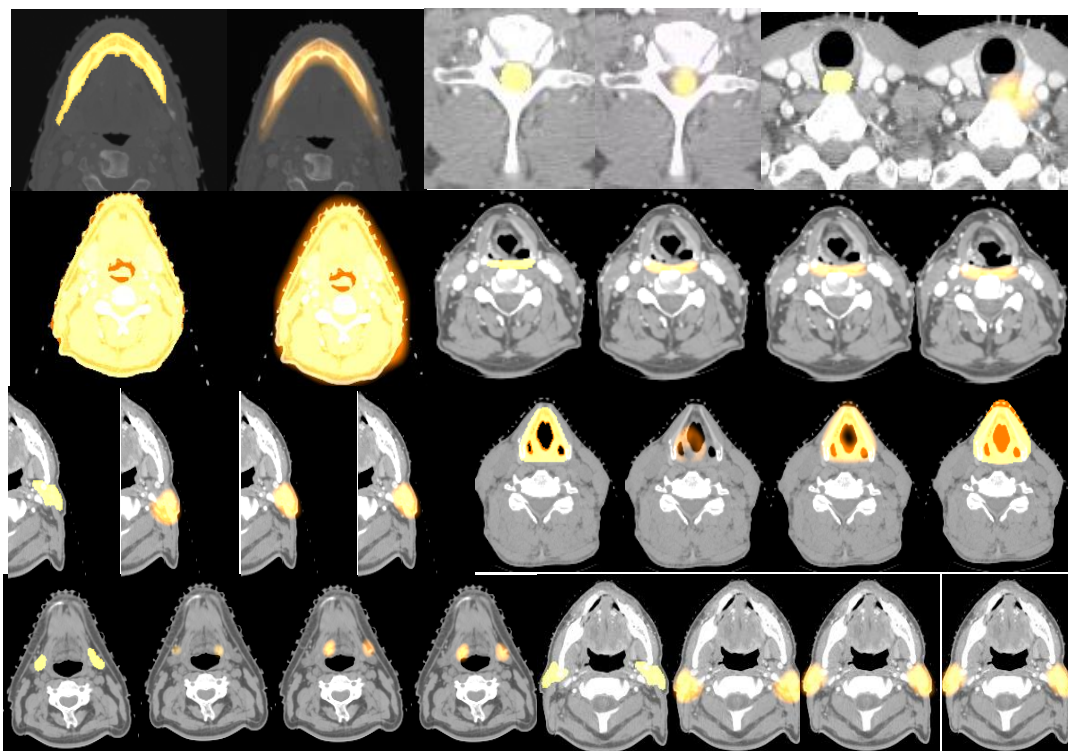


Figure 3. Representative recognition results from L1, L2, and L3. The 1st row and the leftmost two images in the 2nd row are from regular AAR recognition using only intensity information. The right four images of the 2nd row (left to right) show the ground truth and recognition results from L1, L2 and L3 for object OHP. The 3rd and 4th rows display the ground truth and recognition results from L1, L2, and L3 for LPG, LX, SG, and PG.

4) Build a super mask after the images are registered by combining the similarly registered binary images corresponding to each object by taking (**S1**) union of all binary images, (**S2**) intersection among all binary images, or (**S3**) the voting-based fuzzy mask created by adding the binary images.

In this study, we use a directed graph. The super-mask created by one of the three methods **S1-S3** is then used as a fuzzy model. After the one-shot strategy, the super-mask/model is dilated and combined with the thresholded image from intensity and texture for establishing the candidate binary image J in Equation (1) for determining region of search for optimal pose p^* .

3. Results

We test and compare among three methods. Method **L1**: Previous AAR recognition with intensity only. Method **L2**: Previous AAR recognition with intensity and texture. Method **L3**: Super-mask method as described above with three strategies (**S1-S3**) for creating the super-mask.

Among the 216 scans, we used 31 scans plus the associated ground truth segmentations of OARs from group G_{MI} for model building. All remaining 185 data sets were used for testing.

Qualitative results

Figure 3 shows representative recognition results from **L1**, **L2**, and **L3** for some OARs. The first row and the leftmost two images in the second row are from regular AAR recognition using only intensity information. The right four images of the second row (left to right) show the ground truth and recognition results from **L1**, **L2**, and **L3** for object OHP. The third and fourth rows display the ground truth and recognition results from **L1**, **L2**, and **L3** for objects LPG, LX, SG, and PG. Overall, **L2** and **L3** have achieved similar recognition results that are better than those of **L1**.

Quantitative results

Table 1. AAR recognition results by using different strategies for combining texture, intensity, and super-mask.

