



University of Dundee

OME-NGFF

Moore, Josh; Allan, Christopher; Besson, Sebastien; Burel, Jean-Marie; Diel, Erin; Gault, David

Published in:
Nature Methods

DOI:
[10.1038/s41592-021-01326-w](https://doi.org/10.1038/s41592-021-01326-w)

Publication date:
2021

Licence:
CC BY

Document Version
Publisher's PDF, also known as Version of record

[Link to publication in Discovery Research Portal](#)

Citation for published version (APA):

Moore, J., Allan, C., Besson, S., Burel, J.-M., Diel, E., Gault, D., Kozlowski, K., Lindner, D., Linkert, M., Manz, T., Moore, W., Pape, C., Tischer, C., & Swedlow, J. R. (2021). OME-NGFF: a next-generation file format for expanding bioimaging data access strategies. *Nature Methods*, 18, 1496-1498. <https://doi.org/10.1038/s41592-021-01326-w>

General rights

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from Discovery Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain.
- You may freely distribute the URL identifying the publication in the public portal.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

OPEN



OME-NGFF: a next-generation file format for expanding bioimaging data-access strategies

Josh Moore¹, Chris Allan², Sébastien Besson¹, Jean-Marie Burel¹, Erin Diel², David Gault¹, Kevin Kozlowski², Dominik Lindner¹, Melissa Linkert², Trevor Manz³, Will Moore¹, Constantin Pape⁴, Christian Tischer⁴ and Jason R. Swedlow^{1,2} ✉

The rapid pace of innovation in biological imaging and the diversity of its applications have prevented the establishment of a community-agreed standardized data format. We propose that complementing established open formats such as OME-TIFF and HDF5 with a next-generation file format such as Zarr will satisfy the majority of use cases in bioimaging. Critically, a common metadata format used in all these vessels can deliver truly findable, accessible, interoperable and reusable bioimaging data.

Biological imaging is one of the most innovative fields in the modern biological sciences. New imaging modalities, probes and analysis tools appear every few months and often prove decisive for enabling new directions in scientific discovery. One feature of this dynamic field is the need to capture new types of data and data structures. While there is a strong drive to make scientific data findable, accessible, interoperable and reusable (FAIR¹), the rapid rate of innovation in imaging has resulted in the creation of hundreds of proprietary file formats (PFFs) and has prevented the unification and adoption of standardized data formats. Despite this, opportunities for sharing and integrating bioimaging data and, in particular, linking these data to other ‘omics’ datasets have never been greater. Therefore, to every extent possible, increasing ‘FAIRness’ of bioimaging data is critical for maximizing scientific value, as well as for promoting openness and integrity².

When working with a large number of PFFs, interoperability and accessibility are achieved using translation and conversion provided by open-source, community-maintained libraries that produce an open, common data representation. On-the-fly translation produces a transient representation of bioimage metadata and binary data in an open format but must be repeated on each use. In contrast, conversion produces a permanent copy of the data, again in an open format, bypassing bottlenecks in repeated data access. As workflows and data resources emerge that handle terabytes (TB) to petabytes (PB) of data, the costs of on-the-fly translation have become bottlenecks to scientific analysis and the sharing of results. Open formats like OME-TIFF³ and HDF5 (ref. ⁴) are often used for permanent conversion, but both have limitations that make them ill-suited for use cases that depend on very high and frequent levels of access, such as training of artificial intelligence models and publication of reference bioimage datasets in cloud-based resources. For these situations, the community is missing a multidimensional, multi-resolution binary container that provides parallel read-and-write capability that is natively accessible from the cloud (without server infrastructure) and that has a flexible, comprehensive metadata structure (Supplementary Note).

To this end, we have begun building OME’s next-generation file format (OME-NGFF) as a complement to OME-TIFF and HDF5.

Together these formats provide a flexible set of choices for bioimaging data storage and access at scale over the next decade and, potentially, a common, FAIR solution for all members of the biological imaging community (academic and industrial researchers, imaging scientists, and academic and commercial technology developers).

