

1 Editor Summary: This Perspective provides a roadmap toward globally federated open data  
2 resources for bioimaging data.

3  
4 Peer review information:

5 Primary Handling editor: Rita Strack, in collaboration with the Nature Methods team.

6  
7 Peer review information: Nature Methods thanks Kelly Rogers and the other, anonymous,  
8 reviewer(s) for their contribution to the peer review of this work.

## 11 **Enabling Global Image Data Sharing in the Life Sciences**

13 Peter Bajcsy (1) Sreenivas Bhattiprolu (2), Katy Börner (3), Beth A Cimini (4), Lucy Collinson (5), Jan Ellenberg  
14 (6), Reto Fiolka (7), Maryellen Giger (8), Wojtek Gosinski (9), Matthew Hartley (10), Nathan Hotelling (11), Rick  
15 Horwitz (12), Florian Jug (13), Isabel Kemmer (20), Anna Kreshuk (6), Emma Lundberg (14), Aastha Mathur (20),  
16 Kedar Narayan (15), Shuichi Onami (16), Anne L. Plant (1), Fred Prior (17), Jason R. Swedlow (18), Adam Taylor  
17 (19), and Antje Keppler\* (20)

18  
19 Affiliations:

- 20 1. National Institute of Standards and Technology, Gaithersburg, USA.
- 21 2. ZEISS Microscopy Customer Center, Dublin, USA
- 22 3. Intelligent Systems Engineering, Indiana University, Bloomington, USA
- 23 4. Imaging Platform, Broad Institute, Cambridge, MA USA
- 24 5. Electron Microscopy Science Technology Platform, Francis Crick Institute, 1 Midland Road, London, UK
- 25 6. European Molecular Biology Laboratory, Heidelberg, Germany
- 26 7. Lyda Hill Department of Bioinformatics, UT Southwestern Medical Center, Dallas, Texas, USA.
- 27 8. Department of Radiology and Committee on Medical Physics, University of Chicago, Chicago, IL USA
- 28 9. National Imaging Facility, Brisbane, Australia
- 29 10. European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton, Cambridge, UK
- 30 11. National Center for Advancing Translational Science, National Institutes of Health, Rockville, USA
- 31 12. Allen Institute for Cell Science, Seattle, USA
- 32 13. Fondazione Human Technopole, Milan, Italy
- 33 14. Stanford University, California, US and SciLifeLab, KTH Royal Institute of Technology, Stockholm,  
34 Sweden
- 35 15. Center for Molecular Microscopy, Center for Cancer Research, National Cancer Institute, National  
36 Institutes of Health, Bethesda USA. and Cancer Research Technology Program, Frederick National  
37 Laboratory for Cancer Research, Frederick USA.
- 38 16. RIKEN Center for Biosystems Dynamics Research, Kobe, Japan
- 39 17. Department of Biomedical Informatics, Department of Radiology, University of Arkansas for Medical  
40 Sciences, Little Rock, USA
- 41 18. Divisions of Computational Biology and Molecular, Cell and Developmental Biology, University of  
42 Dundee, UK
- 43 19. Sage Bionetworks, Seattle, WA, USA
- 44 20. Euro-BioImaging Bio-Hub, European Molecular Biology Laboratory, Heidelberg, Germany

45  
46 \*Corresponding author: [antje.keppler@eurobioimaging.eu](mailto:antje.keppler@eurobioimaging.eu).

### 47 *Abstract*

48  
49 Despite the importance of imaging in biological and medical research, a large body of informative and  
50 precious image data never sees the light of day. To ensure scientific rigour, as well as the reuse of data  
51 for scientific discovery, image data needs to be made FAIR (Findable, Accessible, Interoperable,  
52 Reusable). Image data experts are working together globally to agree on common data formats,

53 metadata, ontologies and supporting tools towards image data FAIRification. With this paper we call  
54 on public funders to join these efforts to support their national scientists. What researchers most  
55 urgently need are openly accessible resources for image data storage which are operated under long-  
56 term commitments by their funders. Although existing resources in Australia, Japan and Europe are  
57 already collaborating to enable global image data sharing, these efforts will fall short unless more  
58 countries invest in operating and federating their own open data resources. This will allow us to  
59 harvest the enormous potential of existing image data, preventing significant loss of unrealised value  
60 from past investments in imaging acquisition infrastructure.

