1    # Introduction to Special Issue on Datasets hosted in

2    # The Cancer Imaging Archive (TCIA)

3

4    Justin Kirby[1], Fred Prior[2], Nicholas Petrick[3], Lubomir Hadjiski[4], Keyvan Farahani[5], Karen

5    Drukker[6], Jayashree Kalpathy-Cramer[7], Carri Glide-Hurst[8], Issam El Naqa[9]*

6

7    [1]Frederick National Laboratory for Cancer Research, Cancer Imaging Informatics Lab, National

8    Institute of Health, Frederick, MD 2170, USA

9

10    [2]Department of Biomedical Informatics, University of Arkansas for Medical Sciences, Little

11    Rock AR 72205, USA

12

13    [3]Center for Devices and Radiological Health, U.S. Food and Drug Administration, Silver Spring,

14    MD 20993, USA

15

16    [4]Department of Radiology, University of Michigan, Ann Arbor, MI 48109, USA

17

18    [5]Center for Biomedical Informatics and Information Technology, National Cancer Institute,

19    Bethesda, MD 20814, USA

20

21    [6]Department of Radiology, Chicago, IL 60637, USA.

22

23    [7]Department of Radiology, Charlestown, MA 02129, USA

24

25    [8]Department of Radiation Oncology, University of Wisconsin, Madison, WI  53792, USA

26

27    [9]Department of Machine Learning, Moffitt Cancer Center, Tampa, FL 33612, USA

28

29

30  **Corresponding author email:**

31  *Issam.elnaqa@moffitt.org

32

33

35

36

37

38

39

40

41

42

43

44

45

46

47  **Introduction**

48

49  Public datasets play a key role in enabling the medical research community to validate and build

50  upon each other works using data acquired outside of their home institutions.  This is especially

51  critical for stimulating studies utilizing quantitative data analysis (radiomics) or artificial

52  intelligence/machine learning (AI/ML) approaches, for which validation and generalizability on

53  independent external cohorts are essential for acceptance and future clinical translation.

54  Recognizing this fact, the Journal of Medical Physics has introduced a new category of article

55  submissions known as *Medical Physics Dataset Articles* (MPDAs)[1]. MPDAs help facilitate the

56  use of valuable open-access datasets by granting authors the opportunity to publish detailed

57  scientific or clinical descriptions of their data with unique digital object identifiers (DOIs) for

58  future citations.  Unlike traditional manuscripts, these articles would focus on reproducibility and

59  the dataset's potential use cases and details of how it was acquired, curated, and published.

60

61     This special issue was organized in partnership with <u>The Cancer Imaging Archive</u> (TCIA)[2].  TCIA

62     is an official image repository of the National Cancer Institute (NCI), and the preferred digital

63     repository for sharing cancer-related datasets described by the MPDA readership[3]. Its mission is

64     to provide proper de-identification and hosting services to relieve individual researchers of the

65     legal and technical complexities of sharing patient datasets. Image datasets are organized as

66     "collections"; typically focused on a common disease (e.g., lung cancer), image modality (MRI,

67     CT, digital histopathology, etc.) or research focus (e.g., quantitative imaging). TCIA is currently

68     home to 126 datasets[4] collected as part of numerous NCI-funded clinical trials and data sharing

69     initiatives[5,6] as well as datasets proposed by investigators in the broader research community[7].

70

71     In many cases the submitter(s) of TCIA datasets may include radiology or pathology annotations,

72     image classifications, segmentations, radiomics features, or derived/reprocessed images.

73     However, there are often cases where those who access the data on TCIA may perform their own

74     analyses, which can result in additional image labels.  In order to further support the enrichment

75     of existing datasets with these additional labels, TCIA has begun accepting proposals for third

76     party "Analysis Results" based on existing image collections.  Sharing such analyses is critical,

77     not only to enhancing medical studies reproducibility and reusability, but also to providing

78     significant value to the data science community in the form of labeled image sets for training new

79     AI/ML algorithms and other automated analysis approaches. Currently TCIA contains 28 such

80     datasets[8] several of which were submitted in relation to this special issue.

81

82     The aim of this special issue is to highlight valuable examples of MPDAs and publicly available

83     datasets that can be reused for future research endeavors and utilized for addressing emerging

84     scientific or clinical questions.

