

Towards Automatic Abdominal MRI Organ Segmentation: Leveraging Synthesized Data Generated From CT Labels

Cosmin Ciausiu¹, Deepa Krishnaswamy¹, Benjamin Billot², Steve Pieper³, Ron Kikinis¹, and Andrey Fedorov¹

¹ Brigham and Women’s Hospital, Boston MA 02115, USA,

² Massachusetts Institute of Technology, Cambridge MA 02139, USA,

³ Isomics, Cambridge MA 02138, USA

Abstract. Deep learning has shown great promise in the ability to automatically annotate organs in magnetic resonance imaging (MRI) scans, for example, of the brain. However, despite advancements in the field, the ability to accurately segment abdominal organs remains difficult across MR. In part, this may be explained by the much greater variability in image appearance and severely limited availability of training labels. The inherent nature of computed tomography (CT) scans makes it easier to annotate, resulting in a larger availability of expert annotations for the latter. We leverage a modality-agnostic domain randomization approach, utilizing CT label maps to generate synthetic images on-the-fly during training, further used to train a U-Net segmentation network for abdominal organs segmentation. Our approach shows comparable results compared to fully-supervised segmentation methods trained on MR data. Our method results in Dice scores of 0.90 ± 0.08 and 0.91 ± 0.08 for the right and left kidney respectively, compared to a pretrained nnU-Net model yielding 0.87 ± 0.20 and 0.91 ± 0.03 . We will make our code publicly available.

Keywords: image segmentation · domain randomization · computed tomography · magnetic resonance imaging · abdominal

1 Introduction

Accurate segmentation of abdominal organs in magnetic resonance (MR) images would be beneficial for many clinical tasks including liver volumetry [15], kidney disease monitoring [24], adaptive radiotherapy [29]. However, manual delineations of organs by experts are often time-consuming and tedious to perform [17]. The use of supervised convolutional neural network (CNN) methods for segmentation help alleviate the aforementioned issues, being time-efficient and robust to in-domain training data.

Popular supervised CNN methods for segmentation are the U-Net architecture [27] and its variants, including V-Net [25], and others built on top such as nnU-Net [19]. Specifically for abdominal segmentation, multiple approaches have

been developed based on these architectures. A multi-2D slice input approach was used to train a U-Net-based neural network, segmenting ten abdominal organs [5], and others used a 3D U-Net based approach [2] for MR pancreas segmentation. Going beyond the traditional encoder-decoder segmentation models, others have integrated the use of conditional generative adversarial networks (GAN) [8]. There has also been multiple challenges related to the development of supervised algorithms for multi-modality abdominal segmentation, such as AMOS [20] and CHAOS [21–23]. Most of the supervised methods can adapt well to the domain of the training dataset, but can fail on out-of-distribution data [11]. This is a significant concern as MR imaging data is highly heterogeneous with regards to resolution, orientation, and soft tissue contrast. Additionally, these methods require the need for large training datasets, which may not be readily available.

To rectify the need of large training datasets, techniques like data augmentation have been used to increase the heterogeneity of the data. Many techniques have been developed, ranging from basic transforms, to deformable and other learning-based methods [6, 31]. Though they may work well on a downstream segmentation task within a single modality, they may suffer when used in a cross-modality setting. This has led to the development of data augmentation techniques specifically for cross-modality use cases [4]. From that point, the field of domain adaptation has expanded rapidly, with the development of methods to account for domain shift between training and testing datasets [16]. For instance, a cross-modality domain adaptation method such as SIFA [3], modifies the input image domain to appear like the target domain using a GAN approach. Our proposed method uses a domain randomization approach, alleviating the need for defined source and target domain.

The proposed method implements a deep learning-based algorithm specifically for modality agnostic abdominal organ segmentation. The method contributes the following: 1) Leveraging publicly available computed tomography (CT) label maps for synthesizing training data for primarily MR segmentation, using a domain randomization approach 2) Extensive validation and testing of our approach on publicly available datasets 3) Analyzing the effect of the number of labels for contextual information and the level of granularity of the labels on the performance of the network 4) Comparison to multiple methods publicly available in the literature.

