

OOD Generalization: ISNet & Touchstone & AbdomenAtlas 3.0

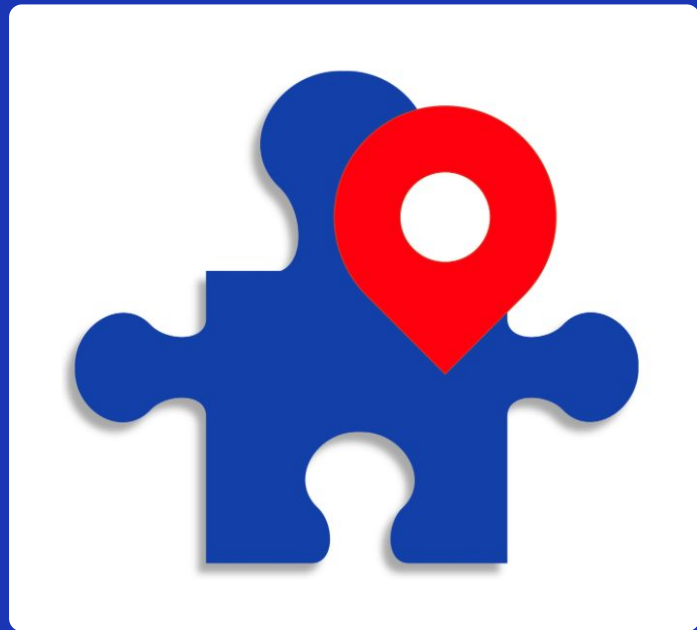
Pedro R. A. S. Bassi
PhD Student University of Bologna, Italy
Visiting PhD Student, CCVL



OOD Generalization: ISNet

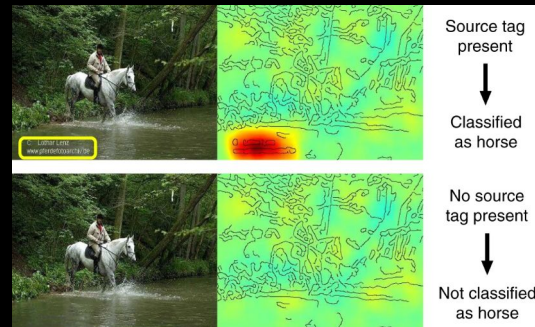
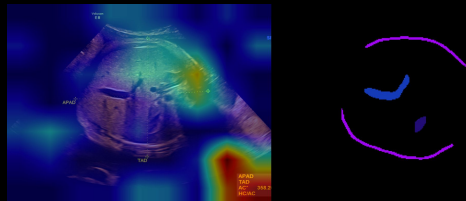
Pedro R. A. S. Bassi

PhD Student University of Bologna, Italy
Visiting PhD Student, CCVL



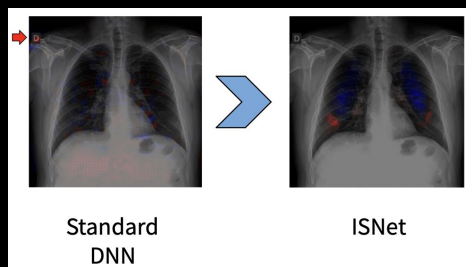
Problem: Bias & Shortcut Learning

- What do DNNs look at?
- Spurious Correlations
- Classification, Segmentation, Medical and Natural
- Mixed Unbalanced Datasets
- IID Vs. OOD



ISNet: A New Training Paradigm

- Objective: Classifier that Ignores Backgrounds
- How: Explanation Heatmap Optimization w/ Foreground Masks
- Segmentation + Classification? Dependency on Segmentation



nature communications



Article

<https://doi.org/10.1038/s41467-023-44371-z>

Improving deep neural network generalization and robustness to background bias via layer-wise relevance propagation optimization

