



JOHNS HOPKINS  
UNIVERSITY



ISTITUTO ITALIANO  
DI TECNOLOGIA



RSNA® 2024  
*Building Intelligent Connections*  
Technical Exhibits: Dec. 1-4

# **Touchstone Benchmark:** **Are we on the right way for evaluating medical segmentation?**

**Pedro R. A. S. Bassi**

Johns Hopkins University

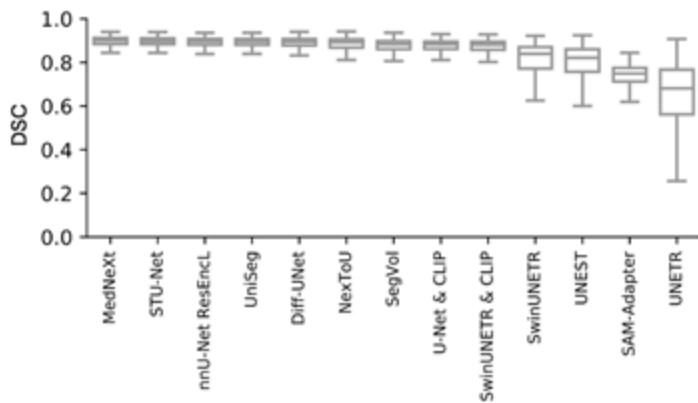
University of Bologna

Italian Institute of Technology

E: [psalvad2@jh.edu](mailto:psalvad2@jh.edu)

## Touchstone Ideals for AI evaluation:

- External (OOD) evaluation
- Large test set
- Analysis by age, sex, race, diagnosis, and more
- AI inventors' participation
- Long-term commitment



14

Research  
teams

29

Institutions

8

Countries

6K

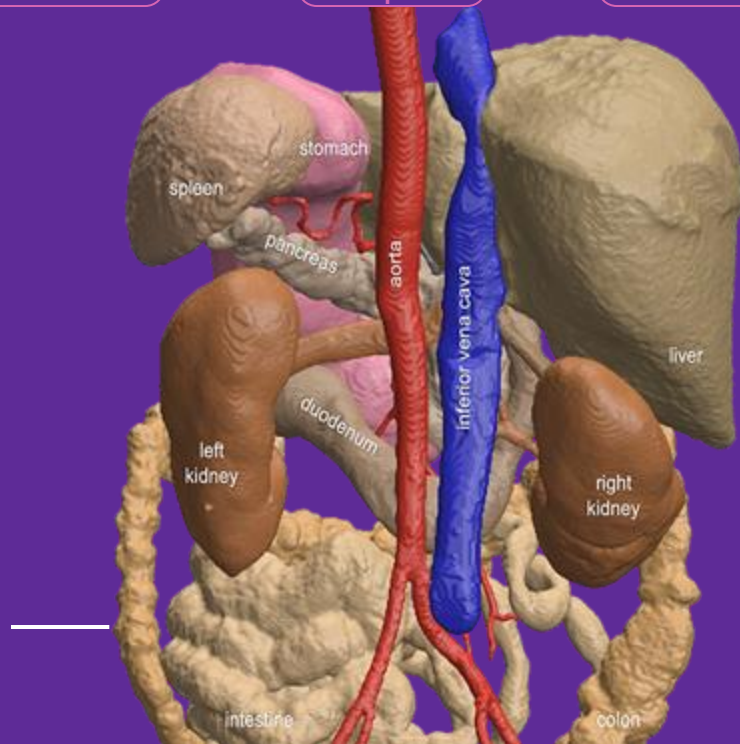
External  
Test CTs

76

Training  
Hospitals

5K

Training  
CTs



# Scale

Table 1: **Related benchmarks & our innovations.** We compare Touchstone with influential CT segmentation benchmarks in light of the five contributions presented in the introduction.

contribution	promoting superior OOD performance with a large and diverse training dataset (#1)			boosting results' significance & large-scale OOD test (#1, #2)	multi-faceted evaluation (#3)	encouraging innovative AI (#4, #5)
benchmark	# CT scans train	# hospitals train	# countries train	# CT scans test	AI consistency analysis	targeted invitation
MSD-CT [2]	947 <sup>†</sup>	1	1	465 IID	none	no
FLARE <sup>22</sup> [53]	2,050 <sup>†</sup>	22	5+	200 IID, 600 OOD	sex, age	no
FLARE <sup>23</sup> [55]	4,000 <sup>†</sup>	30	n/a	n/a	n/a	no
KITS21 [29]	300	50+	1	100 OOD	sex, race	no
AMOS22-CT [38]	200	3	1	78 IID, 122 OOD	none	no
LiTS [9]	130	7	5	70 IID	none	no
BTCV [41]	30	1	1	20 IID	none	no
CHAOS-CT [71]	20	1	1	20 IID	none	no
Touchstone (ours)	<b>5,195</b>	<b>76</b>	<b>8</b>	<b>5,903 OOD</b>	sex, age, race	yes

<sup>†</sup>Partially labeled: annotations for each organ do not cover the entire dataset, and/or may contain unlabeled samples.