		SB	SBs	SBi	PG	LPG	RPG	SG	LSG	RSG	SC	OHP	MD	ES	LX	All
S1	LE	6.52	2.65	4.90	10.59	11.56	11.07	11.27	10.46	11.64	10.44	13.09	5.34	7.39	11.17	9.15
	SE	1.00	0.99	0.97	1.17	1.18	1.24	1.09	1.29	1.24	0.92	1.32	1.02	0.87	1.22	1.11
S2	LE	6.52	2.65	4.90	9.45	10.70	8.56	8.66	7.50	7.95	10.44	15.72	5.34	7.39	12.97	8.63
	SE	1.00	0.99	0.97	1.02	1.00	1.04	0.85	0.81	0.81	0.92	0.89	1.02	0.87	0.85	0.93
S3	LE	6.52	2.65	4.90	8.65	6.99	6.47	8.06	6.87	7.12	10.44	12.18	5.34	7.39	11.19	7.48
	SE	1.00	0.99	0.97	1.03	1.18	1.22	0.87	0.88	0.92	0.92	1.30	1.02	0.87	1.15	1.02

To determine which among **S1**, **S2**, and **S3** would give the best results for approach **L3** and the optimal amount of dilation needed, we used 23 data sets from group G_{MI} . Once the best strategy and values were determined from this preliminary test, they were fixed and the full battery of tests was carried out on all remaining 162 data sets. Table 1 shows recognition results achieved by using the three strategies **S1-S3** in the preliminary test. All strategies used the same dilation parameter ($\Delta = 15$). Strategy (**S3**) achieved the best recognition results among the three strategies.

Location error (LE) and scale error (SE) are used to quantitatively evaluate recognition results. LE is the distance (in mm) of the geometric center of the object model at recognition to the known true geometric center of the object. SE is the ratio of the estimated object size to its true size. The ideal values for these factors are 0 mm and 1, respectively. **S3** achieves the best results among three strategies.

With strategy **S3**, sensitivity of the dilation parameter (Δ) to recognition was tested with results presented in Table 2, where only the most challenging objects are listed. Recognition results of AAR super-mask (**L3**) using **S3** with dilation parameter Δ varied from 5 to 40 are compared. Not surprisingly, different objects seem to need different dilation parameters (see the last row in Table 2) for best performance since different objects show different amount of variation over a population and hence lead to different amount of fuzziness.

Table 2. Recognition results of super-mask approach L3 (with strategy S3) with dilation parameter Δ from 5 to 40.

Δ		PG	LPG	RPG	SG	LSG	RSG	OHP	LX	All
5	LE	8.81	8.53	8.27	7.37	8.14	6.83	14.56	13.33	9.78
	SE	1.03	1.05	1.09	0.86	0.81	0.82	1.18	0.89	0.98
10	LE	8.57	7.95	7.41	7.35	7.68	6.59	13.75	12.41	9.19
	SE	1.03	1.16	1.20	0.86	0.85	0.88	1.26	1.04	1.06
20	LE	9.05	6.87	7.24	10.69	7.02	7.27	12.05	10.08	8.51
	SE	1.03	1.18	1.23	0.86	0.89	0.94	1.32	1.18	1.11
30	LE	8.97	7.88	8.59	10.88	7.25	7.30	10.85	10.23	8.72
	SE	1.03	1.18	1.23	0.87	0.90	0.94	1.32	1.22	1.12
40	LE	8.97	7.88	8.56	10.88	7.25	7.30	10.98	10.42	8.77
	SE	1.03	1.18	1.23	0.87	0.90	0.94	1.32	1.22	1.12
LX, OHP 30; Others 15	LE	8.65	6.99	6.47	8.06	6.87	7.12	10.85	10.23	8.17
	SE	1.03	1.18	1.22	0.87	0.88	0.92	1.32	1.22	1.11

AAR recognition comparison: LE and SE are reported in Table 3 to quantitatively show the recognition results from **L1**, **L2**, and **L3** on all 162 testing data sets. Interestingly, Table 3 also shows that different recognition strategies perform differently for different objects, where some objects such as SB, SBs, SC, and MD have good recognition

Table 3. Recognition results from different methods L1, L2, and L3.

		SB	SBs	SBi	PG	LPG	RPG	SG	LSG	RSG	SC	OHP	MD	ES	LX	Mean
L1	LE	6.52	2.65	4.9	11.73	9.14	8.25	14.89	13.12	10.07	10.44	14.57	5.34	7.39	20.43	9.96
	SE	1.00	0.99	0.97	1.01	1.18	1.23	0.85	0.86	0.85	0.92	1.24	1.02	0.87	0.83	0.99
L2	LE	6.52	2.65	4.9	8.97	7.88	8.56	10.88	7.25	7.3	10.44	12.75	5.34	7.39	11.32	8.01
	SE	1.00	0.99	0.97	1.03	1.18	1.23	0.87	0.9	0.94	0.92	1.32	1.02	0.87	1.22	1.03
L3	LE	6.52	2.65	4.9	8.65	6.99	6.47	8.06	6.87	7.12	10.44	12.18	5.34	7.39	11.19	7.48
	SE	1.00	0.99	0.97	1.03	1.18	1.22	0.87	0.88	0.92	0.92	1.3	1.02	0.87	1.15	1.02

