

# Versatile Medical Image Segmentation Learned from Multi-Source Datasets via Model Self-Disambiguation

## Supplementary Material

### 1. More Details about Datasets

Details of the seven public datasets are provided in their corresponding papers. Regarding the private dataset, it comprises 122 contrast-enhanced CT images from patients undergoing urinary system examinations. The images have a uniform matrix size of  $512 \times 512$ , with a variable number of 2D slices ranging from 62 to 685. Pixel spacing ranges from 0.607 to 0.977 mm, and slice thickness varies from 1.0 to 3.0 mm. Urologists annotated the kidney, bladder, and ureters in each image. For this study, only the masks of the two kidneys and the bladder were retained.

### 2. More Details about Network Architecture

Table 10 presents the architecture for 3D TransUNet. The structure of 3D TransUNet is asymmetric, featuring a greater number of layers in the encoder compared to the decoder. Both the encoder and decoder are composed of 5 stages, wherein spatial sizes progressively decrease by 50% from stage 1 to stage 5 in a sequential manner.

The network’s building blocks are shown in brackets. All blocks, except those in stage 5, comprise two consecutive convolutional layers. The adjacent pair of numbers within each bracket represent the input channels and output channels of a convolutional layer. A skip connection [4] is added when the input channels of the first convolutional layer is different from the output channels of the second convolutional layer within each building block in stages 1–4. In accordance with [2], we employed weight normalization [1] in every convolutional layer to expedite training. Subsequent to each convolution operation, instance normalization [3] and rectified linear unit activation are applied. Downsampling and upsampling are executed through trilinear interpolation.

At the bottleneck, four multi-head attention layers were incorporated, each with eight heads. The size of each attention head for query, key, and value was set to be 512.

### 3. More Results on Partially Labeled Data

**Effect of patch size.** We conducted experiments with two different patch sizes, namely  $96 \times 96 \times 96$  and  $112 \times 112 \times 112$ . Larger patch sizes were not explored due to limitations in GPU memory. As indicated in Table 11, employing a patch size of  $96 \times 96 \times 96$  resulted in an average DSC of 87.9%, which is 0.8% DSC lower than the performance achieved with a patch size of  $112 \times 112 \times 112$ . These findings underscore the advantageous impact of using a larger

Table 10. Network architecture.

	Encoder	Decoder
Stage 1	{1, 32, 64}	{128, 64, 64}
Stage 2	{64, 64, 128} {128, 128, 256}	{128, 64, 64}
Stage 3	{256, 128, 256} {256, 128, 256} {256, 128, 512}	{512, 64, 64}
Stage 4	{512, 256, 512} {512, 256, 512} {512, 256, 512} {512, 256, 1024}	{1024, 256, 256}
Stage 5	{1024, 512}	{512, 512}

patch for abdominal organ segmentation, since increased patch size contributes to a more comprehensive inclusion of contextual information.

Table 11. Performance trained with different patch sizes.

Patch Size	DSC [%]
$96 \times 96 \times 96$	87.9 $\pm$ 8.1
$112 \times 112 \times 112$	<b>88.7<math>\pm</math>7.0</b>

**Effect of voxel spacing.** In our experiments, we standardized the voxel spacing for all images to facilitate model training. To assess the influence of voxel spacing on model performance, we conducted experiments with three different voxel spacings:  $1.5 \times 1.5 \times 1.5 \text{ mm}^3$ ,  $2.0 \times 2.0 \times 2.0 \text{ mm}^3$ , and  $2.5 \times 2.5 \times 2.5 \text{ mm}^3$ . As indicated in Table 12, employing a voxel spacing of  $1.5 \times 1.5 \times 1.5 \text{ mm}^3$  and  $2.5 \times 2.5 \times 2.5 \text{ mm}^3$  led to a performance decrease of 0.3% and 0.6% in terms of average DSC, respectively. The diminished performance with a voxel spacing of  $1.5 \times 1.5 \times 1.5 \text{ mm}^3$  can be attributed to reduced contextual information within the input image patch. Conversely, the inferior performance with a voxel spacing of  $2.5 \times 2.5 \times 2.5 \text{ mm}^3$  is likely due to information loss during downsampling, particularly impacting small structure segmentation.