Next-generation file formats

We use the term next-generation file formats (NGFFs) to denote file formats that can be hosted natively in an object (or cloud) storage for direct access by a large number of users. Our current work, which we refer to as OME-NGFF, is built upon the Zarr format⁵ but heavily informed and connected to both TIFF and HDF5. We have compared the characteristics of these three open formats in Supplementary Table 1.

To date, the development of OME-NGFF has focused on pixel data and metadata specifications for multidimensional, multiscale images, high-content screening datasets and derived labeled images. These specifications include support for ‘chunking’ or storage of parts of the binary pixel data in smaller files that support rapid access to the data from orthogonal views or different resolution levels (also known as pyramidal data). Labeled images, such as segmentation or classification masks can now remain in a common data structure with the original pixel data and metadata, providing a single mechanism for tracking the provenance of original and derived data allowing programmatic rather than manual management.

We have also built multiple implementations of these specifications, demonstrating the usability and performance of these formats. `bioformats2raw` can be used for writing OME-NGFF from standalone Java applications and `omero-cli-zarr` is available for exporting from OMERO⁶. Reading is implemented in `ome-zarr-py`, which has been integrated into the napari viewer⁷, in Fiji via the MoBIE plugin⁸ and finally via Viv-based vizarr for access in the browser⁹. Permissively licensed example datasets from the Image Data Resource (IDR)¹⁰ have been converted into Zarr and stored in an S3-object storage bucket for public consumption (Extended Data Fig. 1). Though OME-NGFF is still in development, each of these implementations is an example of how data access and application is simplified by having a universal data-storage pattern. Current and future specifications are published under <https://ngff.openmicroscopy.org/latest/>.

Bioimage latency benchmark

To demonstrate how NGFFs complement available, open formats, we have built and published a bioimage latency benchmark that compares random, serial-access speeds to uncompressed TIFF, HDF5 and Zarr files. These measurements provide an upper bound on the overhead that a user would experience accessing the formats

¹University of Dundee, Dundee, UK. ²Glencoe Software, Inc., Seattle, WA, USA. ³Harvard Medical School, Boston, MA, USA. ⁴European Molecular Biology Laboratory (EMBL), Heidelberg, Germany. ✉e-mail: j.r.swedlow@dundee.ac.uk