## Touchstone Benchmark: Are We on the Right Way for Evaluating AI Algorithms for Medical Segmentation?

Pedro R. A. S. Bassi<sup>1,2,3\*</sup> Wenxuan Li<sup>1\*</sup> Yucheng Tang<sup>4</sup> Fabian Isensee<sup>5,6</sup>  
Zifu Wang<sup>7</sup> Jieneng Chen<sup>1</sup> Yu-Cheng Chou<sup>1</sup> Saikat Roy<sup>5,8</sup> Yannick Kirchhoff<sup>5,8,9</sup>  
Maximilian Rokuss<sup>5,8</sup> Ziyang Huang<sup>10</sup> Jin Ye<sup>11</sup> Junjun He<sup>11</sup> Tassilo Wald<sup>5,6</sup>  
Constantin Ulrich<sup>5</sup> Michael Baumgartner<sup>5,6</sup> Klaus H. Maier-Hein<sup>5,12</sup> Paul Jaeger<sup>6,13</sup>  
Yiwen Ye<sup>14</sup> Yutong Xie<sup>15</sup> Jianpeng Zhang<sup>16</sup> Ziyang Chen<sup>14</sup> Yong Xia<sup>14</sup>  
Zhaohu Xing<sup>17</sup> Lei Zhu<sup>17,18</sup> Yousef Sadegheih<sup>19</sup> Afshin Bozorgpour<sup>19</sup>  
Pratibha Kumari<sup>19</sup> Reza Azad<sup>20</sup> Dorit Merhof<sup>19,21</sup> Pengcheng Shi<sup>22</sup>  
Ting Ma<sup>22</sup> Yuxin Du<sup>23</sup> Fan Bai<sup>23,24</sup> Tiejun Huang<sup>23,25</sup> Bo Zhao<sup>10,23</sup>  
Haonan Wang<sup>18</sup> Xiaomeng Li<sup>18</sup> Hanxue Gu<sup>26</sup> Haoyu Dong<sup>26</sup>  
Jichen Yang<sup>26</sup> Maciej A. Mazurowski<sup>26</sup> Saumya Gupta<sup>27</sup> Linshan Wu<sup>18</sup>  
Jiixin Zhuang<sup>18</sup> Hao Chen<sup>28</sup> Holger Roth<sup>4</sup> Daguang Xu<sup>4</sup>  
Matthew B. Blaschko<sup>7</sup> Sergio Decherchi<sup>29</sup> Andrea Cavalli<sup>2,29,30</sup>  
Alan L. Yuille<sup>1†</sup> Zongwei Zhou<sup>1†</sup>

<sup>1</sup>Department of Computer Science, Johns Hopkins University

<sup>2</sup>Department of Pharmacy and Biotechnology, University of Bologna

<sup>3</sup>Center for Biomolecular Nanotechnologies, Istituto Italiano di Tecnologia

<sup>4</sup>NVIDIA

<sup>5</sup>Division of Medical Image Computing, German Cancer Research Center (DKFZ)

<sup>6</sup>Helmholtz Imaging, German Cancer Research Center (DKFZ)

Full affiliations are given in Appendix F.

Code, Models & Data: <https://github.com/MrGiovanni/Touchstone>

### Abstract

*How can we test AI performance?* This question seems trivial, but it isn't. Standard benchmarks often have problems such as in-distribution and small-size test sets, oversimplified metrics, unfair comparisons, and short-term outcome pressure. As a consequence, good performance on standard benchmarks does not guarantee success in real-world scenarios. To address these problems, we present Touchstone,

# Results

model	organization	average DSC	paper
MedNeXt	DKFZ	89.2	<a href="#">arXiv 2303.09975</a>
STU-Net-B	Shanghai AI Lab	89.0	<a href="#">arXiv 2304.06716</a>
MedFormer	Rutgers	89.0	<a href="#">arXiv 2203.00131</a>
nnU-Net ResEnCL	DKFZ	88.8	<a href="#">arXiv 1809.10486</a>
UniSeg	NPU	88.8	<a href="#">arXiv 2304.03493</a>
Diff-UNet	HKUST	88.5	<a href="#">arXiv 2303.10326</a>
LHU-Net	UR	88.0	<a href="#">arXiv 2404.05102</a>
NexToU	HIT	87.8	<a href="#">arXiv 2305.15911</a>
SegVol	BAAI	87.1	<a href="#">arXiv 2311.13385</a>
U-Net & CLIP	CityU	87.1	<a href="#">arXiv 2301.00785</a>
Swin UNETR & CLIP	CityU	86.7	<a href="#">arXiv 2301.00785</a>
Swin UNETR	NVIDIA	80.1	<a href="#">arXiv 2211.11537</a>
UNesT	NVIDIA	79.1	<a href="#">arXiv 2303.10745</a>
SAM-Adapter	Duke	73.4	<a href="#">arXiv 2404.09957</a>
UNETR	NVIDIA	64.4	<a href="#">arXiv 2111.04004</a>