Received: 23 February 2022

Accepted: 11 December 2023

Published online: 04 January 2024

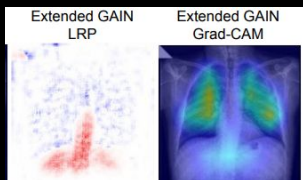
Check for updates

Pedro R. A. S. Bassi^{1,2} , Sergio S. J. Dertkigil³ & Andrea Cavalli^{1,4}

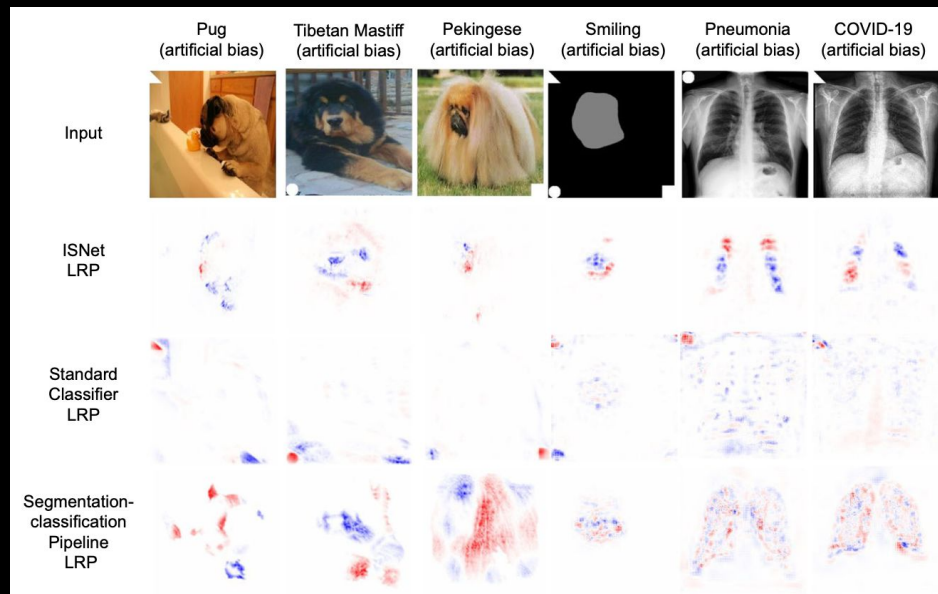
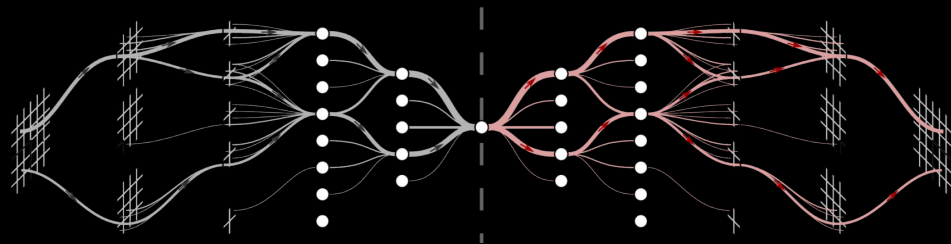
Features in images' backgrounds can spuriously correlate with the images' classes, representing background bias. They can influence the classifier's decisions, causing shortcut learning (Clever Hans effect). The phenomenon generates deep neural networks (DNNs) that perform well on standard evaluation datasets but generalize poorly to real-world data. Layer-wise Relevance Propagation (LRP) explains DNNs' decisions. Here, we show that the optimi-

ISNet: Methodology

- Precise
- Flexible
- Why LRP?
- Technical Challenges:
 - Brackpropagation
 - Loss Design: minimize background, stabilize foreground



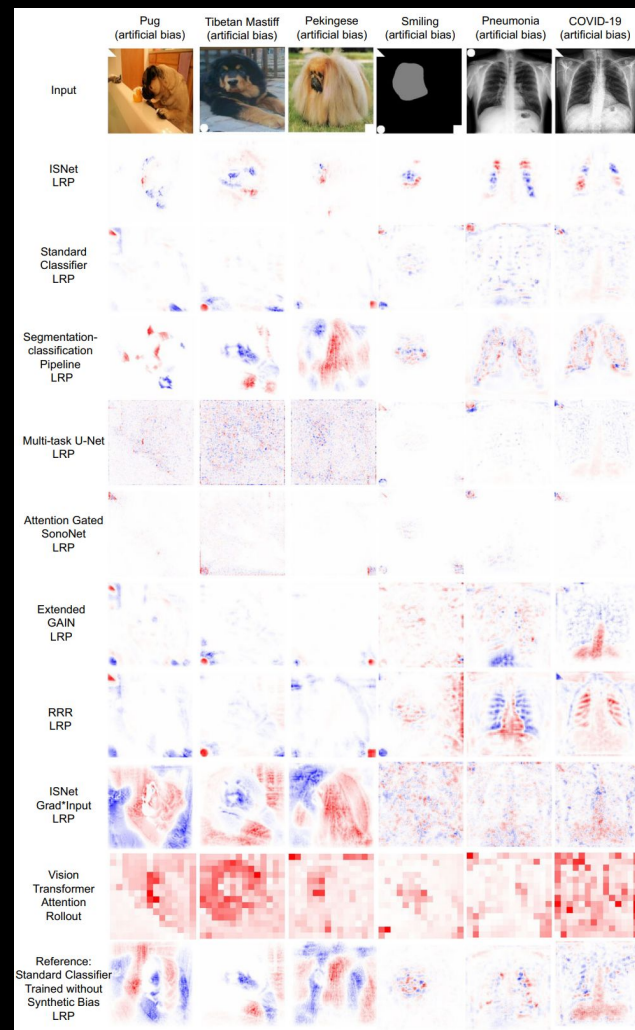
$$R_j = \sum_k \frac{w_{jk} a_j}{z_k + \text{sign}(z_k) \varepsilon} R_k, \text{ where } \varepsilon > 0$$



Synthetic Bias Results

- 8 Baselines, 3 Datasets, Medical & CV
- Not Incremental

Model	Biased test mAF1	Standard test mAF1	Deceiving bias test mAF1
Stanford dogs with synthetic background bias			
ISNet	0.548 ± 0.035	0.553 ± 0.035	0.548 ± 0.035
ISNet Grad*Input	0.55 ± 0.034	0.545 ± 0.034	0.545 ± 0.034
Standard classifier	0.926 ± 0.019	0.419 ± 0.034	0.071 ± 0.017
Segmentation-classification pipeline	0.519 ± 0.035	0.519 ± 0.035	0.518 ± 0.035
COVID-19 detection with synthetic background bias			
ISNet	0.775 ± 0.008	0.775 ± 0.008	0.775 ± 0.008
ISNet Grad*Input	0.542 ± 0.01	0.544 ± 0.01	0.417 ± 0.01
Standard classifier	0.775 ± 0.008	0.434 ± 0.01	0.195 ± 0.004
Segmentation-classification pipeline	0.618 ± 0.009	0.619 ± 0.009	0.618 ± 0.009
Facial attribute estimation with synthetic background bias			
ISNet	0.807 ± 0.027	0.807 ± 0.027	0.807 ± 0.027
ISNet Grad*Input	0.496 ± 0.02	0.499 ± 0.02	0.503 ± 0.021
Standard classifier	0.974 ± 0.012	0.641 ± 0.054	0.398 ± 0.019
Segmentation-classification pipeline	0.794 ± 0.031	0.794 ± 0.031	0.794 ± 0.031