2 Methods

2.1 Approach

We adopt the domain randomization strategy of SynthSeg [1] for the automatic segmentation of abdominal structures in MR data. The method was originally developed for the segmentation of brain structures in MR and CT volumes, and was extended to the delineation of cardiac structures [1]. Using solely label maps as an input, the method generates synthetic data which is then used to train a U-Net model [27] for segmentation. Synthetic scans are generated by

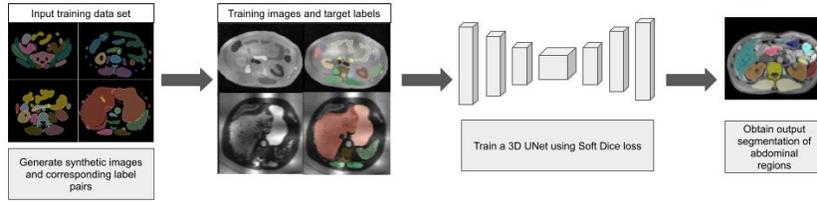


Fig. 1. Flowchart of the method for abdominal segmentation. The model requires as an input training label maps as input which are used to create synthetically generated images and corresponding label maps on the fly. These are used to train a U-Net model using a soft Dice loss. The final segmentation map of the anatomical structures of interests are obtained using the trained model.

sampling a Gaussian mixture model process conditioned on the expert annotated label maps. These parameters include intensity, bias, resolution, and orientation. During the training process, volumes are generated on the fly, where the network sees randomized data at every step, thereby becoming agnostic to resolution and contrast. Please refer to the paper [1] for further details of the adopted domain randomization strategy. Figure 1 describes the overall methodology for our proposed abdominal organ segmentation approach.

2.2 Datasets

Training We use the publicly available data (v1) provided from the TotalSegmentator method for CT structure segmentation [30]. All of the data were re-sampled to 1.5 mm, and includes 104 anatomical regions. The authors used a semi-automatic process to generate ground-truth annotations, combining AI-produced segmentations and expert feedback. For our study, we used 10 subjects for training, as it has previously been shown that the generative model performs well despite a low number of training subjects [1]. We randomly chose 10 patients where the abdomen was present.

Table 1. Dataset split for validation and testing, for AMOS and CHAOS, in terms of number of subjects (subjs). AMOS contains MR and CT modalities, and CHAOS contains T1 in phase, T2 out phase and T2 along with CT.

Dataset	# MR val subjs	# MR test subjs	# CT val subjs	# CT test subjs
AMOS [20]	40	20	25	12
CHAOS [21–23]	24	24	10	10

Validation and testing We used two publicly available collections for validation and testing of our approach. The AMOS dataset [20] is a collection

of patients with abdominal cancer or abnormalities, where both CT and MR data were collected from two medical centers using eight different scanners. The dataset includes labels for 15 abdominal organs. A coarse segmentation was first performed using a model trained on a small sample of the data, and the output of these models were refined by junior radiologists for the remaining data, and reviewed by senior radiologists. We used a subset of this data for validation and testing, as detailed in Table 1.

The second publicly available dataset we used was CHAOS [21–23], which consists of CT and MR images from healthy patients at the Dokuz Eylul University Hospital, in Izmir, Turkey. MR sequences included T1 in phase, T1 out phase and T2 SPIR, and includes labels for the liver, spleen, right kidney and left kidney, while the CT dataset only includes labels for the liver. These images were annotated from three radiologists with a majority voting procedure. We used a subset of the provided training dataset for our validation, and the provided validation set for our testing. We split the data into validation and testing cohorts as described in Table 1.

Pre-processing for target labels selection The original TotalSegmentator training data contains whole-body CT images and label maps, with 104 structures segmented. In our study, we only focus on the abdominal organs, therefore we cropped and/or padded our images and label volumes to a fixed size of 300x300x250. The following organs were selected for segmentation: liver, spleen, kidneys, stomach, duodenum, pancreas, gallbladder, small bowel, colon, adrenal glands, sacrum, hip bone, gluteus maximus, gluteus medius, gluteus minimus, autochthon, iliopsoas. We combined the vertebrae present in the abdominal region into a single segment for prediction. This yielded a total of 26 predicted labels (including left and right separately).