Table 2: External validation on proprietary JHH dataset ( $N=5,160$ ). Performance is given as DSC score (mean $\pm$ s.d.). For each class, we bold the best-performing results and highlight the runners-up, which show no significant difference from the best results at  $p = 0.05$  level, in red. Architectures are grouped by their frameworks and sorted in ascending order based on the number of parameters. CNNs based on the nnU-Net framework have the best performance on most classes, but other models excel at specific structures (e.g., the graph neural network-based NeXTtoU for aorta, and the diffusion-based Diff-UNet for kidneys). The NSD results are reported in Appendix Table 9.

framework	architecture	param	spleen	kidneyR	kidneyL	gallbladder	liver
nnU-Net	UniSeg <sup>1</sup> [83]	31.0M	94.9 $\pm$ 6.0	92.2 $\pm$ 7.2	91.5 $\pm$ 7.0	84.7 $\pm$ 12.6	96.1 $\pm$ 4.4
	MedNeXt [64]	61.8M	95.2 $\pm$ 6.3	92.6 $\pm$ 7.4	91.8 $\pm$ 7.3	85.3 $\pm$ 12.9	96.3 $\pm$ 4.5
	NexToU [66]	81.9M	94.7 $\pm$ 8.1	90.1 $\pm$ 9.5	89.6 $\pm$ 9.3	82.3 $\pm$ 17.0	95.7 $\pm$ 5.5
	STU-Net-B [34]	58.3M	95.1 $\pm$ 6.4	92.5 $\pm$ 7.3	91.9 $\pm$ 7.2	85.5 $\pm$ 12.3	96.2 $\pm$ 4.8
	STU-Net-L [34]	440.3M	95.2 $\pm$ 6.1	92.5 $\pm$ 7.1	91.8 $\pm$ 7.1	85.7 $\pm$ 11.8	96.3 $\pm$ 4.4
	STU-Net-H [34]	1457.3M	95.2 $\pm$ 5.9	92.6 $\pm$ 6.9	91.9 $\pm$ 7.1	<b>86.0<math>\pm</math>11.6</b>	96.3 $\pm$ 4.4
	U-Net [62]	31.1M	95.1 $\pm$ 6.3	92.7 $\pm$ 6.9	91.9 $\pm$ 7.2	84.7 $\pm$ 13.1	96.2 $\pm$ 4.5
	ResEnCL [35, 37]	102.0M	95.2 $\pm$ 6.3	92.6 $\pm$ 7.0	91.9 $\pm$ 6.9	84.9 $\pm$ 13.0	96.3 $\pm$ 4.5
	ResEnCL *	102.0M	95.1 $\pm$ 6.2	92.7 $\pm$ 6.9	91.9 $\pm$ 7.1	84.9 $\pm$ 13.0	96.3 $\pm$ 4.5
	Vision-Language	U-Net & CLIP [46]	19.1M	94.3 $\pm$ 6.9	91.9 $\pm$ 7.8	91.1 $\pm$ 8.8	82.1 $\pm$ 15.4
	Swin UNETR & CLIP [46]	62.2M	94.1 $\pm$ 7.7	91.7 $\pm$ 9.1	91.0 $\pm$ 9.1	80.2 $\pm$ 18.3	95.8 $\pm$ 5.6
MONAI	LJU-Net [65]	8.6M	94.9 $\pm$ 6.3	92.5 $\pm$ 7.0	91.8 $\pm$ 7.4	83.9 $\pm$ 14.5	96.2 $\pm$ 4.3
	UCTransNet [72]	68.0M	90.2 $\pm$ 11.9	86.5 $\pm$ 14.6	86.9 $\pm$ 12.8	77.8 $\pm$ 19.5	93.6 $\pm$ 6.4
	Swin UNETR [68]	72.8M	92.7 $\pm$ 8.8	89.8 $\pm$ 11.1	89.7 $\pm$ 10.2	76.9 $\pm$ 20.7	95.2 $\pm$ 5.3
	UNesT [85]	87.2M	93.2 $\pm$ 7.1	90.9 $\pm$ 8.1	90.1 $\pm$ 8.2	75.1 $\pm$ 21.2	95.3 $\pm$ 5.0
	UNETR [25]	101.8M	91.7 $\pm$ 10.1	90.1 $\pm$ 9.4	89.2 $\pm$ 9.6	74.7 $\pm$ 20.4	95.0 $\pm$ 5.3
	SegVol <sup>2</sup> [18]	181.0M	94.5 $\pm$ 6.9	92.5 $\pm$ 7.1	91.8 $\pm$ 7.3	79.3 $\pm$ 18.8	96.0 $\pm$ 4.7