85

86     In "Head and neck cancer patient images for determining auto-segmentation accuracy in T2-

87     weighted magnetic resonance imaging through expert manual segmentations" by Cardenas, et al.[9]

88     describe a T2-weighted MRI dataset of 55 head and neck cancer patients that can be used to

89     evaluate the accuracy of auto-segmentation systems delineating organs at risk (OAR) through

90     comparisons to expert manual segmentations. The dataset can further complement existing CT

91  datasets, where MR soft tissue discrimination can further aid results for treatment planning, for
92  instance.

93

94  In "FAIR-compliant clinical, radiomics and DICOM metadata of RIDER, Interobserver, Lung1
95  and Head-Neck1 TCIA collections" by Kalendralis, et al.[10] describe updated clinical data,
96  radiomics features and Digital Imaging and Communications in Medicine (DICOM) headers from
97  4 datasets analyzed as part of their Nature Communications radiomics study[11] in order to support
98  repeatability, reproducibility, generalizability and transparency in radiomics research, which can
99  be used as useful benchmark for future CT radiomics studies.

100

101  In "DICOM Re-encoding of Volumetrically Annotated Lung Imaging Data Consortium (LIDC)
102  Nodules" by Fedorov, et al.[12] describe annotations for lung nodules from 875 of the subjects
103  collected by the Lung Imaging Data Consortium and Image Database Resource Initiative (LIDC)
104  converted into standard DICOM objects to simplify reuse of the data with the readily available
105  open-source tools, and to improve adherence to FAIR (Findable, Accessible, Interoperable,
106  Reusable) principles[13].

107

108  In "PleThora: Pleural effusion and thoracic cavity segmentations in diseased lungs for
109  benchmarking chest CT processing pipelines" by Kiser, et al.[14] describe a dataset of thoracic cavity
110  segmentations and discrete pleural effusion segmentations annotated on 402 CT scans acquired
111  from patients with non-small cell lung cancer (NSCLC). These data can be used for developing
112  image analysis pipelines such as lung structure segmentation, lesion detection, and radiomics
113  feature extraction. Combining these pleural effusion segmentations with the gross tumor volume
114  segmentations already available from the "NSCLC Radiomics" dataset, which will also enable
115  investigation of radiomics profile differences between effusion and primary tumors.

116

117  In "Reproducibility analysis of multi-institutional paired expert annotations and radiomic features
118  of the Ivy Glioblastoma Atlas Project (Ivy GAP) dataset" by Pati, et al.[15] describe their analyses
119  and resulting data from 31 IvyGAP subjects including multi-institutional expert annotations for
120  tumor sub-compartments, radiomic features, and the associated reproducibility meta-analysis to

121 facilitate developing image-based biomarkers for prognostic/predictive applications in patients
122 with glioblastoma.
123
124 In "CT images with expert manual contours of thoracic cancer for benchmarking auto-
125 segmentation accuracy" by Yang, et al.[16] describe a well-curated computed tomography (CT)
126 dataset of high-quality manually drawn contours from 60 patients with thoracic cancer that can be
127 used to evaluate the accuracy of thoracic normal tissue auto-segmentation systems.
128
129 In "MRQy: An Open-Source Tool for Quality Control of MR Imaging Data" by Sadri, et al.[17]
130 describe how they used MRQy, an open-source quality control tool to analyze TCIA collections
131 with data that was submitted from multiple sites. The results can be used for: (a) interrogating
132 MRI cohorts for site- or equipment-based differences, and (b) quantifying the impact of MRI
133 artifacts on relative image quality. This information can help determine how to correct for these
134 variations prior to model development and assess future harmonization techniques.
135
136 In summary, this special issue and its related datasets will serve as a valuable resource to help
137 develop benchmarks for a wide variety of imaging applications including image processing,
138 quality assurance, diagnostic, prognostic, and radiomics approaches using rich, annotated CT
139 and/or MR datasets. This will further strengthen the value of these datasets, their utility and
140 potential impact in the field of medical physics with the overarching goal of encouraging the
141 creation of new public datasets through MPDA/TCIA and their dissemination in the field.
142
143

144 **References**

145 1. Williamson, J.F., Das, S.K., Goodsitt, M.S. and Deasy, J.O. (2017), Introducing the
146    Medical Physics Dataset Article. Med. Phys., 44: 349-350. doi:10.1002/mp.12003
147 2. Clark K, Vendt B, Smith K, Freymann J, Kirby J, Koppel P, Moore S, Phillips S, Maffitt
148    D, Pringle M, Tarbox L, Prior F. The Cancer Imaging Archive (TCIA): maintaining and
149    operating a public information repository. J Digit Imaging. 2013 Dec;26(6):1045-57. doi:
150    10.1007/s10278-013-9622-7. PMID: 23884657; PMCID: PMC3824915.
151 3. URL: https://www.aapm.org/pubs/MPJ/policies/details.asp?id=465&type=MP.

