

# Insights from Acquiring Open Medical Imaging Datasets for Foundation Model Development

Stefan Dvoretzki<sup>1,5</sup> // Paul Jäger<sup>1,4</sup> // Fabian Isensee<sup>1,4</sup> // Tassilo Wald<sup>1,4</sup> // Constantin Ulrich<sup>1</sup> // Lucas Kulla<sup>1,5</sup> // Philipp Schader<sup>1,2,5</sup> // Klaus Maier-Hein<sup>1,3</sup> // Josh Moore<sup>6</sup> // Marco Nolden<sup>1,3,5</sup>

<sup>1</sup> Division of Medical Computing, Deutsches Krebsforschungszentrum (DKFZ) Heidelberg German Cancer Research Center (DKFZ) Heidelberg, Division of Medical Image Computing, Germany; <sup>2</sup> Faculty of Mathematics and Computer Science, Heidelberg University, Heidelberg, Germany; <sup>3</sup> Pattern Analysis and Learning Group, Department of Radiation Oncology, Heidelberg University Hospital, Heidelberg, Germany; <sup>4</sup> Helmholtz Imaging, German Cancer Research Center (DKFZ), Heidelberg, Germany; <sup>5</sup> Helmholtz Metadata Collaboration (HMC) Hub Health, German Cancer Research Center (DKFZ), Heidelberg, Germany; <sup>6</sup> German BioImaging, University of Konstanz, Germany

## Background

- Foundation Model for Radiology involves a huge amount of clinical imaging volumes
- We made a list of open clinical imaging datasets – UK Biobank, NAKO, HCP and others
- Some open clinical imaging datasets were not machine-actionable and accessible
- We look into the problems and suggest helpful concepts

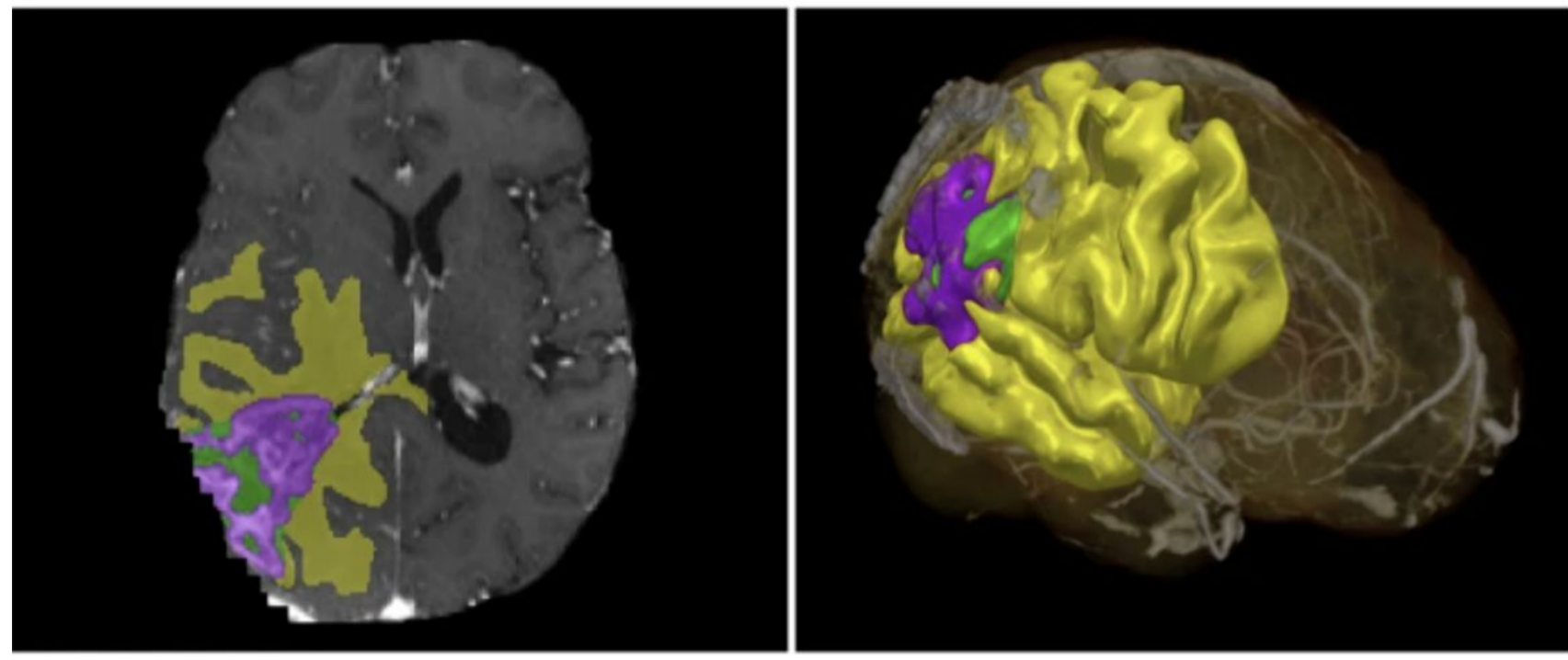
## How to get (meta)data AI-ready?