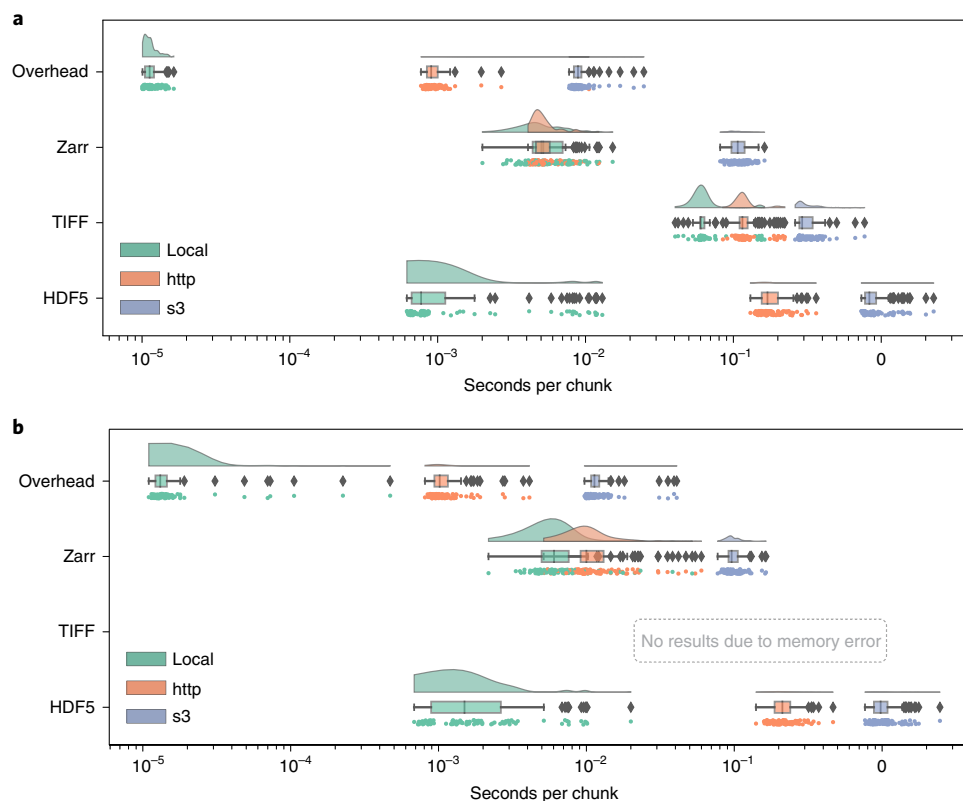


Fig. 1 | Chunk retrieval time is less sensitive to data location with next-generation file formats. a,b, Random sampling of 100 chunks from synthetically generated, five-dimensional images measures access times for three different formats on the same file system (green), over HTTP using the nginx web server (orange) and using Amazon's proprietary S3 object storage protocol (blue) under two scenarios: a whole-slide CycIF imaging dataset with many large planes of data ($x=64,000$, $y=64,000$, $c=8$) and chunks of 256×256 pixels (128 KB) (**a**); and a time-lapse LSM dataset with isotropic dimensions ($x=1,024$, $y=1,024$, $z=1,024$, $t=100$) and chunks of $32 \times 32 \times 32$ pixels (64 KB) (**b**) (Methods).

using common libraries, tiffle, h5py and Zarr-Python, respectively. Though future extensions to the benchmark are intended, we have focused on a single, serverless Python environment because one library (fsspec) can be used to access all three data formats across multiple storage mechanisms without the need for any additional infrastructure.

The benchmark includes instructions for running on Docker or AWS EC2 and contains all necessary code to regenerate representative samples for two established imaging modalities: large multi-channel two-dimensional (2D) images such as the ones produced by cyclic immunofluorescence (CycIF)¹¹ and time-lapse isotropic volumes typically generated by light-sheet microscopy (LSM)¹². Each synthetic HDF5, TIFF and Zarr dataset was generated by first invoking the ImarisWriter, then converting the HDF5-based Imaris files into Zarr with bioformats2raw and finally converting the Zarr to TIFF with raw2ometiff. All three datasets along with a 1byte dummy file for measuring overhead were placed in three types of storage: local disk, a remote server and object storage. We measured the reading time of individual chunks for all four file types across the three storage systems. Figure 1 shows that as the latency of access grows, access times for monolithic formats such as TIFF and HDF5 increase because libraries must seek the appropriate data chunk, whereas NGFF formats such as Zarr provide direct access to individual chunks. In the three-dimensional (3D) case, the TIFF data were too large to fit into local memory and the benchmark errored.

On local storage, access speeds for NGFF files were similar to HDF5 and both substantially outperformed TIFF. This matches previous results showing that a number of factors must be taken into

account to determine the relative performance of HDF5 and Zarr¹³. Together these results partially explain HDF5's popularity for desktop analysis and visualization of LSM datasets.

However, on cloud storage, access speeds for NGFF files are at least an order of magnitude faster than HDF5. Parallel reads¹⁴, supporting streaming of image data files from remote http-based or cloud-based servers give performance similar to local disk access. Data streaming obviates the need for wholesale data download and is especially important for providing performant access to multi-TB datasets.

We note that our benchmark measures direct access to underlying storage. Additional applications, such as HSDS for HDF5 or OMERO for TIFF, may improve the performance of specific use cases, but add complexity to any deployment and make direct comparisons between the different data-access regimes in Fig. 1 difficult. Additionally, a key parameter in overall access times is the size of individual chunks. As chunk sizes decrease, the number of individual chunk files increases rapidly (Extended Data Fig. 2). In this benchmark, we have chosen a compromise between chunk size and number of individual files. This illustrates a primary downside of NGFF formats; as the number of files increases, the time required for copying data between locations increases. Users will need to understand and balance these trade-offs when choosing between open, bioimaging file formats.