n/a	SAM-Adapter <sup>1</sup> [23]	11.6M	90.5 $\pm$ 8.8	90.4 $\pm$ 7.9	87.3 $\pm$ 9.6	49.4 $\pm$ 22.9	94.1 $\pm$ 5.3
	MedFormer [19]	38.5M	<b>95.5<math>\pm</math>6.1</b>	<b>92.8<math>\pm</math>7.3</b>	91.9 $\pm$ 7.4	85.3 $\pm$ 13.6	<b>96.4<math>\pm</math>4.4</b>
	Diff-UNet [81]	434.0M	95.0 $\pm$ 6.9	92.8 $\pm$ 7.4	<b>91.9<math>\pm</math>7.5</b>	83.8 $\pm$ 14.8	96.2 $\pm$ 4.7
framework	architecture	param	stomach	aorta	postcava	pancreas	average
nnU-Net	UniSeg <sup>1</sup> [83]	31.0M	93.3 $\pm$ 6.0	82.3 $\pm$ 10.3	81.2 $\pm$ 8.1	82.7 $\pm$ 10.4	88.8 $\pm$ 5.0
	MedNeXt [64]	61.8M	93.5 $\pm$ 6.0	83.1 $\pm$ 10.2	81.3 $\pm$ 8.3	83.3 $\pm$ 11.0	<b>89.2<math>\pm</math>5.1</b>
	NexToU [66]	81.9M	92.7 $\pm$ 7.5	86.4 $\pm$ 8.7	78.1 $\pm$ 9.1	80.2 $\pm$ 13.5	87.8 $\pm$ 6.2
	STU-Net-B [34]	58.3M	93.5 $\pm$ 6.0	82.1 $\pm$ 10.5	<b>81.3<math>\pm</math>8.2</b>	83.2 $\pm$ 10.7	89.1 $\pm$ 5.3
	STU-Net-L [34]	440.3M	93.7 $\pm$ 5.6	81.0 $\pm$ 10.9	81.3 $\pm$ 8.2	83.4 $\pm$ 10.7	89.0 $\pm$ 5.0
	STU-Net-H [34]	1457.3M	<b>93.7<math>\pm</math>5.7</b>	81.1 $\pm$ 10.9	81.1 $\pm$ 8.2	<b>83.4<math>\pm</math>10.7</b>	89.1 $\pm$ 5.0
	U-Net [62]	31.1M	93.3 $\pm$ 6.0	82.8 $\pm$ 10.2	81.0 $\pm$ 8.2	82.3 $\pm$ 11.4	88.9 $\pm$ 5.1
	ResEnCL [35, 37]	102.0M	93.4 $\pm$ 6.0	81.4 $\pm$ 11.1	80.5 $\pm$ 8.8	82.9 $\pm$ 10.8	88.8 $\pm$ 5.1
	ResEnCL *	102.0M	93.5 $\pm$ 5.9	88.0 $\pm$ 7.3	80.5 $\pm$ 8.7	82.8 $\pm$ 11.1	89.5 $\pm$ 7.8
	Vision-Language	U-Net & CLIP [46]	19.1M	92.4 $\pm$ 6.8	77.1 $\pm$ 12.7	78.5 $\pm$ 9.6	80.8 $\pm$ 11.5
	Swin UNETR & CLIP [46]	62.2M	92.2 $\pm$ 8.3	78.1 $\pm$ 12.6	76.8 $\pm$ 11.0	80.2 $\pm$ 12.5	86.7 $\pm$ 6.3
MONAI	LJU-Net [65]	8.6M	93.0 $\pm$ 6.1	79.5 $\pm$ 11.2	79.4 $\pm$ 9.3	81.0 $\pm$ 11.3	88.1 $\pm$ 5.2
	UCTransNet [72]	68.0M	81.9 $\pm$ 12.9	<b>86.5<math>\pm</math>8.0</b>	68.1 $\pm$ 15.8	59.0 $\pm$ 21.6	81.2 $\pm$ 8.6
	Swin UNETR [68]	72.8M	90.5 $\pm$ 8.6	77.2 $\pm$ 15.1	75.4 $\pm$ 11.8	75.6 $\pm$ 14.5	84.9 $\pm$ 7.1
	UNesT [85]	87.2M	90.9 $\pm$ 7.3	77.7 $\pm$ 16.1	74.4 $\pm$ 11.8	76.2 $\pm$ 12.1	85.0 $\pm$ 6.2
	UNETR [25]	101.8M	88.8 $\pm$ 8.4	76.5 $\pm$ 16.4	71.5 $\pm$ 12.8	72.3 $\pm$ 14.5	83.4 $\pm$ 7.0
	SegVol <sup>2</sup> [18]	181.0M	92.5 $\pm$ 7.0	80.2 $\pm$ 11.3	77.8 $\pm$ 9.7	79.1 $\pm$ 12.4	87.2 $\pm$ 5.6
n/a	SAM-Adapter <sup>1</sup> [23]	11.6M	88.0 $\pm$ 9.3	62.8 $\pm$ 12.2	48.0 $\pm$ 14.2	50.2 $\pm$ 12.6	73.8 $\pm$ 6.3
	MedFormer [19]	38.5M	93.4 $\pm$ 6.4	82.1 $\pm$ 11.7	80.7 $\pm$ 10.1	83.1 $\pm$ 11.2	89.0 $\pm$ 5.4
	Diff-UNet [81]	434.0M	93.1 $\pm$ 6.5	81.2 $\pm$ 11.3	80.8 $\pm$ 8.9	81.9 $\pm$ 11.4	88.6 $\pm$ 5.5