152    4. URL: https://www.cancerimagingarchive.net/collections/

153    5. URL: https://www.cancerimagingarchive.net/imaging-proteogenomics/

154    6. URL: https://wiki.cancerimagingarchive.net/x/BQHDAg

155    7. URL: https://www.cancerimagingarchive.net/primary-data/

156    8. URL: https://www.cancerimagingarchive.net/tcia-analysis-results/

157    9. Cardenas CE, Mohamed ASR, Yang J, Gooding M, Veeraraghavan H, Kalpathy-Cramer
158       J, Ng SP, Ding Y, Wang J, Lai SY, Fuller CD, Sharp G. Head and neck cancer patient
159       images for determining auto-segmentation accuracy in T2-weighted magnetic resonance
160       imaging through expert manual segmentations. Med Phys. 2020 Jun;47(5):2317-2322.
161       doi: 10.1002/mp.13942. PMID: 32418343; PMCID: PMC7322982.

162    10. Kalendralis P, Shi Z, Traverso A, Choudhury A, Sloep M, Zhovannik I, Starmans MPA,
163        Grittner D, Feltens P, Monshouwer R, Klein S, Fijten R, Aerts H, Dekker A, van Soest J,
164        Wee L. FAIR-compliant clinical, radiomics and DICOM metadata of RIDER,
165        interobserver, Lung1 and head-Neck1 TCIA collections. Med Phys. 2020 Jun 10. doi:
166        10.1002/mp.14322. Epub ahead of print. PMID: 32521049.

167    11. Aerts HJWL, Velazquez ER, Leijenaar RTH, et al. Decoding tumour phenotype by
168        noninvasive imaging using a quantitative radiomics approach. Nat Commun.434
169        2014;5(1):4006. doi:10.1038/ncomms5006

170    12. Fedorov, A., Hancock, M., Clunie, D., Brochhausen, M., Bona, J., Kirby, J., Freymann,
171        J., Pieper, S., J. W. L. Aerts, H., Kikinis, R. and Prior, F. (2020), DICOM re-encoding of
172        volumetrically annotated Lung Imaging Database Consortium (LIDC) nodules. Med.
173        Phys. doi:10.1002/mp.14445

174    13. Wilkinson MD, Dumontier M, Aalbersberg IJJ, et al. The FAIR Guiding Principles for
175        scientific data management and stewardship. Sci Data. 2016;3:160018.
176        doi:10.1038/sdata.2016.18

177    14. Kiser, K.J., Ahmed, S., Stieb, S., Mohamed, A.S.R., Elhalawani, H., Park, P.Y.S., Doyle,
178        N.S., Wang, B.J., Barman, A., Li, Z., Zheng, W.J., Fuller, C.D. and Giancardo, L. (2020),
179        PleThora: Pleural effusion and thoracic cavity segmentations in diseased lungs for
180        benchmarking chest CT processing pipelines. Med. Phys.. doi:10.1002/mp.14424

181    15. Pati, S., Verma, R., Akbari, H., Bilello, M., Hill, V.B., Sako, C., Correa, R., Beig, N.,
182        Venet, L., Thakur, S., Serai, P., Min Ha, S., Blake, G.D., Taki Shinohara, R., Tiwari, P.

183        and Bakas, S. (2020), Reproducibility analysis of multi-institutional paired expert

184        annotations and radiomic features of the Ivy Glioblastoma Atlas Project (Ivy GAP)

185        dataset. Medical Physics. doi:10.1002/mp.14556

186    16. Yang, J., Veeraraghavan, H., van Elmpt, W., Dekker, A., Gooding, M. and Sharp, G.

187        (2020), CT images with expert manual contours of thoracic cancer for benchmarking

188        auto-segmentation accuracy. Med. Phys., 47: 3250-3255. doi:10.1002/mp.14107

189    17. Sadri, A.R., Janowczyk, A., Zhou, R., Verma, R., Beig, N., Antunes, J., Madabhushi, A.,

190        Tiwari, P. and Viswanath, S.E. (2020), Technical Note: MRQy — An open-source tool for

191        quality control of MR imaging data, Medical Physics, doi: 10.1002/mp.14593.