Data concepts like **FAIR Digital Object** or a **FAIR Data Point** provide basis for machine data access. Some data portals provide good data overview mechanisms. There may even be FAIR-compliant APIs behind.

## Data for Radiological Foundation Model

## Datasets are not machine-actionable

We ended up doing a lot by hand. This was time and effort consuming. A lot of described problems could be seen through the **FAIR** data guidelines lens.



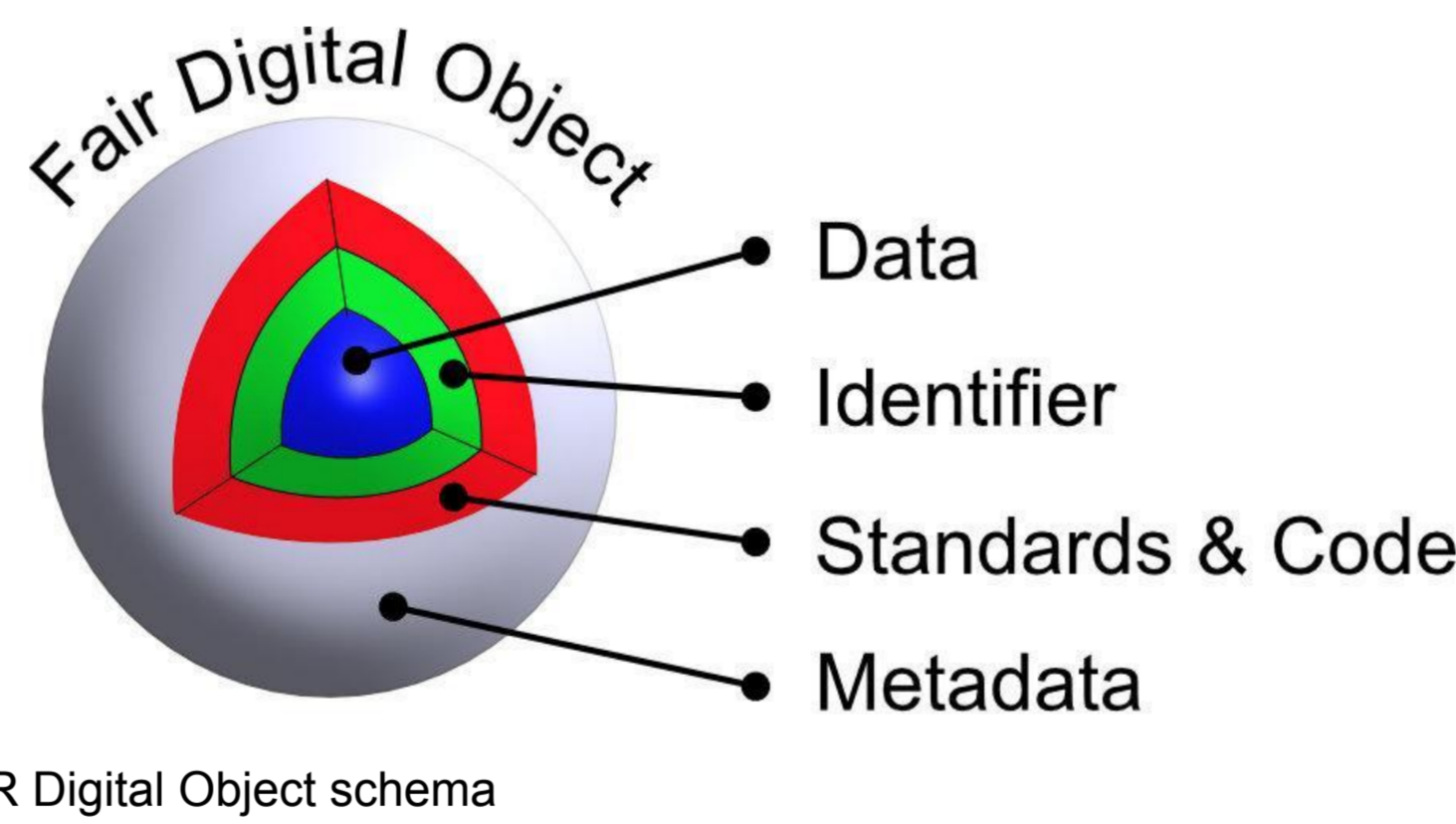
MedSAM CT image datasets (Isensee, Jäger et al. Nat Methods 2021)