COVID-19 & TB: Real Bias Results

- 8 Baselines, 2 Medical Datasets
- Not Incremental

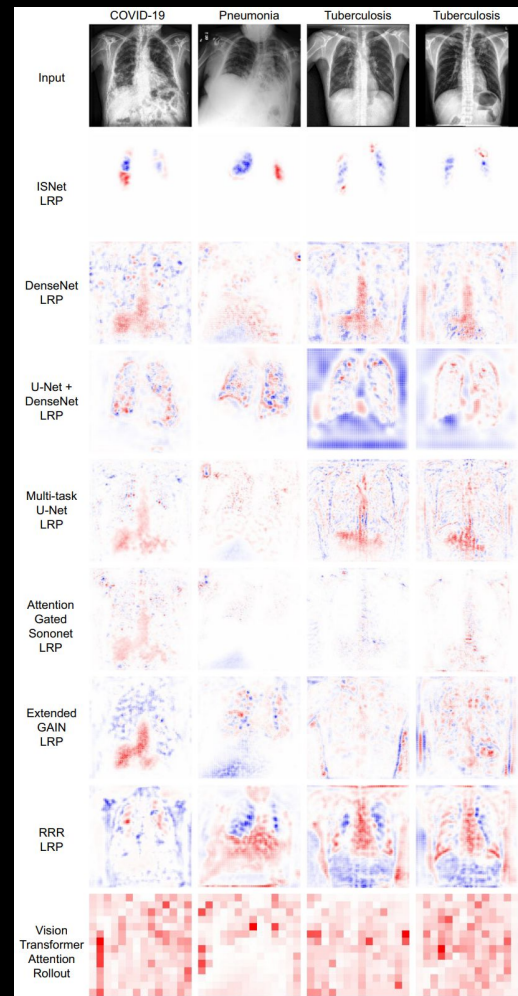
Table 2 | Test F1-Scores and ROC-AUC for the deep neural networks in COVID-19 detection (o.o.d. evaluation)^a

Model and Metric	Normal	Pneumonia	COVID-19	Mean (macro-average)
ISNet F1-Score	0.555 ± 0.022, [0.512,0.597]	0.858 ± 0.007, [0.844,0.871]	0.907 ± 0.006, [0.896,0.918]	0.773 ± 0.009, [0.755,0.791]
U-Net+DenseNet121 F1-Score	0.571 ± 0.018, [0.535,0.607]	0.586 ± 0.013, [0.561,0.611]	0.776 ± 0.008, [0.76,0.792]	0.645 ± 0.009, [0.626,0.663]
DenseNet121 F1-Score	0.444 ± 0.02, [0.403,0.482]	0.434 ± 0.015, [0.405,0.463]	0.76 ± 0.008, [0.744,0.775]	0.546 ± 0.01, [0.527,0.565]
Multi-task U-Net F1-Score	0.419 ± 0.025, [0.369,0.469]	0.119 ± 0.011, [0.098,0.14]	0.585 ± 0.009, [0.566,0.602]	0.374 ± 0.01, [0.355,0.394]
AG-Sononet F1-Score	0.124 ± 0.015, [0.096,0.153]	0.284 ± 0.015, [0.255,0.312]	0.659 ± 0.009, [0.641,0.676]	0.356 ± 0.008, [0.34,0.372]
Extended GAIN F1-Score	0.203 ± 0.019, [0.166,0.24]	0.485 ± 0.013, [0.46,0.511]	0.711 ± 0.009, [0.693,0.728]	0.466 ± 0.009, [0.449,0.485]
RRR F1-Score	0.36 ± 0.018, [0.325,0.394]	0.552 ± 0.013, [0.526,0.577]	0.737 ± 0.009, [0.72,0.755]	0.55 ± 0.009, [0.532,0.568]
Vision Transformer (ViT-B/16) F1-Score	0.382 ± 0.017, [0.348,0.415]	0.474 ± 0.013, [0.448,0.499]	0.525 ± 0.011, [0.503,0.548]	0.46 ± 0.009, [0.443,0.478]

Table 5 | Performance metrics for the deep neural networks in tuberculosis detection (o.o.d. evaluation)^a

Model and metric	Normal	Tuberculosis	Mean (macro-average)
ISNet F1-Score	0.729 ± 0.046	0.748 ± 0.043	0.738 ± 0.044
U-Net+DenseNet121 F1-Score	0.496 ± 0.056	0.656 ± 0.044	0.576 ± 0.05
DenseNet121 F1-Score	0.524 ± 0.052	0.608 ± 0.047	0.566 ± 0.05
Multi-task U-Net F1-Score	0.548 ± 0.049	0.501 ± 0.052	0.524 ± 0.05
AG-Sononet F1-Score	0.522 ± 0.057	0.704 ± 0.04	0.613 ± 0.048
Extended GAIN F1-Score	0.697 ± 0.039	0.501 ± 0.056	0.599 ± 0.048
RRR F1-Score	0.652 ± 0.049	0.674 ± 0.046	0.663 ± 0.048
Vision Transformer (ViT-B/16) F1-Score	0.589 ± 0.046	0.463 ± 0.054	0.526 ± 0.05

- Just our Classifier!
- Segmenter + Classifier



- Just our Classifier!
- Segmenter + Classifier (5x larger, 2x slower)