Pre-processing for synthetic data generation Our U-Net segmentation network is trained on synthetic scans generated from our training label maps. In order to diversify the appearance of our synthetic images, we applied a number of additional processing steps to increase the variety of details generated. From the original 10 label maps, we removed the CT table from 50% of the subjects using a 3DSlicer extension [13]. To further enhance the appearance of fine structures in CT scans later used for clustering, pre-processing steps including Gaussian filter blurring (https://docs.scipy.org/doc/scipy/reference/generated/scipy.ndimage.gaussian_filter.html) and gamma-based contrast-stretching was applied to the CT images used to generate synthetic data during training. Additional labels only used during synthetic data generation were added, to add more randomization to our training images. These labels were obtained by performing intensity clustering on the background and the foreground of our CT training scans. A Gaussian Mixture Model (GMM) with $K_{background} \in [3, 4, 5]$ components was fitted on the background voxels of our training CT scans. This GMM was then optimized using the Expectation-Maximization (EM) algorithm [9], providing K clusters assignments for all of our background CT images

Algorithm 1 Background and foreground labels clustering used for synthetic image generation

```

for  $BG$  in  $[3,4,5]$  do
   $P(x) = \sum_{k=1}^{BG} \mathcal{N}(x | \mu_k, \sigma_k^2) \pi_k$ 
  Optimise  $\theta_{BG} = \{\mu_k, \sigma_k^2, \pi_k, k = 1, \dots, BG\}$  using EM algorithm
  for  $FG$  in  $[1,2,3]$  do
    for segment in  $[1, \dots, N_{seg}]$  do
       $P(x_{segment}) = \sum_{k=1}^{FG} \mathcal{N}(x_{segment} | \mu_k, \sigma_k^2) \pi_k$ 
      Optimise  $\theta_{FG} = \{\mu_k, \sigma_k^2, \pi_k, k = 1, \dots, FG\}$  using EM algorithm
    end for
  end for
end for

```

N_{seg} refers the total number of labels present in a particular CT scan and corresponding label map

voxels. These additional background clusters were used during our synthetic data generation step. The same approach was used to produce additional foreground label clusters, for $K_{foreground} \in [1, 2, 3]$ for each available segment from the TotalSegmentator label maps. Algorithm 1 describes the process of background and foreground clustering GMM process used to generate additional labels for synthetic data generation. Using this process, we obtained 180 training label maps from the original 10 selected label maps.

2.3 Training procedure and evaluation

Following the published method developed by [1] from here (<https://github.com/BBillot/SynthSeg>), we trained the U-Net segmentation network on our synthetic images using the default parameters. We trained our network for two weeks using an NVIDIA A100 GPU with 40 GB RAM [28] for 100 epochs with 5000 steps per epoch. Epoch 10 was chosen for testing. Overlap metrics and distance-based metrics were computed for testing and validation of our methods, namely Dice score [10] and Hausdorff distance [18].

3 Results and Discussion

3.1 Experimental Results

We performed inference on the external testing collections outlined in Table 1. Table 2 provides the quantitative results for the same collections for the liver, spleen and kidneys. We observe high Dice scores for AMOS MR and CT modalities. However, CHAOS MR yields relatively lower results than AMOS. This could be partly due to the differences between the ground truth segmentations, as both the TotalSegmentator and AMOS collections did not contain the renal cavity of the kidney, while CHAOS did. Figure 2 displays our qualitative results on one sample subject, showing the ground truth segments vs predictions from

Table 2. Quantitative results for inference on two collections AMOS and CHAOS for each modality (MR and CT) in terms of mean Dice score (Dice) and mean Hausdorff distance (95th percentile) (HD95) in mm. The standard deviation values are also provided.

Dataset	Segment	MR	MR	CT	CT
		Dice	HD95	Dice	HD95
AMOS	Liver	0.90 ± 0.04	25.13 ± 26.0	0.91 ± 0.05	13.15 ± 15.4
AMOS	Spleen	0.86 ± 0.15	5.34 ± 6.37	0.78 ± 0.22	19.37 ± 26.20
AMOS	Right kidney	0.90 ± 0.08	4.24 ± 5.50	0.91 ± 0.05	2.83 ± 1.84
AMOS	Left kidney	0.91 ± 0.08	2.92 ± 2.04	0.92 ± 0.03	3.55 ± 2.84
CHAOS	Liver	0.87 ± 0.05	5.77 ± 3.27	0.91 ± 0.10	21.08 ± 23.95
CHAOS	Spleen	0.78 ± 0.13	10.46 ± 15.06	—	—
CHAOS	Right kidney	0.80 ± 0.14	3.53 ± 2.29	—	—
CHAOS	Left kidney	0.68 ± 0.31	7.79 ± 11.19	—	—