Dataset Name	Modality	Segmentation Targets	# of scans
AbdomenCT-1K [1], [2]	CT	Liver, kidneys, pancreas, spleen	1056
Adrenal-ACC-Ki67* [3]-[5]	CT	Adrenocortical carcinoma	53
AMOS-CT [6]	CT	Abdominal organ	200
AutoPET [7]	PET-CT	Whole-body tumor	900
COVID-19 Seg. Challenge [8], [9]	CT	COVID-19 infections	199
COVID-19-CT-Seg [10]	CT	COVID-19 infections, left lung, and right lung	20
GLIS-RT [11]	CT	Head tumor	75
HCC-TACE-Seg* [5], [12]	CT	Liver cancer	70
HECKTOR [13]	PET-CT	Head and neck tumor	524
INSTANCE [14]	CT	Hematoma	100
KiPA [15], [16]	CT	Kidney, tumor, renal artery, renal vein	70
KiTS [17]	CT	Kidney, tumor, cyst	489
LNQ2023* [18]	CT	Mediastinal lymph node	393
Lymph Nodes [19], [20]	CT	Lymph nodes	176
MSD-Colon Tumor [21]	CT	Colon tumor	126
MSD-Hepatic Tumor [21]	CT	Hepatic tumor	303
MSD-Lung Tumor [21]	CT	Lung tumor	96
MSD-Pancreas [21]	CT	Pancreas, pancreas tumor	281
MSD-Spleen [21]	CT	Spleen	61
NSCLC Pleural Effusion [5], [22], [23]	CT	Pleural effusion	78
NSCLC Radiogenomics [24]	CT	Lung Tumor	88
ORG [25]	CT	Whole-body organs	140
SegTHOR [26]	CT	Esophagus, heart, aorta, trachea	40
StructSeg* [27]	CT	Nasopharyngeal cancer and lung cancer, with OAR and GTV	50
TotalSegmentor [28]	CT	Whole body organs	1204
WORD* [29]	CT	Abdominal organs	150