<sup>1</sup>These architectures were pre-trained (Appendix B.3).

\*These architectures were trained on AbdomenAtlas 1.0 with enhanced label quality for the aorta and kidney classes (discussed in §4).

# Results

Table 3: **Validation on TotalSegmentator** ( $N=743$ ). Performances given as DSC score (mean $\pm$ s.d.). For each class, we bold the best-performing results and highlight the runners-up, which show no significant difference from the best results at  $p = 0.05$  level, in red. To ease the direct comparison with other literature, we also reported the *official* test set performance in Appendix Tables 11–12.

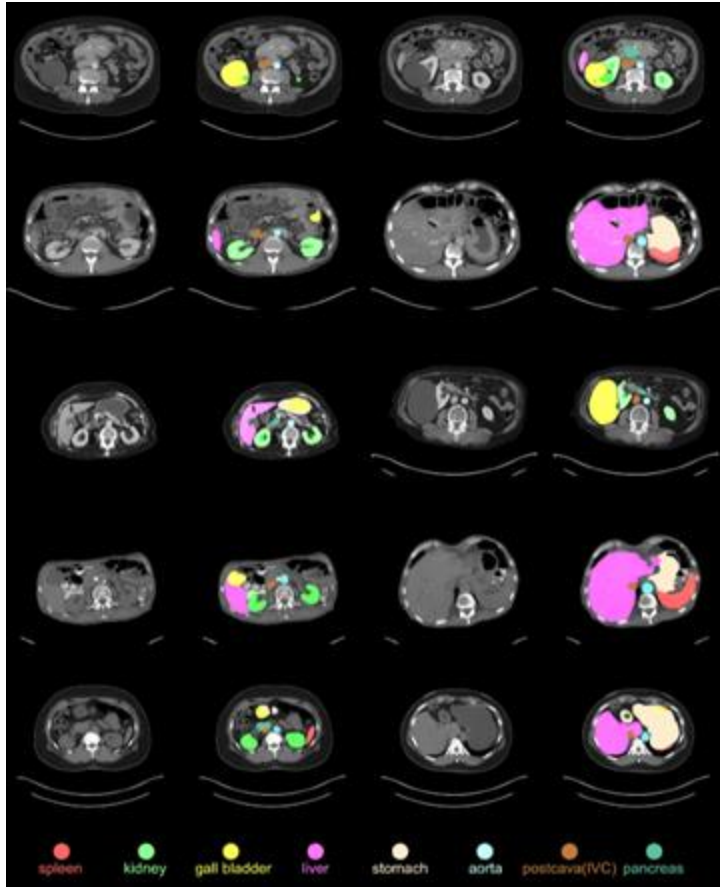
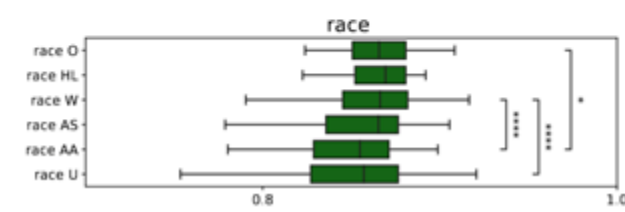
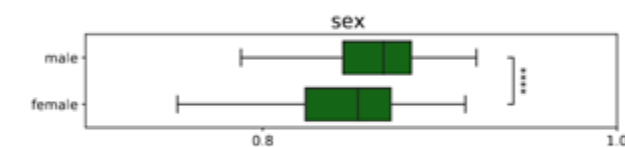
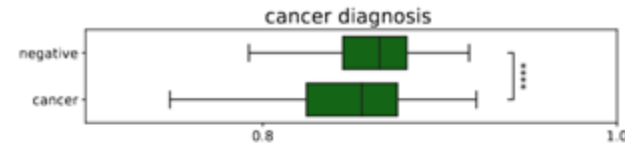
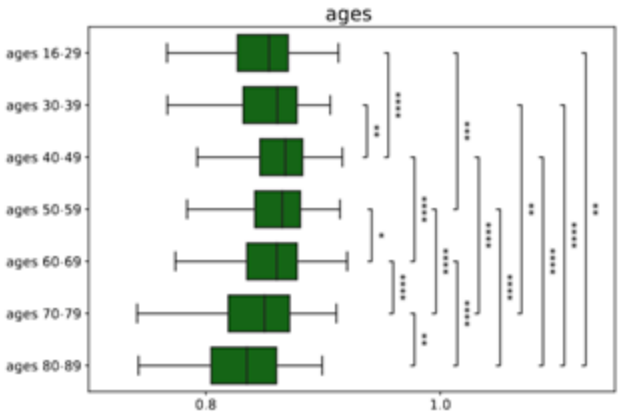
framework	architecture	param	spleen	kidneyR	kidneyL	gallbladder	liver
nnU-Net	UniSeg <sup>1</sup> [83]	31.0M	89.4 $\pm$ 19.4	84.5 $\pm$ 23.8	81.9 $\pm$ 27.9	74.6 $\pm$ 27.3	91.7 $\pm$ 16.5
	ModNeXt [64]	61.8M	91.6 $\pm$ 18.2	85.5 $\pm$ 24.7	86.0 $\pm$ 23.8	75.8 $\pm$ 28.4	93.0 $\pm$ 15.8
	NexToU [66]	81.9M	83.0 $\pm$ 29.5	78.2 $\pm$ 32.7	78.7 $\pm$ 30.8	72.0 $\pm$ 31.1	87.6 $\pm$ 23.0
	STU-Net-B [34]	58.3M	92.3 $\pm$ 15.3	87.1 $\pm$ 20.2	86.8 $\pm$ 22.1	<b>78.5<math>\pm</math>24.9</b>	93.0 $\pm$ 13.9
	STU-Net-L [34]	440.3M	91.6 $\pm$ 17.8	88.2 $\pm$ 18.5	86.3 $\pm$ 22.9	78.1 $\pm$ 24.6	<b>94.2<math>\pm</math>11.2</b>
	STU-Net-H [34]	1457.3M	<b>92.4<math>\pm</math>14.6</b>	<b>88.9<math>\pm</math>16.2</b>	86.5 $\pm$ 23.4	77.7 $\pm$ 25.3	94.0 $\pm$ 11.4
	U-Net [62]	31.1M	91.2 $\pm$ 17.8	88.4 $\pm$ 18.3	87.7 $\pm$ 20.8	78.3 $\pm$ 25.5	93.4 $\pm$ 13.8
	ResEncl. [35, 37]	102.0M	91.8 $\pm$ 17.5	88.9 $\pm$ 18.0	<b>88.2<math>\pm</math>20.5</b>	78.0 $\pm$ 25.1	91.7 $\pm$ 18.4