61

62

### 63 *Introduction*

64 Image data are not only intriguing to look at but they are also the fastest growing data resource in the  
65 life sciences. Biological and medical images are complex and deep in information content based on  
66 the nature of many imaging modalities making the data conducive to reuse. With simultaneous rapid  
67 advances in and access to computation power, the scientific community is only beginning to tap into  
68 this invaluable resource. Acquiring and managing image data according to the FAIR principles <sup>1</sup>  
69 (Findable, Accessible, Interoperable, Reusable) is therefore being requested by major public funders  
70 such as the European Union (<https://digital-strategy.ec.europa.eu/en/policies/data-act>) and NIH  
71 (<https://grants.nih.gov/grants/guide/notice-files/NOT-OD-21-013.html>). Beyond supporting their Open  
72 Science policy objectives, many of the most exciting opportunities for deriving value from sharing and  
73 reusing imaging data at scale rely on the ability to curate and centrally access those data. This  
74 aggregation based on common standards not only enables very large scale AI/ML model training, but  
75 also the creation of cross-modality and cross-domain benchmarks and reference datasets. These gains  
76 are not possible when FAIR data are siloed and fragmented in many different places, often without  
77 public access or consistent organization. The key to easy data access lies in consolidating data in a  
78 limited number of locations that are connected with each other and which follow compatible  
79 standards for data and metadata storage, i.e., they are interoperable.

80 Australia, Japan, and the European Union have established nationally funded, coordinated research  
81 infrastructures (RIs) for life science and biomedical imaging (collectively "BioImaging"). These  
82 BioImaging RIs are the output of concerted long-term strategic planning efforts from key stakeholders  
83 that include academic, funding, political as well as industrial partners. They are now in operation and  
84 deliver technology and community support for quality image data acquisition, management, analysis,  
85 and publication for their stakeholders. They are supported by dedicated image data repositories, like  
86 the BioImage Archive<sup>2</sup>, IDR<sup>3</sup> and EMPIAR<sup>4</sup> in Europe, SSBD<sup>5</sup> in Japan, and AIS repository in Australia  
87 (<https://australian-imaging-service.github.io/>) that together hold around 6 PB of image data. Publicly  
88 funded projects and community initiatives are also actively building and supporting image data  
89 formats (<https://ngff.openmicroscopy.org/>) and applications for managing, analysing, and sharing data  
90 using the most advanced technologies, including AI (<https://ai4life.eurobioimaging.eu/>). Collaboration  
91 between academic and commercial organizations is maturing and beginning to provide a powerful  
92 ecosystem for development and eventual scaling of new technologies and products.

93 However, these regionally limited efforts are only a starting point of a future global ecosystem for  
94 sharing image data. We therefore urge funders to invest in open access infrastructure supporting  
95 image data acquisition, management, and sharing, which then can be coordinated around the world.  
96 This is a necessary step towards realizing the significant value that is currently lying dormant in  
97 underutilized image data and instead can catalyze the exploitation of this invaluable resource which  
98 was largely generated through previous public investments.