COVID-19 & TB: Real Bias Results

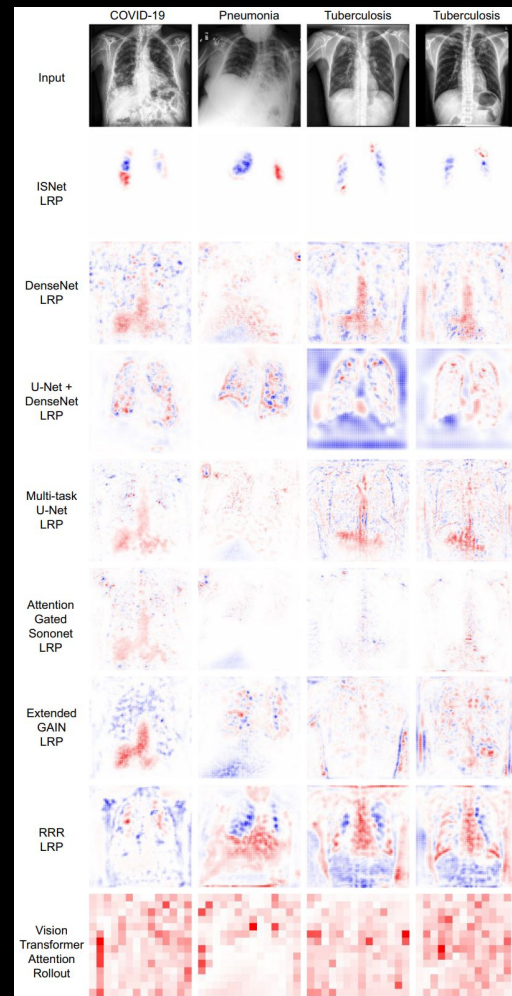
- 8 Baselines, 2 Medical Datasets
- Not Incremental

Table 2 | Test F1-Scores and ROC-AUC for the deep neural networks in COVID-19 detection (o.o.d. evaluation)^a

Model and Metric	Normal	Pneumonia	COVID-19	Mean (macro-average)
ISNet F1-Score	0.555 ± 0.022, [0.512,0.597]	0.858 ± 0.007, [0.844,0.871]	0.907 ± 0.006, [0.896,0.918]	0.773 ± 0.009, [0.755,0.791]
U-Net+DenseNet121 F1-Score	0.571 ± 0.018, [0.535,0.607]	0.586 ± 0.013, [0.561,0.611]	0.776 ± 0.008, [0.76,0.792]	0.645 ± 0.009, [0.626,0.663]
DenseNet121 F1-Score	0.444 ± 0.02, [0.403,0.482]	0.434 ± 0.015, [0.405,0.463]	0.76 ± 0.008, [0.744,0.775]	0.546 ± 0.01, [0.527,0.565]
Multi-task U-Net F1-Score	0.419 ± 0.025, [0.369,0.469]	0.119 ± 0.011, [0.098,0.14]	0.585 ± 0.009, [0.566,0.602]	0.374 ± 0.01, [0.355,0.394]
AG-Sononet F1-Score	0.124 ± 0.015, [0.096,0.153]	0.284 ± 0.015, [0.255,0.312]	0.659 ± 0.009, [0.641,0.676]	0.356 ± 0.008, [0.34,0.372]
Extended GAIN F1-Score	0.203 ± 0.019, [0.166,0.24]	0.485 ± 0.013, [0.46,0.511]	0.711 ± 0.009, [0.693,0.728]	0.466 ± 0.009, [0.449,0.485]
RRR F1-Score	0.36 ± 0.018, [0.325,0.394]	0.552 ± 0.013, [0.526,0.577]	0.737 ± 0.009, [0.72,0.755]	0.55 ± 0.009, [0.532,0.568]
Vision Transformer (ViT-B/16) F1-Score	0.382 ± 0.017, [0.348,0.415]	0.474 ± 0.013, [0.448,0.499]	0.525 ± 0.011, [0.503,0.548]	0.46 ± 0.009, [0.443,0.478]

- Just our Classifier!

- Segmenter + Classifier
(5x larger, 2x slower)



Explanation Distillation

- Unbiased Teacher + Biased Data = Biased Student
- Unbiased Teacher + Biased Data + Expl. Distillation = Unbiased Student
- Any Bias
- No Other Loss

Method	COLOURED MNIST			DogsWithTies		
	IID	OOD	Shift	IID	OOD	Shift
ERM	100	16.93	0.18	92.1	78.2	62.4
LRP Distill	99.2	98.2	98.3	79.2	81.2	77.2
Gradient*Input Distill	98.7	97.4	94.4	70.3	67.3	67.3
Input Gradient Distill	98.2	95.6	87.7	68.3	64.4	72.3
Grad-CAM Distill	N/A	N/A	N/A	70.3	69.3	66.3
Attention Distill	N/A	N/A	N/A	82.2	80.2	78.2
Features Distill	99.9	92.1	85.1	51.5	51.5	53.5
Output Distill	100	92.6	87.3	86.1	80.2	75.2
IRM*	99.7	60.2	53.2	N/A	N/A	N/A
GroupDRO*	99.6	52.2	40.3	N/A	N/A	N/A
PGI*	99.7	63.6	58.2	N/A	N/A	N/A
ISNet**	N/A	N/A	N/A	N/A	N/A	N/A
RRR**	N/A	N/A	N/A	N/A	N/A	N/A
GAIN**	N/A	N/A	N/A	N/A	N/A	N/A
GALS**	N/A	N/A	N/A	75.2	73.3	71.3
Teacher	98.6	98.7	98.7	96	96	96