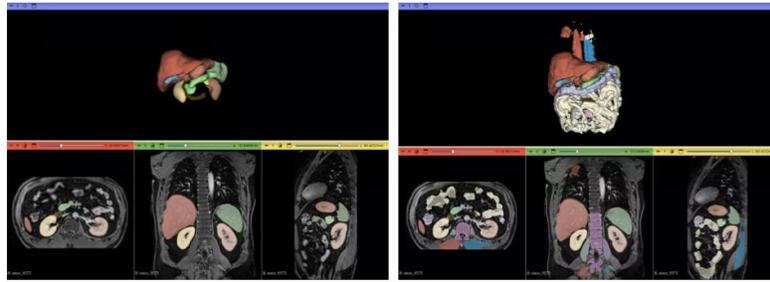


Fig. 2. Qualitative results of the proposed method on a subject from AMOS MR. The left view shows the expert radiologist annotations, and the right one shows the results from our baseline method. Here we have high agreement between many abdominal regions, including the liver in red, spleen in green, kidneys in yellow and brown. Our method segments 26 structures in the abdominal area.

our proposed method. Please refer to the supplementary material for additional box plots results.

3.2 Ablation study results

In order to assess the performance of the model and to understand the effect of the synthetic generation process, we performed a number of ablation studies. We studied the effect of foreground labels clustering and the effect of segmenting one label versus all abdominal labels. The three studies we performed were 1) No foreground clustering and segmenting all original abdominal labels, 2) Include foreground clustering and predict a single label (left and right kidneys), 3) No foreground clustering and predict a single label (left and right kidneys).

After examining the validation curves, epoch 5 was chosen for both A. and B. The last epoch was chosen for C. Please refer to the supplementary material for analysis of the validation loss curves. Overall our baseline method performed

Table 3. Dice score and HD95 quantitative testing results for inference on AMOS and CHAOS for MR and CT modalities. The last four rows are our baseline method. Subjects with undefined Hausdorff distance values were removed for the overall metrics calculation. Values in bold are the best performing. Kruskal-Wallis statistical testing was performed comparing our baseline approach to the three ablation studies. Values with an asterisk indicate a significant difference of $p < 0.01$. C for the first column refers to clusters (Yes/No). The L column refers to the the number of labels, all or 1, RK refers to right kidney segmentation, LK refers to left kidney segmentation.

C	L	Dataset	RK Dice	RK HD95	LK Dice	LK HD95
N	1	AMOS MR	0.58 ± 0.34	162.87 ± 42.94	0.25 ± 0.15	203.44 ± 59.35
N	1	AMOS CT	0.0 ± 0.0	250.08 ± 31.09	0.02 ± 0.06	285.07 ± 19.28
N	1	CHAOS MR	0.29 ± 0.36	94.12 ± 36.03	0.14 ± 0.16	111.78 ± 15.88
Y	1	AMOS MR	0.82 ± 0.22	5.46 ± 6.68	0.85 ± 0.19	5.06 ± 6.90
Y	1	AMOS CT	0.81 ± 0.21	5.55 ± 8.12	0.84 ± 0.23	4.82 ± 7.36
Y	1	CHAOS MR	0.51 ± 0.30	12.48 ± 10.73	0.76 ± 0.27	5.98 ± 7.20
N	all	AMOS MR	0.81 ± 0.27	8.05 ± 9.89	0.87 ± 0.08	6.65 ± 4.79
N	all	AMOS CT	0.76 ± 0.26	10.69 ± 16.03	0.68 ± 0.28	37.75 ± 97.07
N	all	CHAOS MR	0.63 ± 0.25	7.88 ± 6.75	0.62 ± 0.26	6.73 ± 4.90
Y	all	AMOS MR	$0.90 \pm 0.08^*$	$4.24 \pm 5.50^*$	$0.91 \pm 0.08^*$	$2.92 \pm 2.04^*$
Y	all	AMOS CT	$0.91 \pm 0.05^*$	$2.83 \pm 1.84^*$	$0.92 \pm 0.03^*$	$3.55 \pm 2.84^*$
Y	all	CHAOS MR	$0.80 \pm 0.14^*$	$3.53 \pm 2.29^*$	$0.68 \pm 0.31^*$	$7.79 \pm 11.19^*$

the best compared to the three ablation studies. Table 3 shows that our baseline experiment including all abdominal organs and the additional clustering method for the synthetic generation step is significantly better than the ablation methods, across all collections. One could argue that the clustering step adds more diversity to the synthesized images used for training, and including more segmentation labels during training helps to add contextual information. Please refer to the supplementary material for qualitative results.