99

100 *Image data sharing – Cui bono?*

101 Everyone can benefit tremendously if open access image data resources were available to them  
102 (Figure 1). Life scientists, as producers of large image datasets, can strengthen their research output  
103 by opening up the data that supports their publications and exchange the data more easily with  
104 existing and future researchers. AI experts, data analysts, and software developers can use the  
105 available image datasets for training, testing, and developing their tools. Trainers and educators can  
106 expand their portfolio of training material and expedite the sharing of knowledge and new tools.  
107 Researchers can follow their data management plan and upload their final data for publication in these  
108 resources, meaningfully fulfilling the requirements of their funders and publishers. Funders can  
109 sustain their research investments and enable the reuse of existing research data for creation of new  
110 knowledge. Finally, the image data becomes accessible to the general public, which can lead to  
111 increased trust in and support of public research and education, as well as new and unanticipated  
112 discoveries.

113  
114 As a key technology in the biological and biomedical sciences and beyond, imaging is enabling  
115 scientists to map and understand the complex spatiotemporal patterns of biological processes in  
116 greater detail than ever before. While other fields like genomics, proteomics etc., have benefited  
117 heavily from investing into extensive data sharing mechanisms and practices, recent advances in  
118 computational infrastructure have only now given the imaging community the same level of aspiration  
119 to make image data an accessible, shared global resource.

120  
121 We are convinced that sharing of FAIR image data from all domains will significantly increase  
122 interoperability among the different research fields in the life sciences. Most importantly, it will allow  
123 a comprehensive view of biological processes at different levels of organization, from molecular to  
124 cellular to whole organism and thereby will advance both basic biological as well as translational  
125 research. The integration of microscopy and medical imaging can provide a better understanding of  
126 disease mechanisms, particularly in the context of complex diseases that involve multiple organ  
127 systems and cellular processes. Preclinical and clinical imaging techniques, such as magnetic  
128 resonance imaging (MRI), provide valuable data on anatomical structures and physiological functions.  
129 Integrating this data with light microscopy images can provide a more accurate and detailed  
130 understanding of the underlying biology. Such image data integration will support and grow the field  
131 of bioimage informatics and trigger advancements in drug discovery and development and help to  
132 identify potential therapeutic targets, as well as to evaluate drug efficacy and safety. If we foster image  
133 data sharing across domains, then we will not only democratize science but also enable new  
134 discoveries by linking different types of image data sets across the entire spectrum of life.

135  
136 *What are the technical challenges?*

137 There are several technical challenges that make sharing image data more difficult than sharing  
138 genomic data<sup>6</sup> or protein structures<sup>7</sup>, which are currently the best examples of open shared data and  
139 data repositories. One is that the data may be very large with individual files of sizes from 100 GBs to  
140 10 TBs, and beyond<sup>8</sup> and require either high-speed internet connections to be uploaded and  
141 downloaded or physically transported on external drives. These practices lead to duplication of data  
142 making data management and sharing expensive and complex for large datasets over time. For  
143 supporting interconnectivity between research institutions and data repositories, a robust technical  
144 infrastructure, including large data storage capacity and powerful computing resources is therefore  
145 required, as many research institutions still do not have the necessary infrastructure to support image  
146 data sharing.

147  
148 Another key challenge is the need for careful management and organisation to ensure that image data  
149 are accessible and, ideally, reusable by other researchers. While medical images typically are easily  
150 shared, since most are acquired in DICOM format<sup>9</sup>, biological imaging data sharing is still limited  
151 because of the lack of standardisation in image data formats, ontologies, and metadata. This makes it

152 difficult for researchers to compare and analyse data from different sources. Ensuring that the data  
153 are properly labelled, annotated, and stored in a standardised format can require significant time and  
154 resources<sup>10</sup>, however, it is an essential, yet often underappreciated process. In addition, proper  
155 curation and harmonization of imaging data and associated metadata are needed to ensure machine  
156 actionability in future use of data sets at scale.

157  
158 Finally, depending on the context, intellectual property issues, questions of ownership, and data  
159 privacy concerns may arise. Because of this, researchers may be reluctant to share their image data,  
160 and in case of medical imaging data, they must take steps to ensure that the data are de-identified  
161 and cannot be traced back to individual patients. Medical image datasets must remain in their country  
162 of origin, and the researchers together with the legislative authorities need to find common legal and  
163 ethical solutions. Also, there may be cultural barriers, such as a lack of trust or reluctance to share  
164 image data due to concerns about competition or ownership.