OOD Generalization: Touchstone

Pedro R. A. S. Bassi

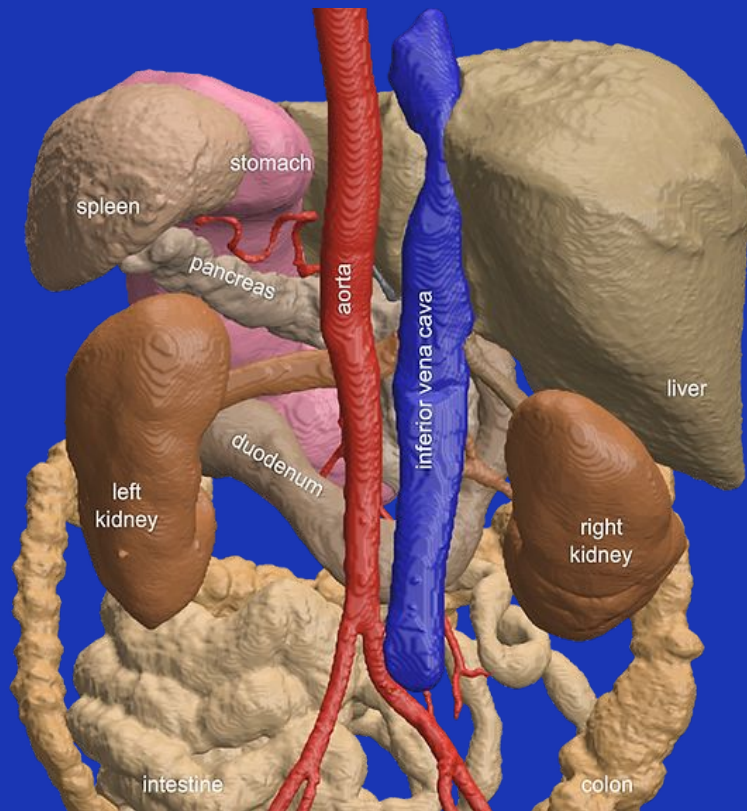
PhD Student University of Bologna, Italy

Visiting PhD Student, CCVL



Pitfalls in Medical Segmentation Evaluation

- 1- IID testing
- 2- Small test set
- 3- Focus on average dice
- 4- Unfair comparisons
- 5- Short-term outcome pressure

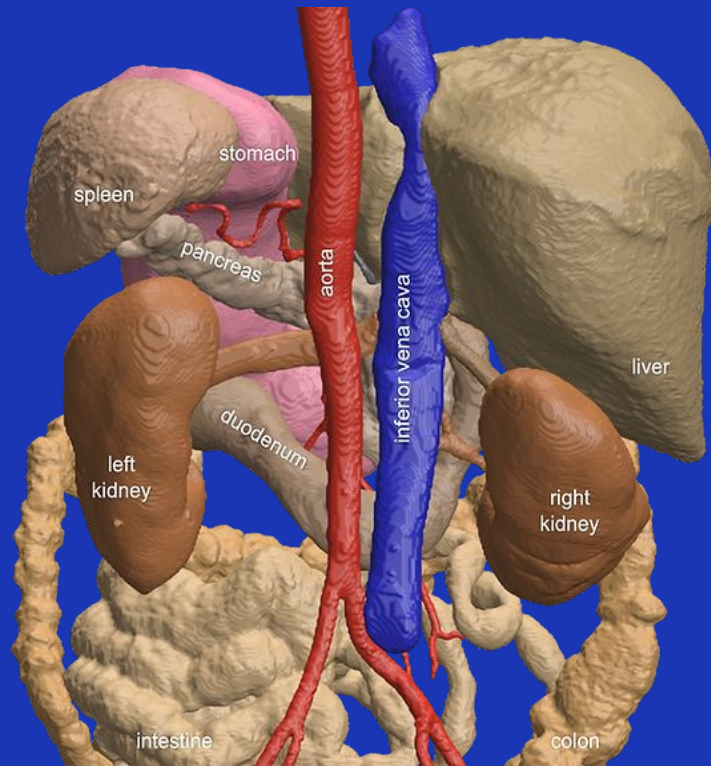


Touchstone Benchmark

Task: 9 Organ Segmentation

Our **ideals** for AI evaluation:

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



Scale & Contributions

14
Research
teams

29
Institutions

8
Countries

6K
External
Test CTs

76
Training
Hospitals

5K
Training
CTs

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 [†]	1	1	465 IID	none	no
FLARE'22 [53]	2,050 [†]	22	5+	200 IID, 600 OOD	sex, age	no
FLARE'23 [55]	4,000 [†]	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)	5,195	76	8	5,903 OOD	sex, age, race	yes

[†]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^{1,2,3*} Wenxuan Li^{1*} Yucheng Tang⁴ Fabian Isensee^{5,6}
 Zifu Wang⁷ Jieneng Chen¹ Yu-Cheng Chou¹ Yannick Kirchhoff^{5,8,9}
 Maximilian Rokuss^{5,8} Ziyang Huang¹⁰ Jin Ye¹¹ Junjun He¹¹ Tassilo Wald^{5,6}
 Constantin Ulrich⁵ Michael Baumgartner^{5,6} Saikat Roy^{5,8}
 Klaus H. Maier-Hein^{5,12} Paul Jaeger^{6,13} Yiwen Ye¹⁴ Yutong Xie¹⁵ Jianpeng Zhang¹⁶
 Ziyang Chen¹⁴ Yong Xia¹⁴ Zhaohu Xing¹⁷ Lei Zhu^{17,18} Yousef Sadegheh¹⁹
 Afshin Bozorgpour¹⁹ Pratibha Kumari¹⁹ Reza Azad²⁰ Dorit Merhof^{19,21}
 Pengcheng Shi²² Ting Ma²² Yuxin Du²³ Fan Bai^{23,24} Tiejun Huang^{23,25}
 Bo Zhao^{10,23} Haonan Wang¹⁸ Xiaomeng Li¹⁸ Hanxue Gu²⁶
 Haoyu Dong²⁶ Jichen Yang²⁶ Maciej A. Mazurowski²⁶ Saumya Gupta²⁷
 Linshan Wu¹⁸ Jiaxin Zhuang¹⁸ Hao Chen²⁸ Holger Roth⁴ Daguang Xu⁴
 Matthew B. Blaschko⁷ Sergio Decherchi²⁹ Andrea Cavalli^{2,29,30}
 Alan L. Yuille^{1†} Zongwei Zhou^{1†}