3.3 Comparison to publicly available methods

We also compare our baseline method to two other publicly available abdominal segmentation methods. We use data from 23 patients from the TCGA-LIHC collection [7, 12] as part of NCI Imaging Data Commons (IDC) [14], which have AI-generated annotations of the liver from BAMF Health (Grand Rapids, MI) [26] that were assessed by a radiologist to be reasonable. We computed the same overlap and distance metrics between our baseline method segmentations and the ones available in IDC [14]. For the second comparison, we used a pre-trained model available from the nnU-Net [19]. The nnU-Net framework provides a large number of pre-trained models, including a model trained on CHAOS MR data for segmentation of the liver, spleen, left and right kidneys. We compare the performance of our method to the nnU-Net model on the AMOS MR dataset, using expert annotated ground-truth segmentations available for AMOS MR.

Table 4 displays the quantitative results between our baseline approach and the BAMF method on IDC data. AMOS MR results are also shown for

Table 4. Quantitative results comparing publicly available methods to our approach using the Dice score and HD95.

Method	Dataset	Segment	Dice	HD95
BAMF [26]	TCGA-LIHC (MR)	Liver	0.92 ± 0.02	7.30 ± 2.67
nnU-Net model [19]	AMOS MR	Liver	0.83 ± 0.23	31.04 ± 57.29
nnU-Net model [19]	AMOS MR	Spleen	0.88 ± 0.20	14.99 ± 52.49
nnU-Net model [19]	AMOS MR	Right kidney	0.87 ± 0.20	16.04 ± 40.12
nnU-Net model [19]	AMOS MR	Left kidney	0.91 ± 0.03	7.62 ± 1.82

the evaluated nnU-Net model and our baseline method, compared to expert annotations. We observe a high overlap between the liver segmentations from BAMF and our baseline segmentations, on IDC MR TCGA-LIHC data (0.92 ± 0.02 DSC). On AMOS MR, we can observe comparable results to the nnU-Net pre-trained model, when evaluated against expert annotations (0.91 ± 0.08 vs 0.91 ± 0.03 DSC for the left kidney, ours and nnU-Net, respectively). Please refer to the supplementary material for qualitative results of the evaluated methods.

4 Conclusion

We proposed a modality-agnostic deep learning method for abdominal organ segmentation using a domain randomization strategy trained on CT label maps. Our method shows promising results when validating and testing on publicly available datasets, as well as a comparison to publicly available fully-supervised segmentation methods. Additionally, we performed an ablation study to understand the effect of prediction of multiple labels and the addition of clustering. Our baseline method, including additional labels and generation-only foreground and background clustering labels, performed the best for both the overlap and distance metrics, namely Dice score and Hausdorff distance (95th percentile), compared to the ablation studies methods. Our study presents a few limitations, as abdominal structures outside of the liver, spleen and kidneys did not perform well. This could be due to the higher heterogeneity in those regions with regards to texture, appearance and location, and requires further investigation. Future work includes modification of the training data to better represent MR-specific features, and refinement of the ground truth labels in order to offer a better label anatomical consensus between collections.

References

1. Billot, B., Greve, D., Puonti, O., Thielscher, A., Van Leemput, K., Fischl, B., Dalca, A., Iglesias, J.: Synthseg: Segmentation of brain mri scans of any contrast and resolution without retraining. *Medical image analysis* **86**, 102789 (2023)
2. Bobo, M., Bao, S., Huo, Y., Yao, Y., Virostko, J., Plassard, A., Lyu, I., Assad, A., Abramson, R., Hilmes, M., Landman, B.: Fully convolutional neural networks improve abdominal organ segmentation. In *Medical Imaging 2018: Image Processing*, SPIE **10574**, 750–757 (2018)