165  
166  
167 *Which solutions are currently developed?*

168 Despite these challenges, community initiatives around the world are actively working on solutions to  
169 share image data globally<sup>11</sup>. However, those developments aren't sustainable without funding to  
170 support open access image data storage. We must move forward with the input from all stakeholders  
171 involved in the process of image data production, sharing and analysis, including scientists acquiring  
172 data, operators of large-scale data infrastructures, and data experts. Educational opportunities  
173 addressing the importance of quality-managed and harmonized data are increasing, but they need to  
174 be scaled and embedded in each researcher's training curriculum from the beginning. It is essential to  
175 recognise and reward the efforts of everyone involved in producing and curating FAIR data, including  
176 specialized data support staff, such as FAIR data stewards, data curators, research software engineers,  
177 etc. Many current funding plans to promote data sharing fail to adequately include these roles,  
178 resulting in poor implementation and limited uptake of developed solutions. While establishment of  
179 these positions is on the rise, wider recognition and sustained funding is essential to support and  
180 secure these invaluable positions. Efforts similar to those driven by Global BioImaging for imaging  
181 scientists working in core facilities<sup>12</sup>, are needed to define the career progression and to develop  
182 measures to support professional development and job security of these important positions.

183  
184 In the long term, global sharing of image data can only be achieved with a globally federated model.  
185 This requires that added value data resources dedicated to particular use cases or domains are stored  
186 and made available through a small number of centralized data archives or other open access servers.  
187 These archives must provide both direct data hosting as well as indexing and search capabilities. This  
188 model must be built on shared core metadata models for interoperability based on international  
189 guidelines, such as REMBI<sup>13</sup>, which can be extended for domain-specific applications such as MITI for  
190 highly multiplexed tissue images<sup>14</sup>. In addition to metadata formats, effective use of data requires  
191 standardised image formats. In light of the exponential growth in the volume of new image data, it is  
192 critical that the imaging community collectively agrees on standard image formats with immediate  
193 effect. To guarantee the widest possible adoption, it is essential to define these formats in discussion  
194 with vendors, a process that has long been carried out by tireless community efforts driven by  
195 initiatives like OME<sup>15</sup>, which is now working towards future proof Next Generation File Formats  
196 designed for large-scale, cloud-enabled use, such as OME-Zarr<sup>16,17</sup>. Allowing data to be easily located,  
197 accessed through consistent mechanisms in standardised formats, and coupled with rich metadata  
198 unlocks valuable large-scale applications. These include training universal image classifiers to support  
199 automated quality control pipelines for researchers and instrument manufacturers, automated  
200 segmentation of common cell structures for use as biomarkers of disease or to better understand  
201 mechanisms of action, and prediction of physiological or biomolecular properties of cells to speed up  
202 bio-manufacturing and reduce assay variability and consumable use etc. By adopting these  
203 approaches, we can transform petabytes to exabytes of immensely valuable scientific data,

204 representing enormous funder investments, into a global resource that will be maximally useful for  
205 the global community.

206

207 *How to go forward?*

208 While enabling global image data sharing requires a concerted effort and endurance from researchers,  
209 institutions, publishers, and industry to establish internationally recognized best practices, it foremost  
210 requires a long-term commitment by public funders and their investments in essential technical  
211 infrastructure and data resources. Funders need to recognise that they currently lose significant value  
212 from their previous investments into imaging acquisition infrastructure unless they provide the  
213 essential resources for storing and making the resulting image data accessible. It is essential that the  
214 imaging community and funders join forces to establish a sustainable and profitable image data  
215 ecosystem ensuring optimal output for all stakeholders (Figure 2). The authors of this paper are  
216 looking forward to an exciting and fruitful dialogue with national funders on the suggested action  
217 items presented here.