¹Department of Computer Science, Johns Hopkins University

²Department of Pharmacy and Biotechnology, University of Bologna

³Center for Biomolecular Nanotechnologies, Istituto Italiano di Tecnologia

⁴NVIDIA

⁵Division of Medical Image Computing, German Cancer Research Center (DKFZ)

⁶Helmholtz 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, a large-scale collaborative segmentation benchmark of 9 types of abdominal organs. This benchmark is based on 5,195 training CT scans from 76 hospitals around the world and 5,903 testing CT scans from 11 additional hospitals. This diverse test set enhances the statistical significance of benchmark results and rigorously evaluates AI algorithms across out-of-distribution scenarios. We invited 14 inventors of 19 AI algorithms to train their algorithms, while our team, as a third party, independently evaluated these algorithms. In addition, we also evaluated pre-existing AI frameworks—which, differing from algorithms, are more flexible and can support different algorithms—including MONAI from NVIDIA, nnU-Net from DKFZ, and numerous other open-source frameworks. We are committed to expanding this benchmark to encourage more innovation of AI algorithms for the medical domain.

Results

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

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 [†] [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	86.0\pm11.6	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 12.8	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	96.0 \pm 4.3
	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	LHU-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 [†] [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 [†] [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	95.5\pm6.1	92.8\pm7.3	91.9 \pm 7.4	85.3 \pm 13.6	96.4\pm4.7
	Diff-UNet [81]	434.0M	95.0 \pm 6.9	92.8 \pm 7.4	91.9\pm7.5	83.8 \pm 14.8	96.2 \pm 4.4
framework	architecture	param	stomach	aorta	postcava	pancreas	average
nnU-Net	UniSeg [†] [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	89.2\pm5.1
	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	81.3\pm8.2	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	93.7\pm5.7	81.1 \pm 10.9	81.1 \pm 8.2	83.4\pm10.7	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	87.2 \pm 5.0
	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	LHU-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	86.5\pm8.0	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 [†] [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 [†] [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

[†]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).

Total Segmetator

- Large DSC Change Across OOD Datasets
- Max: 40%
- Usual: 5-15%

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 [†] [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
	MedNeXt [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	78.5\pm24.9	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	94.2\pm11.2
	STU-Net-H [34]	1457.3M	92.4\pm14.6	88.9\pm16.2	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	88.2\pm20.5	78.0 \pm 25.1	91.7 \pm 18.4
Vision-Language	ResEncl.*	102.0M	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
	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
MONAI	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
	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 [‡] [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 [†] [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
	MedFormer [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 [‡]	pancreas	average
nnU-Net	UniSeg [†] [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
	MedNeXt [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	75.7\pm26.9	77.6\pm18.7	75.2 \pm 27.0	78.9\pm21.5
	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	76.3\pm25.8	77.8 \pm 21.8
Vision-Language	ResEncl.*	102.0M	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
	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
MONAI	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
	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 [‡] [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 [†] [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
	MedFormer [19]	38.5M	80.4\pm23.6	70.3 \pm 28.0	70.0 \pm 23.4	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 23.3	69.7 \pm 29.7	62.5 \pm 31.8

[†]These architectures were pre-trained (Appendix B.3).

[‡]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).

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 NeXTou 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 [†] [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 12.0	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 17.3	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	86.0\pm11.6	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 12.8	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	96.0 \pm 4.3
	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	LHU-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	76.9 \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 12.2	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 [†] [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 [†] [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	95.5\pm6.1	92.8\pm7.3	91.9 \pm 7.4	85.3 \pm 13.6	96.4\pm4.4
	Diff-UNet [81]	434.0M	95.0 \pm 6.9	92.8 \pm 7.4	91.9\pm7.5	83.8 \pm 14.8	96.2 \pm 4.7

framework	architecture	param	stomach	aorta	postcava	pancreas	average
nnU-Net	UniSeg [†] [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	89.2\pm5.1
	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	81.3\pm8.2	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	93.7\pm5.7	81.1 \pm 10.9	81.1 \pm 8.2	83.4\pm10.7	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	87.2 \pm 5.0
	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	LHU-Net [65]	8.6M	93.0 \pm 6.1	79.5 \pm 11.2	79.4 \pm 9.3	80.3 \pm 12.3	88.1 \pm 5.2
	UCTransNet [72]	68.0M	81.9 \pm 12.9	86.5\pm8.0	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 [†] [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 [†] [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

[†]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).