3. Chen, C., Dou, Q., Chen, H., Qin, J., Heng, P.: Synergistic image and feature adaptation: Towards cross-modality domain adaptation for medical image segmentation. In Proceedings of the AAAI conference on artificial intelligence **33**(1), 865–872 (2019)
4. Chen, X., Lian, C., Wang, L., Deng, H., Kuang, T., Fung, S., Gateno, J., Shen, D., Xia, J., Yap, P.: Diverse data augmentation for learning image segmentation with cross-modality annotations. Medical image analysis **71**, 102060 (2021)
5. Chen, Y., Ruan, D., Xiao, J., Wang, L., Sun, B., Saouaf, R., Yang, W., Li, D., Fan, Z.: Fully automated multiorgan segmentation in abdominal magnetic resonance imaging with deep neural networks. Medical physics **47**(10), 4971–82 (2020)
6. Chlap, P., Min, H., Vandenberg, N., Dowling, J., Holloway, L., Haworth, A.: A review of medical image data augmentation techniques for deep learning applications. Journal of Medical Imaging and Radiation Oncology **65**(5), 545–63 (2021)
7. Clark, K., Vendt, B., Smith, K., Freymann, J., Kirby, J., Koppel, P., Moore, S., Phillips, S., Maffitt, D., Pringle, M., Tarbox, L.: The cancer imaging archive (tcia): maintaining and operating a public information repository. Journal of digital imaging **26**, 1045–57 (2013)
8. Conze, P., Kavur, A., Cornec-Le Gall, E., Gezer, N., Le Meur, Y., Selver, M., Rousseau, F.: Abdominal multi-organ segmentation with cascaded convolutional and adversarial deep networks. Artificial Intelligence in Medicine **117**, 102109 (2021)
9. Dempster, A., Laird, N., Rubin, D.: Maximum likelihood from incomplete data via the em algorithm. Journal of the royal statistical society: series B (methodological) **39**(1), 1–22 (1977)
10. Dice, L.: Measures of the amount of ecologic association between species. Ecology **26**(3), 297–302 (1945)
11. Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: Decaf: A deep convolutional activation feature for generic visual recognition. In International conference on machine learning pp. 647–655 (2014)
12. Erickson, B., Kirk, S., Lee, Y., Bathe, O., Kearns, M., Gerdes, C., Rieger-Christ, K., Lemmerman, J.: The cancer genome atlas liver hepatocellular carcinoma collection (tcga-lihc) (version 5) [data set] (2016). <https://doi.org/https://doi.org/10.7937/K9/TCIA.2016.IMMQW8UQ>, <https://wiki.cancerimagingarchive.net/pages/viewpage.action?pageId=6885436>
13. Fedorov, A., Beichel, R., Kalpathy-Cramer, J., Finet, J., Fillion-Robin, J., Pujol, S., Bauer, C., Jennings, D., Fennessy, F., Sonka, M., Buatti, J.: 3d slicer as an image computing platform for the quantitative imaging network. Magnetic resonance imaging **30**(9), 323–41 (2012)
14. Fedorov, A., Longabaugh, W., Pot, D., Clunie, D., Pieper, S., Gibbs, D., Bridge, C., Herrmann, M., Homeyer, A., Lewis, R., Aerts, H.: National cancer institute imaging data commons: Toward transparency, reproducibility, and scalability in imaging artificial intelligence. Radiographics **43**(12), e230180 (2023)
15. Gotra, A., Sivakumaran, L., Chartrand, G., Vu, K., Vandembroucke-Menu, F., Kauffmann, C., Kadoury, S., Gallix, B., de Guise, J., Tang, A.: Liver segmentation: indications, techniques and future directions. Insights into imaging **8**(4), 377–392 (2017)
16. Guan, H., Liu, M.: Domain adaptation for medical image analysis: a survey. IEEE Transactions on Biomedical Engineering **69**(3), 1173–85 (2021)
17. Heerkens, H., Hall, W., Li, X., Knechtges, P., Dalah, E., Paulson, E., van den Berg, C., Meijer, G., Koay, E., Crane, C., Aitken, K.: Recommendations for mri-