	ResEncl.*		92.0 $\pm$ 16.7	89.9 $\pm$ 15.3	89.5 $\pm$ 18.3	78.0 $\pm$ 24.7	92.4 $\pm$ 17.4
Vision-Language	U-Net & CLIP [46]	19.1M	87.4 $\pm$ 23.8	83.6 $\pm$ 25.5	82.7 $\pm$ 26.6	73.1 $\pm$ 29.0	91.6 $\pm$ 14.8
	Swin UNETR & CLIP [46]	62.2M	87.1 $\pm$ 22.4	81.1 $\pm$ 28.9	77.0 $\pm$ 32.3	70.3 $\pm$ 30.9	91.6 $\pm$ 16.0
MONAI	LHU-Net [65]	8.6M	86.0 $\pm$ 25.7	81.8 $\pm$ 29.3	82.4 $\pm$ 26.9	71.3 $\pm$ 32.0	87.7 $\pm$ 22.9
	UCTransNet [72]	68.0M	76.4 $\pm$ 34.5	74.3 $\pm$ 35.1	62.0 $\pm$ 41.4	69.6 $\pm$ 31.8	82.6 $\pm$ 28.1
	Swin UNETR [68]	72.8M	66.3 $\pm$ 36.4	59.7 $\pm$ 39.3	58.5 $\pm$ 40.1	50.6 $\pm$ 40.5	80.2 $\pm$ 28.7
	UNesT [85]	87.2M	79.5 $\pm$ 26.6	73.8 $\pm$ 32.3	72.0 $\pm$ 33.8	50.3 $\pm$ 39.9	87.6 $\pm$ 20.8
	UNETR [25]	101.8M	60.4 $\pm$ 37.9	47.9 $\pm$ 39.5	41.9 $\pm$ 39.7	40.0 $\pm$ 36.7	78.1 $\pm$ 29.8
	SegVol <sup>2</sup> [18]	181.0M	87.1 $\pm$ 23.0	82.8 $\pm$ 23.4	82.6 $\pm$ 24.8	68.1 $\pm$ 29.2	89.4 $\pm$ 20.4
n/a	SAM-Adapter <sup>3</sup> [23]	11.6M	53.5 $\pm$ 33.3	8.5 $\pm$ 11.1	19.9 $\pm$ 22.0	11.5 $\pm$ 17.5	66.4 $\pm$ 35.4
	ModFormer [19]	38.5M	90.7 $\pm$ 15.0	85.5 $\pm$ 18.4	84.0 $\pm$ 21.5	74.1 $\pm$ 26.7	92.8 $\pm$ 12.4
	Diff-UNet [81]	434.0M	88.3 $\pm$ 23.5	81.3 $\pm$ 27.9	81.0 $\pm$ 28.3	71.8 $\pm$ 29.9	92.4 $\pm$ 14.8
framework	architecture	param	stomach	aorta	IVC <sup>2</sup>	pancreas	average
nnU-Net	UniSeg <sup>1</sup> [83]	31.0M	74.0 $\pm$ 29.5	69.2 $\pm$ 31.5	72.8 $\pm$ 25.8	70.3 $\pm$ 30.9	71.8 $\pm$ 28.0
	ModNeXt [64]	61.8M	77.2 $\pm$ 28.7	71.9 $\pm$ 30.1	75.2 $\pm$ 23.5	71.6 $\pm$ 31.4	73.9 $\pm$ 27.3
	NexToU [66]	81.9M	69.0 $\pm$ 34.7	61.5 $\pm$ 33.0	59.4 $\pm$ 32.7	66.8 $\pm$ 31.9	61.4 $\pm$ 31.8
	STU-Net-B [34]	58.3M	78.6 $\pm$ 26.5	74.2 $\pm$ 28.9	77.3 $\pm$ 19.5	74.9 $\pm$ 27.4	76.6 $\pm$ 24.9