218

219

## 220 **Acknowledgements**

221 K. N. is funded by Federal funds from the National Cancer Institute, National Institutes of  
222 Health, under Contract No. 75N91019D00024.

223 R.F. is funded by the National Cancer Institute, (grant R33 CA235254 and U54CA268072)  
224 and the National Institute of General Medical Sciences (grant R35 GM133522).

225 I.K. receives funding from EU Horizon 2020: Grant Agreement No. 101046203 (BY-COVID)  
226 and from the European Union under Grant Agreement Number 101130986 (EVOLVE).

227

228 This manuscript is based on the more extensive white paper “Enabling Global Image Data  
229 Sharing in the Life Sciences” ([arXiv:2401.13023](https://arxiv.org/abs/2401.13023)) which is published with a closely related  
230 companion entitled, *Harmonizing the Generation and Pre-publication Stewardship of FAIR*  
231 *Image Data*, which can be found at [arXiv:2401.13022](https://arxiv.org/abs/2401.13022).

232

233

## 234 **Author Contributions**

235 This manuscript is authored and reviewed by members of an ad-hoc International  
236 Working Group who are listed alphabetically as co-authors of the manuscript.. I.K.  
237 prepared the figures.

238

## 239 **Ethics Declaration**

240 Disclaimer: Commercial products are identified in this document in order to specify the  
241 experimental procedure adequately. Such identification is not intended to imply  
242 recommendation or endorsement by the National Institute of Standards and Technology,  
243 nor is it intended to imply that the products identified are necessarily the best available for  
244 the purpose.

245 K. N. is funded by Federal funds from the National Cancer Institute, National Institutes of  
246 Health, under Contract No. 75N91019D00024. The content of this publication does not  
247 necessarily reflect the views or policies of the Department of Health and Human Services,  
248 nor does mention of trade names, commercial products, or organizations imply  
249 endorsement by the U.S. Government.

250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275

## References

1. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* **3**, 160018 (2016).
2. Hartley, M. *et al.* The BioImage Archive - Building a Home for Life-Sciences Microscopy Data. *J. Mol. Biol.* **434**, 167505 (2022).
3. Williams, E. *et al.* Image Data Resource: a bioimage data integration and publication platform. *Nat. Methods* **14**, 775–781 (2017).
4. Iudin, A. *et al.* EMPIAR: the Electron Microscopy Public Image Archive. *Nucleic Acids Res.* **51**, D1503–D1511 (2023).
5. Tohsato, Y., Ho, K. H. L., Kyoda, K. & Onami, S. SSBD: a database of quantitative data of spatiotemporal dynamics of biological phenomena. *Bioinformatics* **32**, 3471–3479 (2016).
6. Arita, M., Karsch-Mizrachi, I. & Cochrane, G. The international nucleotide sequence database collaboration. *Nucleic Acids Res.* **49**, D121–D124 (2021).
7. wwPDB consortium *et al.* Protein Data Bank: the single global archive for 3D macromolecular structure data. *Nucleic Acids Res.* **47**, D520–D528 (2019).
8. Shapson-Coe, A. *et al.* A petavoxel fragment of human cerebral cortex reconstructed at nanoscale resolution. *Science* **384**, eadk4858 (2024).
9. Larobina, M. Thirty Years of the DICOM Standard. *Tomography* **9**, 1829–1838 (2023).
10. Swedlow, J. R. *et al.* A global view of standards for open image data formats and repositories. *Nat. Methods* **18**, 1440–1446 (2021).
11. Kemmer, I. *et al.* Building a FAIR image data ecosystem for microscopy communities. *Histochem. Cell Biol.* (2023) doi:10.1007/s00418-023-02203-7.
12. Wright, G. D. *et al.* Recognising the importance and impact of Imaging Scientists:

- 276 Global guidelines for establishing career paths within core facilities. *J. Microsc.* **294**, 397–  
277 410 (2024).
- 278 13. Sarkans, U. *et al.* REMBI: Recommended Metadata for Biological Images—enabling  
279 reuse of microscopy data in biology. *Nat. Methods* **18**, 1418–1422 (2021).
- 280 14. Schapiro, D. *et al.* MITI minimum information guidelines for highly multiplexed tissue  
281 images. *Nat. Methods* **19**, 262–267 (2022).
- 282 15. Linkert, M. *et al.* Metadata matters: access to image data in the real world. *J. Cell*  
283 *Biol.* **189**, 777–782 (2010).
- 284 16. Moore, J. *et al.* OME-NGFF: a next-generation file format for expanding bioimaging  
285 data-access strategies. *Nat. Methods* **18**, 1496–1498 (2021).
- 286 17. Moore, J. *et al.* *OME-Zarr: A Cloud-Optimized Bioimaging File Format with*  
287 *International Community Support.*  
288 <http://biorxiv.org/lookup/doi/10.1101/2023.02.17.528834> (2023)  
289 doi:10.1101/2023.02.17.528834.

290

291

## 292 **Figure Legends**

293

294 Figure 1. Data sharing: too impactful not to do. An outline of the tremendous benefits of globally  
295 shared image data. Data sharing drives Open Science, empowers life scientists, and enables reuse of  
296 valuable FAIR image datasets for creating new knowledge.

297 Figure 2: How public decision-makers, funders and the scientific community together can enable  
298 Global Image Data Sharing.