### 4. More Results on Sparsely Labeled Data

Tables 13 and 14 provide a detailed comparison of the performance for each anatomical structure and across each dataset, respectively. Visual results of a randomly selected

Table 12. Performance trained with different patch sizes.

Voxel Size	DSC [%]
$1.5 \times 1.5 \times 1.5$	$88.4 \pm 7.7$
$2.0 \times 2.0 \times 2.0$	<b><math>88.7 \pm 7.0</math></b>
$2.5 \times 2.5 \times 2.5$	$88.1 \pm 7.8$

subject from each dataset are presented in the second and third columns of Fig. 5. These results align with the findings in Table 8, highlighting the consistent success of our method across different views. Notably, even with the utilization of only 20% of incompletely annotated slices for training, our method demonstrates commendable performance across the structures of interest and datasets.

## 5. More Results on Hybrid Data

Tables 15 and 16 present a comprehensive comparison of performance for each anatomical structure and across each dataset, respectively. Visual results of a randomly selected subject from each dataset are displayed in the fourth column of Fig. 5. These results concur with the findings in Table 9, underscoring the effectiveness of our method in utilizing a mixture of partially and sparsely labeled data for model training.

## References

- [1] Salimans, T. et al. (2016). Weight normalization: A simple reparameterization to accelerate training of deep neural networks. *Advances in neural information processing systems*, 29. 1
- [2] Chen, Jieneng, et al. (2021). Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*. 3, 1
- [3] Ulyanov, D. et al. (2016). Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*. 1
- [4] He, K. et al. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778). 1

Table 13. Performance (DSC, %) comparison on each anatomical structure using different portions of annotated slices.

Setting	View	Sp	RK	LK	GB	Eso	L	St	A	PC	Pan	RAG	LAG	Duo	B	PU	PSV	Average
20%	Axial	93.5	92.6	92.4	74.3	78.8	95.4	89.5	89.4	83.2	83.1	70.8	72.4	70.5	85.7	71.1	70.7	82.0
		$\pm 9.3$	$\pm 10.9$	$\pm 9.9$	$\pm 26.4$	$\pm 12.0$	$\pm 9.6$	$\pm 13.7$	$\pm 12.1$	$\pm 15.0$	$\pm 11.4$	$\pm 14.1$	$\pm 16.3$	$\pm 23.5$	$\pm 17.2$	$\pm 22.9$	$\pm 20.1$	$\pm 15.3$
	Sagittal	94.5	92.5	91.9	75.6	76.9	96.1	90.5	91.0	85.0	83.5	71.6	73.1	71.7	88.6	74.2	72.3	83.1
		$\pm 8.1$	$\pm 11.1$	$\pm 9.0$	$\pm 25.5$	$\pm 11.3$	$\pm 7.1$	$\pm 12.2$	$\pm 8.7$	$\pm 12.9$	$\pm 9.3$	$\pm 13.0$	$\pm 14.6$	$\pm 19.7$	$\pm 14.6$	$\pm 24.5$	$\pm 18.5$	$\pm 13.8$
	Coronal	94.8	92.6	92.7	75.1	76.5	96.0	90.7	90.6	85.5	83.5	71.3	71.6	73.3	88.2	73.3	74.3	83.1
		$\pm 6.8$	$\pm 10.8$	$\pm 8.0$	$\pm 24.3$	$\pm 12.6$	$\pm 7.9$	$\pm 12.0$	$\pm 7.8$	$\pm 13.5$	$\pm 10.5$	$\pm 14.6$	$\pm 16.6$	$\pm 18.6$	$\pm 15.3$	$\pm 20.2$	$\pm 17.6$	$\pm 13.6$
100%	Axial	95.0	93.5	93.9	76.1	80.7	96.6	91.3	91.7	86.2	84.6	74.5	74.9	75.0	89.2	76.2	76.8	84.8
		$\pm 7.0$	$\pm 8.4$	$\pm 5.8$	$\pm 26.1$	$\pm 10.2$	$\pm 5.4$	$\pm 11.5$	$\pm 9.0$	$\pm 12.3$	$\pm 8.8$	$\pm 12.0$	$\pm 14.3$	$\pm 18.3$	$\pm 13.6$	$\pm 20.1$	$\pm 15.8$	$\pm 12.4$
	Sagittal	94.6	93.1	93.7	76.1	80.4	96.2	91.2	91.1	86.4	84.6	72.7	73.9	73.9	89.0	76.9	75.8	84.3
		$\pm 8.6$	$\pm 8.7$	$\pm 5.1$	$\pm 25.1$	$\pm 8.6$	$\pm 7.1$	$\pm 12.0$	$\pm 8.8$	$\pm 10.4$	$\pm 9.9$	$\pm 14.2$	$\pm 15.2$	$\pm 18.2$	$\pm 14.4$	$\pm 23.8$	$\pm 17.8$	$\pm 13.0$
	Coronal	95.2	93.5	93.7	77.7	81.3	96.5	91.1	91.7	86.4	84.4	73.0	74.6	74.8	89.6	73.7	75.2	84.5
		$\pm 6.6$	$\pm 7.6$	$\pm 5.6$	$\pm 23.5$	$\pm 8.8$	$\pm 6.4$	$\pm 12.2$	$\pm 9.1$	$\pm 12.0$	$\pm 9.7$	$\pm 13.6$	$\pm 14.7$	$\pm 17.1$	$\pm 12.8$	$\pm 19.1$	$\pm 16.6$	$\pm 12.2$