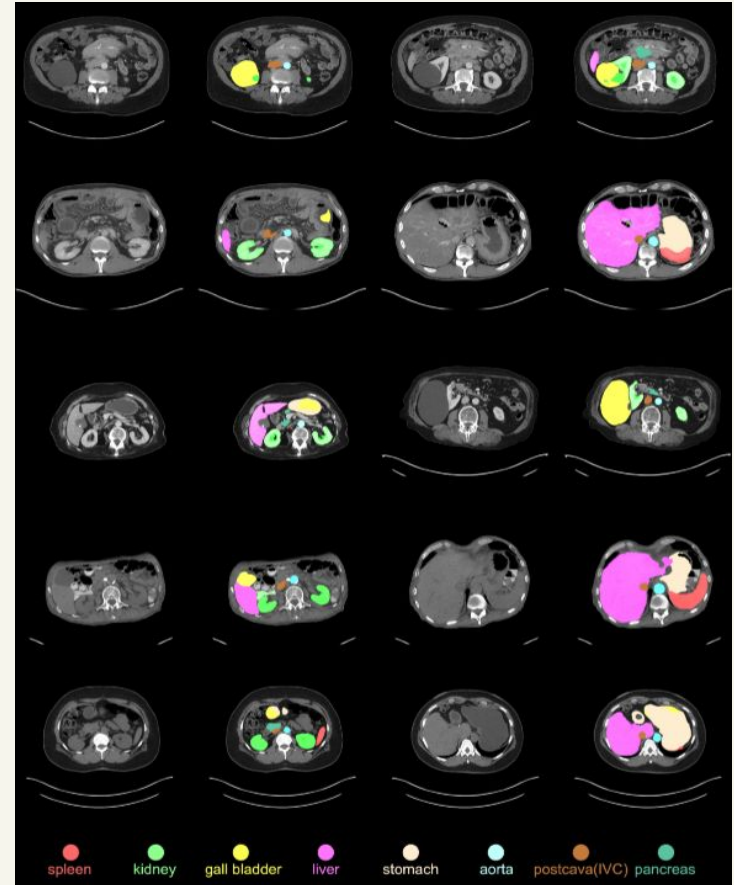
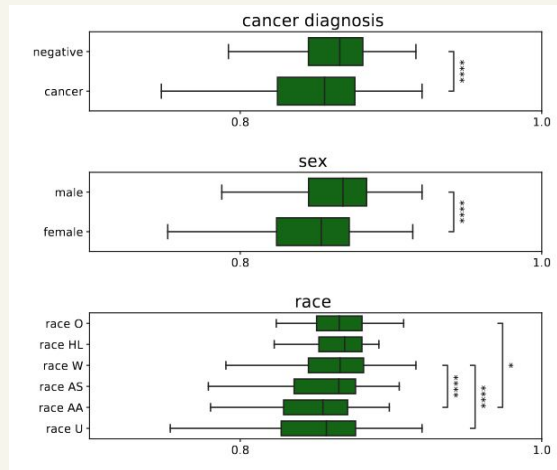
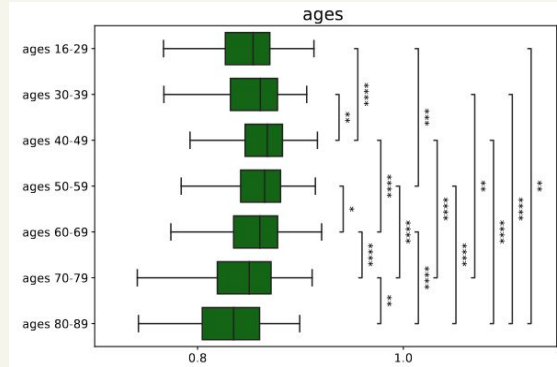
Table 3: **Validation on TotalSegmentator ($N=743$).** Performances given as DSC score (mean \pm s.d.). For each class, we bold the best-performing results and highlight the runners-up, which show no significant difference from the best results at $p = 0.05$ level, in red. To ease the direct comparison 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 [†] [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
	MedNeXt [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	78.5\pm24.9	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	77.1 \pm 24.6	94.2\pm11.2
	STU-Net-H [34]	1457.3M	92.4\pm14.6	88.9 \pm 16.2	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	88.2\pm20.5	78.0 \pm 25.1	91.7 \pm 18.4
	ResEncl.*	102.0M	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 37.7	40.0 \pm 36.7	78.1 \pm 29.8
	SegVol [†] [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 [†] [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
	MedFormer [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