299

300

301

302

303

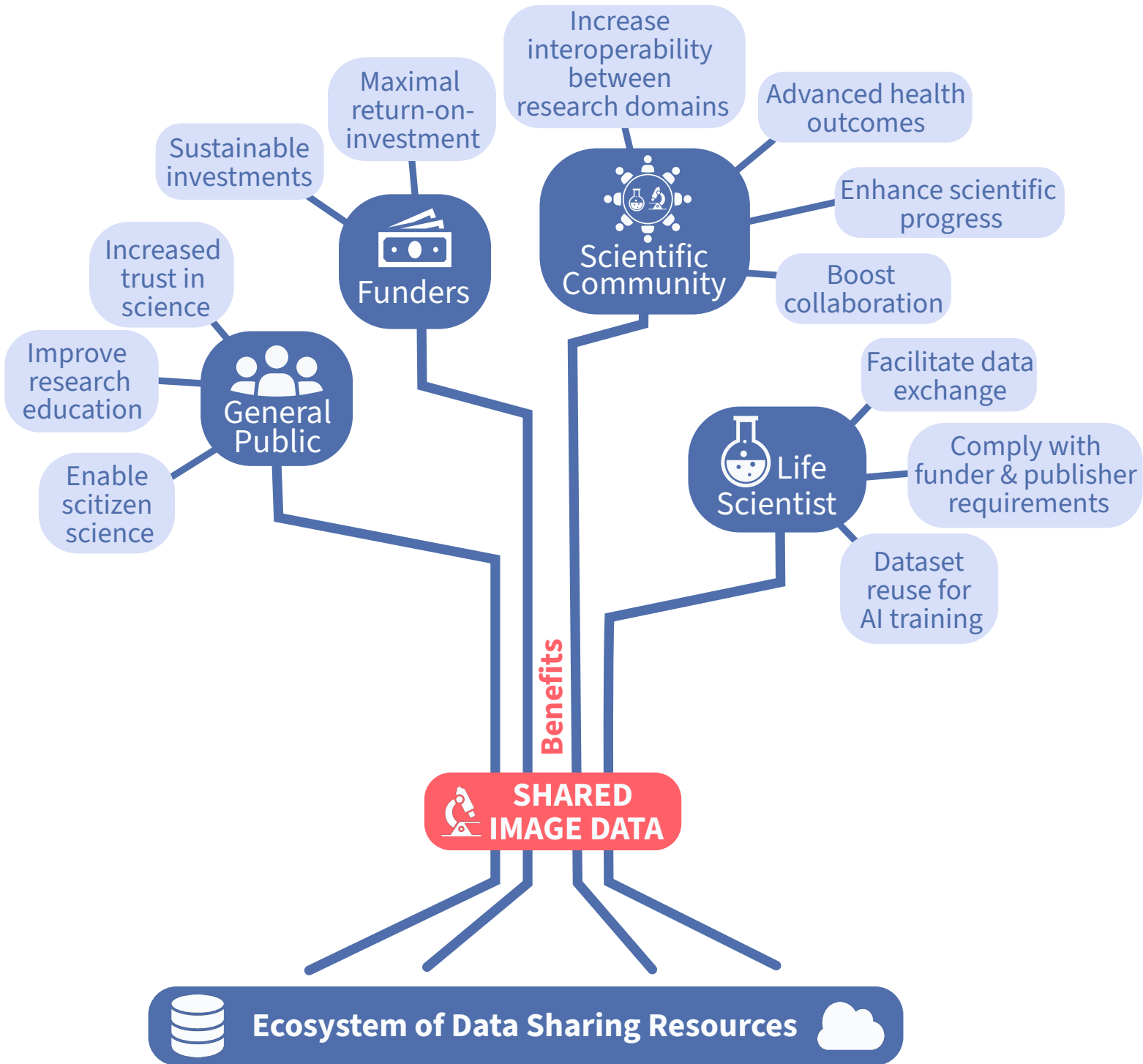
304

305

306



**Data Sharing: Too impactful not to do**



# Image Data Sharing ToDo List for Funders

To maximize the value of investments in the imaging sector we need funding for:

## Open Access Image Data Repositories

Launch, fund and support key public image data repositories long-term, creating an international federated archival system of repositories based on common data architectures supporting benchmark datasets in key areas (e.g. disease models).



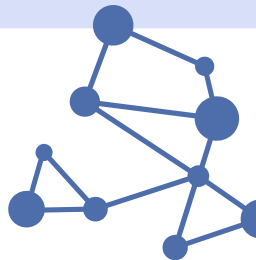
## Open Access Tools Ecosystem

Fund the creation of a tools ecosystem with access to compute resources to create, analyse and interact with FAIR image data



## Open Access Imaging Infrastructures

Fund novel approaches to image data acquisition and sharing, such as coordinated national or regional imaging infrastructures, to maximise return-on-investment.



## Specialised Staff

Develop new funding mechanisms for critical but non-traditional work for data support staff in this space e.g., to develop, maintain, and upgrade data sharing pipelines and ecosystems, as well as for software developers, data wranglers, data curators, and FAIR image data stewards



## Training and Education

Fund training and education of researchers at all career stages for acquiring and working with FAIR research data.



# Image Data Sharing ToDo List for the Imaging Community

To maximize the value of investments in the imaging sector we need to work on:

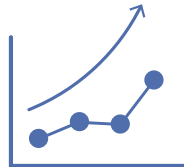
## Metadata and Ontologies

Build useful, accessible, and community recommend ontologies and implementations of metadata standards, and incentivize their uptake



## Metrics and Recognition

Develop metrics to recognize and incentivize high-quality, structured data generation that allows easy sharing and re-use.



## Training and Education

Educate researchers on benefits of image data sharing and reuse, FAIR data, as well as training and dissemination of tools.



## Data Formats

Decide on, implement and maintain standardized data formats that allow data streaming to support navigation and processing of massive image datasets.



## Data in Journals

Encourage journals to establish guidelines and develop additional specialist "Data Journals" to support re-use and citation of high value datasets, similar to Methods and Protocols journals



## Vendor Engagement

Engage with vendors to mandate that image data is accompanied by well-structured metadata in accepted formats and within data models that thoroughly describe images, with the option to export both image data and metadata into open formats.