Table 14. Performance (DSC, %) comparison on each dataset using different portions of annotated slices.

Setting	View	AbCT-1K	AMOS-CT	AMOS-MRI	BTCV	FLARE22	NIH-Pan	TotalSeg	Urogram	WORD
20%	Axial	92.7	82.7	80.0	76.3	89.8	83.5	79.3	92.7	82.0
		$\pm 2.0$	$\pm 8.1$	$\pm 10.1$	$\pm 5.8$	$\pm 1.5$	$\pm 4.7$	$\pm 15.2$	$\pm 2.6$	$\pm 4.4$
	Sagittal	92.7	83.4	80.5	76.5	89.4	83.2	82.0	92.3	82.7
		$\pm 1.9$	$\pm 7.6$	$\pm 8.5$	$\pm 6.5$	$\pm 1.5$	$\pm 5.3$	$\pm 11.3$	$\pm 3.0$	$\pm 4.4$
	Coronal	92.8	83.5	79.2	76.1	89.8	83.3	81.6	92.6	82.5
		$\pm 2.3$	$\pm 7.2$	$\pm 9.8$	$\pm 7.8$	$\pm 1.6$	$\pm 5.8$	$\pm 13.1$	$\pm 3.0$	$\pm 4.2$
100%	Axial	93.4	85.1	80.9	77.9	90.5	83.8	84.7	93.0	83.7
		$\pm 1.7$	$\pm 6.3$	$\pm 9.0$	$\pm 6.4$	$\pm 1.8$	$\pm 5.4$	$\pm 9.7$	$\pm 3.0$	$\pm 4.1$
	Sagittal	93.0	84.2	81.5	78.7	90.4	84.5	84.2	93.0	83.5
		$\pm 2.2$	$\pm 7.7$	$\pm 7.7$	$\pm 5.5$	$\pm 1.0$	$\pm 5.2$	$\pm 9.8$	$\pm 2.9$	$\pm 4.4$
	Coronal	93.3	84.5	81.1	78.1	90.3	84.4	85.0	93.0	83.7
		$\pm 2.0$	$\pm 6.9$	$\pm 8.3$	$\pm 5.7$	$\pm 1.3$	$\pm 4.2$	$\pm 9.5$	$\pm 3.2$	$\pm 4.0$

Table 15. Performance (DSC, %) comparison on each anatomical structure using mixed training.