- based contouring of gross tumor volume and organs at risk for radiation therapy of pancreatic cancer. *Practical radiation oncology* **7**(2), 126–36 (2017)
18. Huttenlocher, D., Klanderman, G., Rucklidge, W.: Comparing images using the hausdorff distance. *IEEE Transactions on pattern analysis and machine intelligence* **15**(9), 850–63 (1993)
 19. Isensee, F., Jaeger, P., Kohl, S., Petersen, J., Maier-Hein, K.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* **18**(2), 203–11 (2021)
 20. Ji, Y., Bai, H., Ge, C., Yang, J., Zhu, Y., Zhang, R., Li, Z., Zhanng, L., Ma, W., Wan, X., Luo, P.: Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *Advances in Neural Information Processing Systems* **35**, 36722–32 (2022)
 21. Kavur, A., Gezer, N., Barış, M., Şahin, Y., Özkan, S., Baydar, B., Yüksel, U., Kılıkçer, c., Olut, c., Akar, G., Ünal, G.: Comparison of semi-automatic and deep learning-based automatic methods for liver segmentation in living liver transplant donors. *Diagnostic and Interventional Radiology* **26**(1), 11 (2020)
 22. Kavur, A., Gezer, N., Barış, M., Aslan, S., Conze, P., Groza, V., Pham, D., Chatterjee, S., Ernst, P., Özkan, S., Baydar, B.: Chaos challenge-combined (ct-mr) healthy abdominal organ segmentation. *Medical Image Analysis* **69**, 101950 (2021)
 23. Kavur, A., Selver, A., Dicle, O., Barış, M., Gezer, N.: Chaos - combined (ct-mr) healthy abdominal organ segmentation challenge data (2019). <https://doi.org/doi:10.5281/zenodo.3431873>, <https://zenodo.org/records/3431873>
 24. Kline, T., Edwards, M., Garg, I., Irazabal, M., Korfiatis, P., Harris, P., King, B., Torres, V., Venkatesh, S., Erickson, B.: Quantitative mri of kidneys in renal disease. *Abdominal Radiology* **43**, 629–38 (2018)
 25. Milletari F, Navab N, A.S.: Fully convolutional neural networks for volumetric medical image segmentation. In 2016 fourth international conference on 3D vision (3DV) pp. 65–571 (2016)
 26. Murugesan, G., McCrumb, D., Aboian, M., Verma, T., Soni, R., Memon, F., Van Oss, J.: The aimi initiative: Ai-generated annotations for imaging data commons collections (2023). <https://doi.org/arXiv:2310.14897>, <https://arxiv.org/abs/2310.14897>
 27. Ronneberger, O., Fischer, P., Brox, T.: -net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, Proceedings, Part III* 18 pp. 234–241 (2015)
 28. Stewart, C., Cockerill, T., Foster, I., Hancock, D., Merchant, N., Skidmore, E., Stanzione, D., Taylor, J., Tuecke, S., Turner, G., Vaughn, M.: Jetstream: a self-provisioned, scalable science and engineering cloud environment. In *Proceedings of the 2015 XSEDE Conference: Scientific Advancements Enabled by Enhanced Cyberinfrastructure* pp. 1–8 (2015)
 29. Tetar, S., Bruynzeel, A., Lagerwaard, F., Slotman, B., Bohoudi, O., Palacios, M.: Clinical implementation of magnetic resonance imaging guided adaptive radiotherapy for localized prostate cancer. *Physics and imaging in radiation oncology* **9**, 69–76 (2019)
 30. Wasserthal, J., Breit, H., Meyer, M., Pradella, M., Hinck, D., Sauter, A., Heye, T., Boll, D., Cyriac, J., Yang, S., Bach, M.: Totalsegmentator: Robust segmentation of 104 anatomic structures in ct images. *Radiology: Artificial Intelligence* **5**(5), 102789 (2023)

31. Yun, S., Han, D., Oh, S., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In Proceedings of the IEEE/CVF international conference on computer vision pp. 6023–6032 (2019)

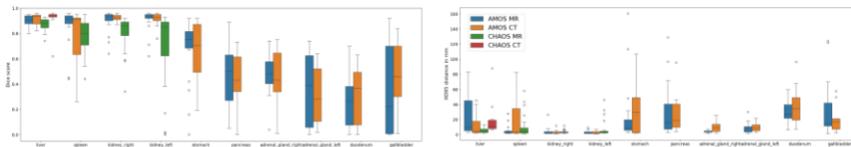


Fig. 1. Boxplots for the distribution of Dice scores (left) and the distribution of Hausdorff distance values (95th percentile) for all segments, between our proposed method and the expert annotations for the AMOS and CHAOS collections. Note the high agreement for the liver, spleen and kidneys, and the low agreement for the others. This may be due to the reduced ability of the network to capture variable appearance of those regions. Predicted segments with undefined Hausdorff distance values (due to empty segmentation sets) were removed from all of the metrics calculation. 9 segment predictions for AMOS MR/CT.

Table 1. Table for the distribution of Dice scores (left) and the distribution of Hausdorff distance values (95th percentile) for all segments, for the AMOS and CHAOS collections. Subjects with undefined Hausdorff distance values were removed for the overall metrics calculation. Predicted segments with undefined Hausdorff distance values (due to empty segmentation sets) were removed from all of the metrics calculation. 9 segment predictions for AMOS MR and 9 segments for AMOS CT were removed.