results with only **L1** (intensity based), whereas other objects such as LX and PG had much improvement in terms of location error with **L2** and **L3**. **L2** achieved a 19.6% reduction in LE compared with that of **L1**, and **L3** achieved the best recognition results with around 24.9 % improvement of LE compared with that of **L1**. If we only consider the most challenging objects LX, PG (LPG, RPG), and SG (LSG, RSG) as shown in Table 3, the recognition performance of **L1** can be improved by up to 45.2% for LX and by up to 45.9% for SG by **L3**. Compared with **L2**, **L3** on average improves the performance by around 10% on those objects. We believe that these results on these most challenging objects are excellent. In the segmentation literature, some objects such as LX and OHP are not dealt with at all but are frequently needed in RT planning.

Recognition considering image quality: Image quality score (IQS) was introduced in segmentation evaluation in [11], where object quality score (OQS) was first derived from several quality factors such as streak artifacts, extent of pathology, intensity deviation, image noise, etc., and then OQS values were combined via a logical predicate to generate IQS for an image. For all 216 subjects in our cohort, images with low quality comprised 38.4% for the male group and 52.1% for the female group. AAR-R achieves a location error of less than 4 mm (~1.5 voxels in our studies) for good quality images. We observed that for the low-quality cases, 48% of object samples can still achieve an average location error within 5.4 mm (which is ~2 voxels) in the male group, and 60% of object samples can similarly achieve an average location error within 4.8 mm in the female group. We conclude that the proposed method works very well even on low quality clinical CT images with severe artifacts and/or pathology.

Most of the current research on image segmentation is focused on object delineation, and without specific work on H&N object recognition/localization. Compared with the current literature [12-15], our study is much larger than any other study reported in terms of object localization, and deals with data sets that constitute the real-world heterogeneity that exists in clinical H&N cases. Instead of only using a minority of images for testing, most of the images (185/216 = 86%) are used for testing in this study. Ref. [12] performed testing on 559 CT images, which only considered larger objects such as heart, liver kidneys in thoracic and abdominal regions and achieved location error of 5.4 voxels, and Ref [13-15] used less than 100 images for testing. Our performance based on a large number of testing data sets is comparable to or better than the results reported in those papers. More details are found in Table 4.

We are in the process of combining this well-localized object information with deep-learning strategies confined to recognized objects to significantly improve delineation accuracy for these very challenging H&N objects.

4. Conclusions

In this paper, we propose a new super-mask-based object localization approach within the AAR framework and improve recognition performance, especially for the most challenging head and neck OARs. We will further test the proposed approach on more H&N data sets in the future, as well as on data sets from other body regions such as thorax and abdomen for the RT application [4], including over 200 scans with more than 2000 3D object samples with ground truth and IQS and OQS.

Table 4. Summary and comparison of recent work on object recognition on CT images.

Approaches	Modality	No. of cases train/test	Target objects	No. of testing object samples	Voxel size (mm ³)	Image/ Object quality artifacts	Location error (mm)
2D bounding box detection on ensemble learning, CMIG, 2012 [12]	Body torso CT	101/559	heart, liver, spleen, left/right kidney	2795	0.625×0.625×0.625	Not mentioned, not illustrated	5.4 voxels
Voxel- and slice-based deep learning for localization, DLMIA, 2016 [11]	Body CT	405/49	Right kidney	49	0.5×0.5×1.5	Not mentioned, not illustrated	7.8±9.4
Three independent ConvNets for boundary detection, SPIE Medical Imaging, 2016 [15]	Cardiac/ chest/abdomen CT	50/50; 50/49; 100/100	Left ventricle; liver; heart, aortic arch, descending aorta	300-600	0.86×0.86×(1.00-3.20); 0.55×0.55×(0.9-2.00)	Not mentioned, not illustrated	4.9±2.8; 19.0±10.5; 5.3±3.2
Single ConvNet trained to localize the 2D bounding box, IEEE TMI 2017 [14]	Cardiac/ chest/abdominal CT	50/50; 50/49; 100/100	Left ventricle; liver; heart, aortic arch, descending aorta	300-600	0.86×0.86×(1.00-3.20); 0.55×0.55×(0.9-2.00)	Not mentioned, not illustrated	4.5±3.4; 16.9±11.5; 5.3±3.7
Improved AAR-recognition via intensity, texture, super-mask	H&N planning CT and re-planning CT	36/262	14 objects Head & neck	2610	0.93×0.93×2 to 1.6×1.6×3	Only 2 images were entirely free from streak artifacts; results are reported according to image/object quality	< 4 mm (~1.5 voxels) for good quality images; 6-12 mm (4-5 voxels) for low quality image

The current study has made the following contributions:

- 1) A new super-mask-based object localization approach is proposed for auto-contouring in H&N RT planning on CT images.