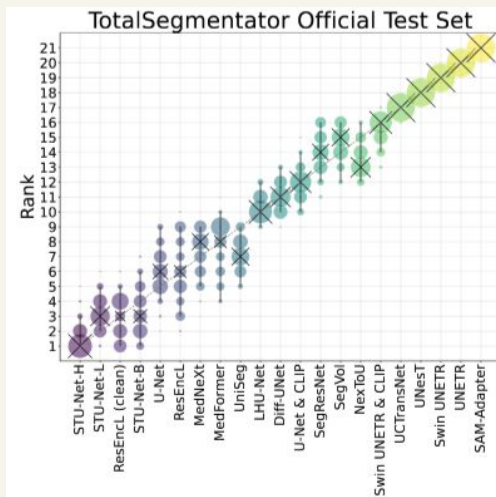
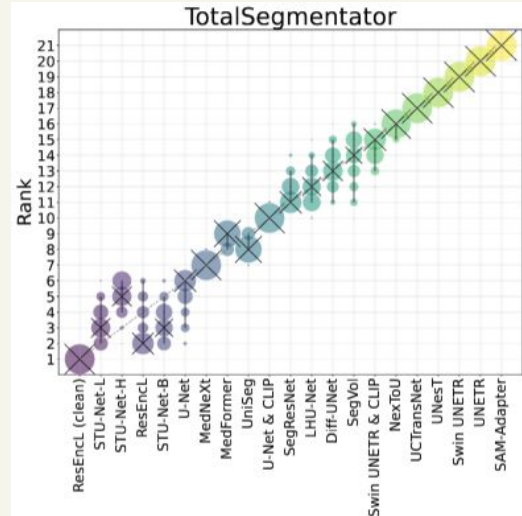
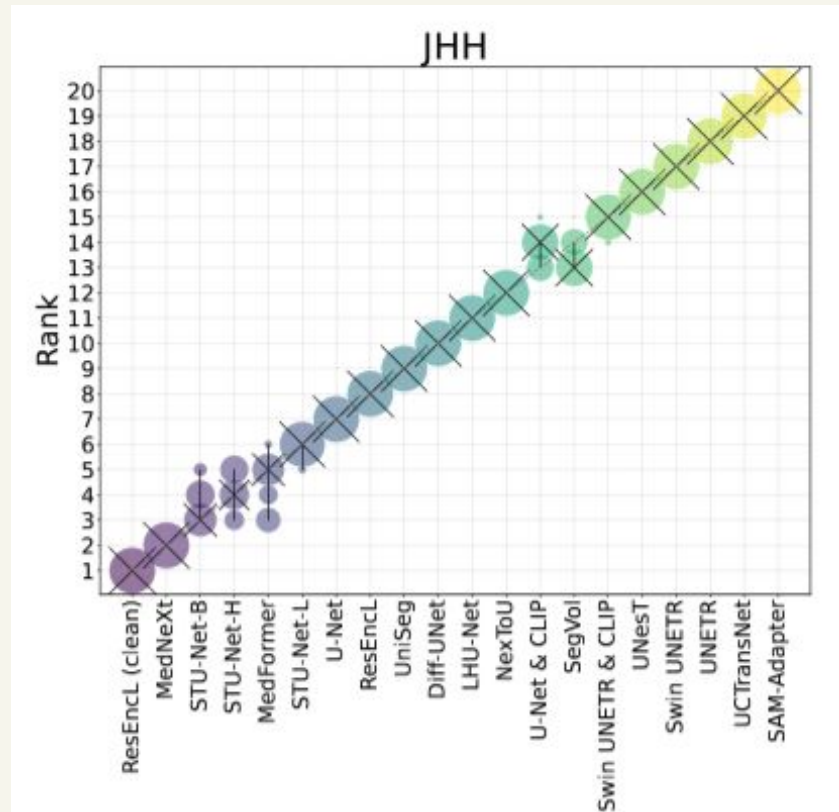
framework	architecture	param	stomach	aorta	IVC [†]	pancreas	average
nnU-Net	UniSeg [†] [83]	31.0M	74.0 \pm 29.5	69.2 \pm 31.5	72.8 \pm 25.8	70.3 \pm 30.9	71.8 \pm 28.0
	MedNeXt [64]	61.8M	77.2 \pm 28.7	71.9 \pm 30.1	75.2 \pm 31.5	71.6 \pm 31.4	73.9 \pm 27.3
	NexToU [66]	81.9M	69.0 \pm 34.7	61.5 \pm 33.0	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	73.7 \pm 19.5	74.9 \pm 27.4	76.6 \pm 24.9
	STU-Net-L [34]	440.3M	79.7 \pm 24.6	75.7\pm26.9	77.6\pm18.7	75.2 \pm 27.0	78.9\pm21.5
	STU-Net-H [34]	1457.3M	78.5 \pm 25.5	74.7 \pm 28.0	76.9 \pm 19.0	74.5 \pm 27.7	78.5 \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	76.3\pm25.8	77.8 \pm 21.8
	ResEncl.*	102.0M	80.9 \pm 23.0	84.2 \pm 20.5	76.3 \pm 20.0	77.3 \pm 24.9	84.3 \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 28.3	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 25.5	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 [†] [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 [†] [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
	MedFormer [19]	38.5M	80.4\pm23.6	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

Biases

- Biases that **Impact** Performance:
 - Diagnosis
 - Sex
 - Age
 - Race



Test Set Size is Key



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

OOD Generalization: AbdomenAtlas 3.0

Pedro R. A. S. Bassi

PhD Student University of Bologna, Italy
Visiting PhD Student, CCVL



AbdomenAtlas 3.0

Liver:

Normal size (volume: 1249.9 cm^3).
Mean HU value: 131.4 ± 23.2 .

Liver lesions:

Liver tumor 1:

Location: hepatic segment 7/8.
Size: $3.1 \times 2.7 \text{ cm}$ (image 387). Volume: 10.9 cm^3 .
Enhancement relative to liver: Hypoattenuating (HU value is 58.7 ± 29.7).

Pancreas:

Pancreas is enlarged (volume: 84.6 cm^3).
Mean HU value: 105.7 ± 33.1 .

Pancreas lesions:

Pancreas tumor 1:

Location: pancreas head/body.
Size: $2.9 \times 2.2 \text{ cm}$ (image 298). Volume: 8.2 cm^3 .
Tumor Stage (T stage): T2.
Enhancement relative to pancreas:
Hypoattenuating (HU value is 52.6 ± 26.8).

Kidney:

Normal size (right kidney volume: 148.6 cm^3 ; left kidney volume: 166.6 cm^3 ; total kidney volume: 315.3 cm^3).
Mean HU value: 172.9 ± 57.4 .

Kidney lesions:

Kidney tumor 1:

Location: left kidney.
Size: $0.6 \times 0.3 \text{ cm}$ (image 321). Volume: 0.1 cm^3 .
Enhancement relative to kidney: Hypoattenuating (HU value is 108.5 ± 49.1).

Kidney tumor 2:

Location: left kidney.
Size: $0.5 \times 0.4 \text{ cm}$ (image 321). Volume: 0.1 cm^3 .
Enhancement relative to kidney: Hypoattenuating (HU value is 97.8 ± 54.5).

Kidney tumor 3:

Location: left kidney.
Size: $0.6 \times 0.4 \text{ cm}$ (image 283). Volume: 0.1 cm^3 .
Enhancement relative to kidney: Hypoattenuating (HU value is 49.5 ± 48.4).

AbdomenReport

- 9,262 CT-Report pairs
- 1.2 million tokens
- 7,960 tumor instances
- 5,262 small tumors (<20mm)
- 3 pancreatic sub-segments
- 8 hepatic sub-segments
- 4 pancreatic tumor stages

Future Directions

- Collaborations for OOD evaluation of our models and benchmark models
- Large, diverse and unbiased datasets for OOD accuracy
- Language is less strong than per-voxel annotations: architectures and training strategies to avoid bias and improve OOD accuracy in medical VLMs