	STU-Net-L [34]	440.3M	79.7 $\pm$ 24.6	<b>75.7<math>\pm</math>26.9</b>	<b>77.6<math>\pm</math>18.7</b>	75.2 $\pm$ 27.0	<b>78.9<math>\pm</math>21.5</b>
	STU-Net-H [34]	1457.3M	78.5 $\pm$ 25.5	74.7 $\pm$ 28.0	76.9 $\pm$ 19.0	74.5 $\pm$ 27.5	77.6 $\pm$ 23.8
	U-Net [62]	31.1M	78.9 $\pm$ 26.3	71.0 $\pm$ 28.4	76.4 $\pm$ 21.8	75.2 $\pm$ 26.9	74.4 $\pm$ 26.1
	ResEncl. [35, 37]	102.0M	78.9 $\pm$ 25.3	73.8 $\pm$ 25.9	76.4 $\pm$ 20.1	<b>76.3<math>\pm</math>25.8</b>	77.8 $\pm$ 21.8
	ResEncl.*		80.9 $\pm$ 23.0	84.2 $\pm$ 20.5	76.3 $\pm$ 20.0	77.3 $\pm$ 24.9	84.5 $\pm$ 20.1
Vision-Language	U-Net & CLIP [46]	19.1M	77.7 $\pm$ 26.7	59.0 $\pm$ 32.8	65.8 $\pm$ 27.2	74.6 $\pm$ 25.7	67.7 $\pm$ 28.4
	Swin UNETR & CLIP [46]	62.2M	71.2 $\pm$ 30.6	58.6 $\pm$ 34.5	63.6 $\pm$ 27.3	70.3 $\pm$ 28.8	64.6 $\pm$ 30.7
MONAI	LHU-Net [65]	8.6M	71.3 $\pm$ 31.8	63.0 $\pm$ 34.0	67.5 $\pm$ 28.5	68.6 $\pm$ 32.5	65.6 $\pm$ 31.8
	UCTransNet [72]	68.0M	61.6 $\pm$ 36.1	49.7 $\pm$ 34.8	49.3 $\pm$ 36.4	59.0 $\pm$ 35.1	48.5 $\pm$ 34.4
	Swin UNETR [68]	72.8M	52.2 $\pm$ 35.1	54.5 $\pm$ 36.9	38.1 $\pm$ 34.6	42.3 $\pm$ 34.4	45.4 $\pm$ 31.1
	UNesT [85]	87.2M	63.9 $\pm$ 31.4	54.7 $\pm$ 36.9	38.9 $\pm$ 36.2	50.0 $\pm$ 32.9	49.4 $\pm$ 32.3
	UNETR [25]	101.8M	42.1 $\pm$ 32.0	41.0 $\pm$ 31.3	41.3 $\pm$ 32.3	28.2 $\pm$ 29.1	37.3 $\pm$ 27.9
	SegVol <sup>2</sup> [18]	181.0M	71.6 $\pm$ 29.8	60.8 $\pm$ 29.8	63.0 $\pm$ 24.3	66.3 $\pm$ 28.0	66.8 $\pm$ 26.2
n/a	SAM-Adapter <sup>3</sup> [23]	11.6M	48.4 $\pm$ 30.9	15.2 $\pm$ 18.6	4.8 $\pm$ 8.1	30.9 $\pm$ 21.7	23.1 $\pm$ 19.7
	ModFormer [19]	38.5M	<b>80.4<math>\pm</math>23.6</b>	70.3 $\pm$ 28.0	70.0 $\pm$ 24.4	72.5 $\pm$ 27.9	75.1 $\pm$ 24.1
	Diff-UNet [81]	434.0M	73.4 $\pm$ 29.7	61.0 $\pm$ 34.5	60.7 $\pm$ 33.3	69.7 $\pm$ 29.7	62.5 $\pm$ 31.8