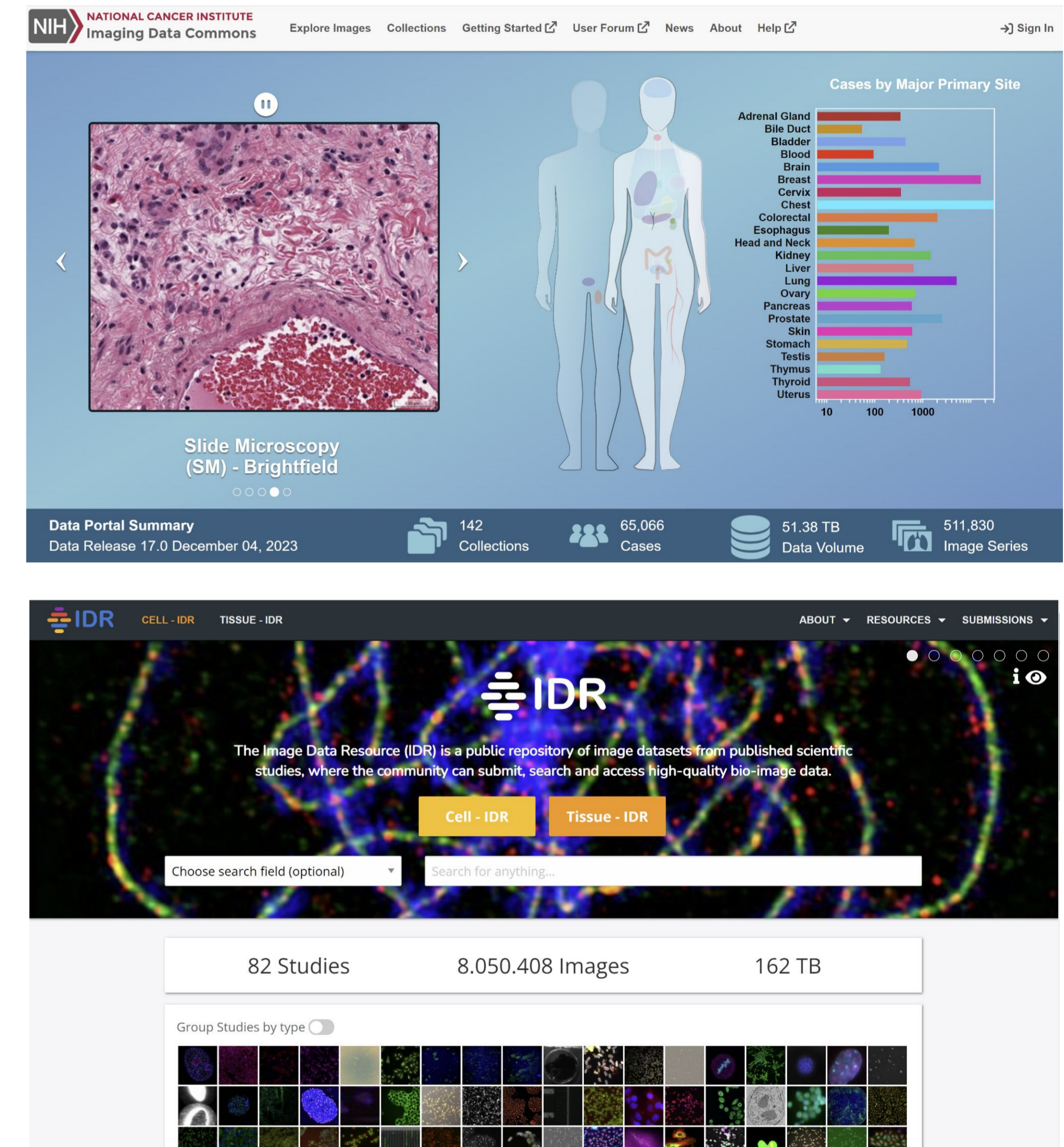
MedSAM CT image datasets (Ma et al. Nat. Communications, 2024)



One of the FAIR guidelines summaries. Source: Australian Research Data Commons



FAIR Digital Object schema



Our data list contained observations count, imaging targets, anatomical regions, pathology, provisioning and license information; data owner contact, imaging modalities information, number of sequences, geographical region and scanner metadata, as well as free-format remarks.

Foundation Models require a lot of data. It is instructive that public data is machine-actionable. With modern developments like **data spaces** and **trusted research environments**, we should define what AI-ready means and how to organize data curation for all stakeholders.