View	Sp	RK	LK	GB	Eso	L	St	A	PC	Pan	RAG	LAG	Duo	B	PU	PSV	Average
Axial	95.1	93.4	93.8	75.8	81.6	96.4	90.7	91.7	86.1	84.5	74.6	75.7	73.4	89.0	76.6	73.2	84.5
	$\pm 9.3$	$\pm 10.9$	$\pm 9.9$	$\pm 26.4$	$\pm 12.0$	$\pm 9.6$	$\pm 13.7$	$\pm 12.1$	$\pm 15.0$	$\pm 11.4$	$\pm 14.1$	$\pm 16.3$	$\pm 23.5$	$\pm 17.2$	$\pm 22.9$	$\pm 20.1$	$\pm 15.3$
Sagittal	95.1	93.2	93.7	77.6	80.6	96.6	91.6	92.1	86.5	84.9	74.1	75.5	75.3	89.2	76.6	75.9	84.9
	$\pm 7.1$	$\pm 10.3$	$\pm 6.6$	$\pm 24.0$	$\pm 10.7$	$\pm 6.5$	$\pm 11.5$	$\pm 7.1$	$\pm 11.2$	$\pm 8.5$	$\pm 13.7$	$\pm 14.3$	$\pm 17.5$	$\pm 14.7$	$\pm 26.7$	$\pm 17.4$	$\pm 13.0$
Coronal	95.1	93.8	93.8	77.2	80.6	96.4	91.3	91.9	86.3	84.7	73.9	75.7	75.1	89.5	77.6	76.7	85.0
	$\pm 6.8$	$\pm 6.4$	$\pm 6.0$	$\pm 25.0$	$\pm 11.3$	$\pm 7.3$	$\pm 12.4$	$\pm 7.8$	$\pm 12.4$	$\pm 8.9$	$\pm 14.2$	$\pm 13.8$	$\pm 18.2$	$\pm 14.2$	$\pm 23.6$	$\pm 17.4$	$\pm 12.9$

Table 16. Performance (DSC, %) comparison on each dataset using mixed training.

View	AbCT-1K	AMOS-CT	AMOS-MRI	BTCV	FLARE22	NIH-Pan	TotalSeg	Urogram	WORD
Axial	93.5	85.1	80.9	76.5	90.8	84.3	84.3	93.3	83.2
	$\pm 1.9$	$\pm 6.5$	$\pm 8.9$	$\pm 5.7$	$\pm 1.0$	$\pm 5.2$	$\pm 10.0$	$\pm 3.0$	$\pm 4.2$
Sagittal	93.3	85.1	80.7	77.1	90.8	84.2	84.7	93.1	83.5
	$\pm 1.8$	$\pm 7.4$	$\pm 8.5$	$\pm 7.2$	$\pm 0.9$	$\pm 4.4$	$\pm 10.6$	$\pm 3.1$	$\pm 4.0$
Coronal	93.3	85.1	80.8	77.6	90.5	84.5	84.1	93.3	83.6
	$\pm 1.7$	$\pm 7.2$	$\pm 9.3$	$\pm 6.6$	$\pm 1.1$	$\pm 4.8$	$\pm 12.7$	$\pm 2.7$	$\pm 4.1$