- 2) Within the AAR framework, a super-mask recognition approach is compared with previous methods and is shown to achieve the best results, especially for the most challenging H&N OARs.
- 3) We evaluate AAR-R performance in the context of image quality on a large number of studies (216 CT images with ~2200 3D object samples) and demonstrate that the proposed method works well even on low quality images.

Acknowledgement:

This work was supported by grants from the National Science Foundation [IIP1549509] and National Cancer Institute [R41CA199735-01A1].

References

- [1] Jayaram K. Udupa, Dewey Odhner, et al., Body-Wide Hierarchical Fuzzy Modeling, Recognition and Delineation of Anatomy in Medical Images. *Medical Image Analysis* 18, 752-771, 2014.
- [2] Yubing Tong, Jayaram K. Udupa et al., Recognition of Upper Airway and Surrounding Structures at MRI in Pediatric PCOS and OSAS, *Proceeding of SPIE, Medical Imaging*, Vol.8670, 86702S1-7, 2013.
- [3] Yubing Tong, Jayaram K. Udupa et al., Fat segmentation on chest CT images via Fuzzy models, *Proceeding of SPIE, Medical Imaging*, Vol. 9786, 9786091-6, 2016.
- [4] Xingyu Wu, Jayaram K. Udupa, et al., Auto-contouring via automatic anatomy recognition of organs at risk in head and neck cancer on CT images, *Proceeding of SPIE, Medical Imaging*, Vol. 10576, 10576171-7, 2018.
- [5] Yubing Tong, Jayaram K. Udupa, et al., Hierarchical model-based object localization for auto-contouring in head and neck radiation therapy planning, *Proceeding of SPIE, Medical Imaging*, Vol. 10578, 105781-6, 2018.
- [6] Fengxia Yan, Jayaram K. Udupa, et al., Automatic anatomy recognition using neural network learning of object relationships via virtual landmarks, *Proceeding of SPIE, Medical Imaging*, Vol. 10574, 105742O1-6, 2018.
- [7] Guoping Xu, Jayaram K. Udupa, et al., Thoracic lymph node station recognition on CT images based on automatic anatomy recognition with an optimal parent strategy, *Proceeding of SPIE, Medical Imaging*, Vol. 10574, 105742F1-6, 2018.
- [8] Monica M. S. Matsumoto; Niha G. Beig, et al., Automatic localization of IASLC-defined mediastinal lymph node stations on CT images using fuzzy models, Vol. 90350, 90350J1-6, 2014.
- [9] <http://www.mipg.upenn.edu/Vnews/BodyRegionsObjects/HeadNeckObjects.pdf>.
- [10] George J. Grevera, Jayaram K. Udupa, et al., Optimal atlas construction through hierarchical image registration, *Proceeding of SPIE, Medical Imaging*, Vol.9786, 97861-7, 2018.
- [11] Gargi V. Pednekar, Jayaram K. Udupa et al., Image quality and segmentation, *Proceeding of SPIE, Medical Imaging*, Vol. 10576, 105761-6, 2018.
- [12] Zhou X, Wang S, Chen H, Hara T, Yokoyama R, Kanematsu M, Fujita H. Automatic localization of solid organs on 3D CT images by a collaborative majority voting decision based on ensemble learning. *Comput Med Imaging Graph.* 2012;36(4):304-13. doi: 10.1016/j.compmedimag. 2011.12.004. PubMed PMID: 22421130.
- [13] de Vos BD, Wolterink JM, de Jong PA, Leiner T, Viergever MA, Isgum I. ConvNet-Based Localization of Anatomical Structures in 3-D Medical Images. *IEEE Trans Med Imaging.* 2017;36(7):1470-81. doi: 10.1109/TMI.2017.2673121. PubMed PMID: 28252392.
- [14] de Vos BD, Wolterink JM, de Jong PA, Viergever MA, and Igum I, 2D image classification for 3D anatomy localization: employing deep convolutional neural networks. *SPIE Med. Imag.*, vol. 9784, 2016, pp. 97 841Y–97 841Y–7.
- [15] Lu X, Xu D, and Liu D, Robust 3D Organ Localization with Dual Learning Architectures and Fusion, in *Deep Learning and Data Labeling for Medical Applications*. Ser. Lecture Notes in Computer Science. Carneiro G, Mateus D, Peter L, Bradley A, Tavares JMRS, Belagiannis V, Papa JP, Nascimento JC, Loog M, Lu Z, Cardoso JS, and Cornebise J, Eds. Springer International Publishing, Oct. 2016, no. 10008, pp. 12–20.