Dataset Name (might be ambiguous)	axes (N)	Anatomical Region	Main dataset topic	Targets at pixel level	Targets at instance level	Targets at image level	Any targets not correctly annotated? (e.g. only a subset lesions)	Label types (entity check)	Pathologies annotated?	Non-Pathological ROI annotated?	Provision Status	License	Usable for Scientific publication	Weighted Publishable	Data Owner (Email)	Platform	Access
1 UNIC	1,810	Lung	Lung nodule detection		Nodules	multiple/benign		Pixel-level	Y	N	Found online (directly downloadable)	CC BY 3.0				TCIA	Public (download)
2 North American Prostate Longitudinal Study (NAPLS)																	
3 ACRN 667	16,104	Breast	Breast MRI								Found online (directly downloadable)	CC BY 4.0				TCIA	Public (download)
4 OBIA	4,136	Whole Body	Chinese Cancer Imaging								Found online (access request required)					Open Website	On Request
5 Duke Breast Cancer MRI	5,161	Breast	Breast Cancer		Breast			Pixel level									
6 AOMIC ID190C		Brain									Downloaded	CC BY-SA 4.0	Yes	Yes		Open Website	
7 Duke Breast Cancer MRI	5,161	Breast	Invasive breast cancer		Resecting boxes of primary lesion			Bounding box	Y	Y	Found online (directly downloadable)	CC BY-NC 4.0				TCIA	Public (download)
8 Longitudinal evaluation of frontalotemporal dementia																	
9 AMIPET 2022 Challenge	1,014	Whole Body	Melanoma, lymphoma, lung cancer segmentation		Melanoma, lymphoma, lung cancer segmentation			Pixel-level	Y		Found online (access request required)	TCIA Research				TCIA	On Request
10 LUNA16	868	Lung	Lung nodule detection		Nodules			Bounding box	N	N	Not Defined						
11 International Consortium for Brain Mapping (ICBM)																	
12 Max. Z. not sure if interesting since labels CT CBM ONCOGRAPHY ACRN 667		Abdominal	Polyps in colon				Site labels by polyp locations	Weak (eg. scribbles)			Found online (directly downloadable)	CC BY 3.0				TCIA	Public (download)
13 ISPY2 Breast Dynamic Contrast Enhanced MRI Trial	2,658	Breast	Drug response in breast cancer		Tumor			Pixel level	Y		Found online (directly downloadable)	CC BY 4.0				TCIA	Public (download)
14 Pediatric Brain Tumors U11 D		Brain	Tumor segmentation		edema, ce tumor			Pixel level	Y	N	Not Defined	Custom DUA				Clinical Cooperation	Private
15 FastMRI Knee	1,040	Knee	make MRI scans right 10X faster					External Labels			Downloaded	Custom DUA				Open Website	On Request
16 Hector MICCAI Challenge	668	Head and Neck	Tumor segmentation and outcome prediction					Pixel level	Y	N	Downloaded & preprocessed/annotated	Custom DUA				Open Website	On Request
17 BraTS2020	2,645	Brain	Brain Tumor Segmentation		whole tumor, tumor core, enhancing tumor			Pixel level	Y	N	Found online (directly downloadable)	Custom DUA	Yes	Yes		Open Website	Public (download)
18 OASIS4		Brain									Downloaded	Custom DUA				Open Website	Public (download)
19 CT Images in COVID-19 Care-CAN Dataset		Lung	Covid 19					No Label			Found online (directly downloadable)	CC BY 4.0				TCIA	Public (download)
20		Brain	Healthy Brain					No Label			Downloaded	Custom DUA				Open Website	On Request
21 Advanced MRI Breast Lesions	4,011	Breast	Breast Lesions					Pixel level + Image level			Found online (directly downloadable)	CC BY 4.0				TCIA	Public (download)
22 (New) UPPEN GBM	3,301	Brain	Glioblastoma					Pixel level	Y	Y	Found online (directly downloadable)	CC BY-NC 4.0				TCIA	Public (download)
23 TCIA-HNSCC		Head and Neck	Head and neck squamous cell carcinoma (HNSCC)			Clinical data		Image-level			Found online (directly downloadable)	TCIA Research				TCIA	Public (download)
24																	
25 AMOS 2022	660	Abdominal	Abdominal multi-organ segmentation		spleen, right kidney, left kidney, gallbladder, esophagus, liver, stomach, aorta, inferior vena cava, pancreas, right adrenal gland, left adrenal gland, testicles, prostate/uterus			Pixel-level	N	N	Downloaded	CC BY-NC-SA	Yes	Yes		Open Website	Public (download)
26 DU Dataset		Brain	Healthy Brain					No Label			Downloaded	CC BY-SA 3.0				Open Website	Public (download)

Our dataset list snippet. Some columns and rows, like scanner and geographical information, are not shown.