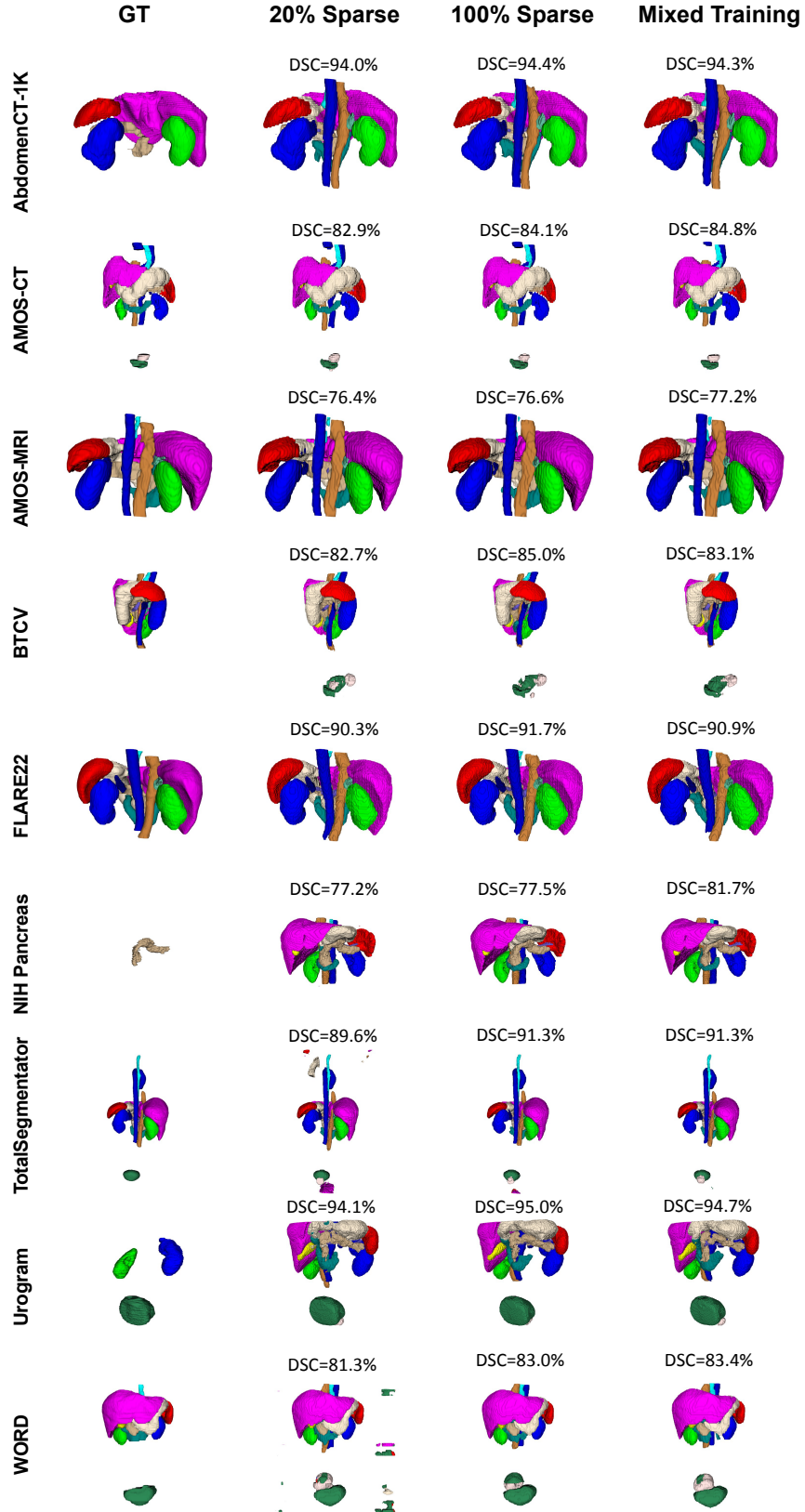


Figure 5. Visual comparisons between the ground truth and predictions from models trained with 20% slices of the axial view, 100% slices of the axial view (loss is computed slice-wise to emulate sparsely labeled data), and hybrid data (the entirety of AMOS, BTCV, and FLARE22 is utilized, while 20% slices of the axial view are taken from other datasets for training) on subjects from various datasets. For a clearer view of detailed differences, zoom in to closely examine the results.