Dataset	Segment	MR	MR	CT	CT
		Dice	HD95	Dice	HD95
AMOS	Liver	0.90 ± 0.04	25.13 ± 26.05	0.91 ± 0.05	13.15 ± 15.43
AMOS	Spleen	0.86 ± 0.15	5.34 ± 6.37	0.78 ± 0.22	19.37 ± 26.20
AMOS	Right kidney	0.90 ± 0.08	4.24 ± 5.50	0.91 ± 0.05	2.83 ± 1.84
AMOS	Left kidney	0.91 ± 0.08	2.92 ± 2.04	0.92 ± 0.03	3.55 ± 2.84
AMOS	Stomach	0.68 ± 0.24	25.38 ± 39.55	0.64 ± 0.24	34.06 ± 30.26
AMOS	Pancreas	0.46 ± 0.23	30.62 ± 34.36	0.44 ± 0.21	27.97 ± 25.82
AMOS	Duodenum	0.27 ± 0.21	31.32 ± 15.53	0.30 ± 0.22	37.88 ± 25.63
AMOS	Gallbladder	0.34 ± 0.33	33.66 ± 38.24	0.46 ± 0.27	21.19 ± 18.69
AMOS	Right adrenal gland	0.47 ± 0.15	4.21 ± 1.42	0.45 ± 0.22	9.89 ± 7.04
AMOS	Left adrenal gland	0.36 ± 0.27	9.19 ± 7.09	0.31 ± 0.23	9.99 ± 6.22
CHAOS	Liver	0.87 ± 0.05	5.77 ± 3.27	0.91 ± 0.10	21.08 ± 23.95
CHAOS	Spleen	0.78 ± 0.13	10.46 ± 15.06	—	—
CHAOS	Right kidney	0.80 ± 0.14	3.53 ± 2.29	—	—
CHAOS	Left kidney	0.68 ± 0.31	7.79 ± 11.19	—	—

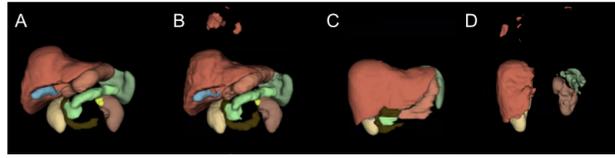


Fig. 2. Comparison of the segmentation results of the expert annotations to the proposed method, for two sample subjects from AMOS MR. A: Expert annotations; B: Proposed method on the subject as (A) with high agreement; C: Expert annotations; D: Proposed method on the same subject as (C) with low agreement.

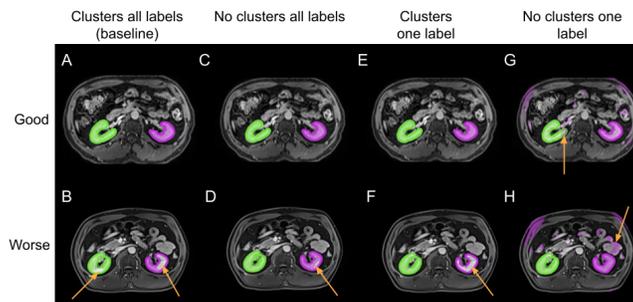


Fig. 3. displays the qualitative results of the kidney segmentation for the baseline method along with the three ablation studies models. The ground truth expert segmentations are given by the thick segment boundary line, while the AI predictions are filled. The right kidney is green, and the left kidney is purple. Yellow arrows indicate missing segments or incorrect segmentations. Note the higher overlap of the expert segmentations versus our AI models for methods that include all labels.

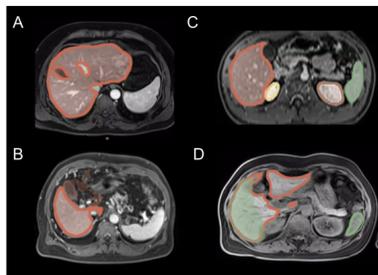


Fig. 4. Comparison of the proposed method with publicly available methods. A) Proposed method (filled) vs results from the BAMF method for liver segmentation (border) for a subject with high agreement. B) Proposed method (filled) vs results from the BAMF method for liver segmentation (border) for a subject with low agreement. C) Results from the nnU-Net model (filled) vs the ground truth (border) for a subject with high agreement. D) Results from the nnU-Net model (filled) vs the ground truth (border) for a subject with low agreement. Note the model predicts the liver (red) as the spleen (green) and vice versa.