<sup>1</sup>These architectures were pre-trained (Appendix B.3).

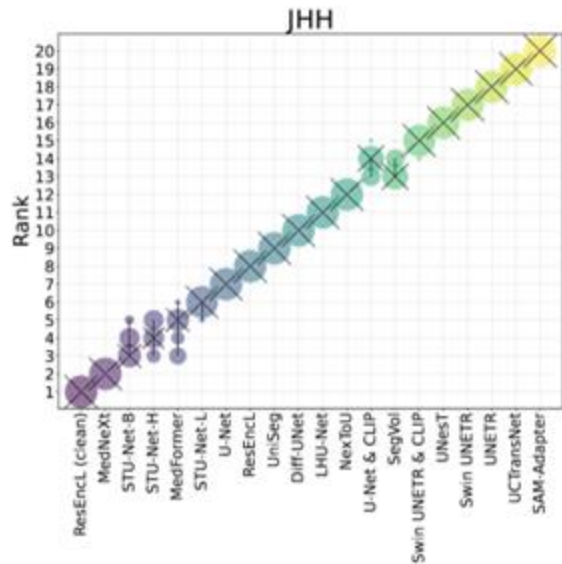
<sup>2</sup>The class IVC (inferior vena cava) shares the same meaning as the class postcava in other datasets (e.g., AbdomenAtlas 1.0 and JHH).

\*These architectures were trained on AbdomenAtlas 1.0 with enhanced label quality for the aorta and kidney classes (discussed in §4).

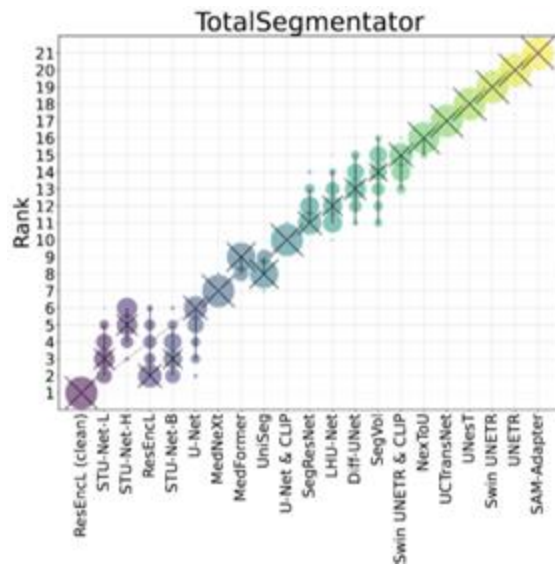
# Potential Confounders Significantly Impact AI



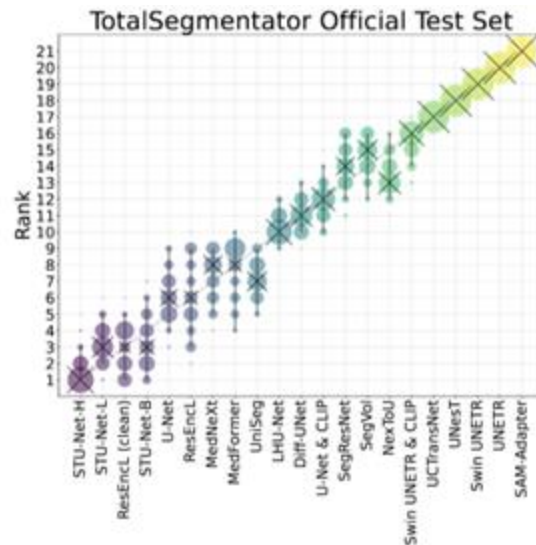
# Test Set Size is Key



N = 5,160



N = 743



N = 59

# Conclusions

1. OOD evaluation: AI performance varies significantly across OOD datasets
2. Large test datasets: more meaningful rankings and nuanced analysis
3. Per-organ analysis revealed AI strengths obscured by mean results
4. Per-group analysis revealed AI biases
5. With creator invitation and third-party evaluation, we establish a fair reference point for future AI algorithms

Touchstone 2.0 is accepting submissions, now with 9K+ CTs



Participate in Touchstone 2.0!





**Thank You!**