Outlook: community adoption

We assert that together low-latency, cloud-capable NGFF, TIFF and HDF5 can provide a balanced set of options that the community can converge upon and slow the development of ever more file

formats. To this end, OME is committed to building an interoperable metadata representation across all three file formats to ensure ease of adoption and data exchange (Supplementary Note).

When data are frequently accessed, for example, as a public resource or a training dataset, upfront conversion will lead to overall time savings. In situations where object storage is mandated, as in large-scale public repositories, we encourage the use of OME-NGFF today. Alternatively, users needing to transfer their images may choose to store their data in a large single file such as HDF5. OME-TIFF remains a safe option for those who rely on proprietary software for visualization and analysis, especially in digital pathology and other whole-slide image applications, as many have been extended to both read and write this open standard. Each choice comes with benefits and costs and individual scientists, institutions, global collaborations and public data resources need the flexibility to decide which approach is suitable. We encourage the community to choose from the most appropriate of the formats described above, secure in the knowledge that conversion is possible if it becomes necessary.

We foresee this being a critical strategy where data generated in advanced bioimaging applications is converted into an optimized format for downstream processing, analysis, visualization and sharing. All subsequent data access occurs via open data formats without the need for repeated, on-the-fly translation. We have begun implementing this workflow in the IDR (Extended Data Fig. 1), alleviating the need for time-consuming downloads and cross-referencing metadata and resulting in substantially more accessible and interoperable data. We look forward to working with other resources to further develop this policy. Further, as adoption of public image data resources increases, commercial vendors will hopefully engage with these efforts to support their customers, who are increasingly required to publish datasets as supplementary material. Moreover, some commercial imaging companies are themselves building cloud-based data handling and analysis solutions (for example, <https://www.apeer.com>), thus broadening the community of users who need cloud-competent file formats.

Ultimately, we hope to see digital imaging systems producing open, transparent (in other words FAIR), data without the need for further conversion. Until that time, we are committed to providing the data conversion needs of the community. Following the same pattern established by `bioformats2raw` and `raw2ometiff`, we propose to meet this challenge via a set of migration tools allowing efficient data transformations between all data formats contained in this suite of interoperable formats. Additionally, as the specification evolves based on community feedback, the same migration tools will allow upgrading the scientific data generated by the bioimaging community to prevent the need for long-term maintenance of older data. Upcoming specifications include geometric descriptions of regions of interest, meshes and transformations for correlative microscopy.

To provide the best chance of wide adoption and engagement, we are developing the formats in the open, with frequent public announcements of progress and releases of reference software and examples (<https://forum.image.sc/tag/ome-ngff>) and regular community meetings where we present work, source feedback and encourage community members, including vendors, to participate

in the specification and implementation. The community process is being developed and we welcome contributions from all interested parties on <https://github.com/ome/ngff>.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41592-021-01326-w>.

Received: 14 April 2021; Accepted: 19 October 2021;
Published online: 29 November 2021

References

1. Wilkinson, M. D. et al. The FAIR guiding principles for scientific data management and stewardship. *Sci. Data* **3**, 160018 (2016).
2. Ellenberg, J. et al. A call for public archives for biological image data. *Nat. Methods* **15**, 849–854 (2018).
3. Linkert, M. et al. Metadata matters: access to image data in the real world. *J. Cell Biol.* **189**, 777–782 (2010).
4. *The HDF5 Library and File Format* (The HDF Group, accessed 18 October 2021); <https://www.hdfgroup.org/solutions/hdf5/>
5. Miles, A. et al. `zarr-developers/zarr-python`: v.2.5.0 <https://doi.org/10.5281/zenodo.4069231> (2020).
6. Allan, C. et al. OMERO: flexible, model-driven data management for experimental biology. *Nat. Methods* **9**, 245–253 (2012).
7. Sofroniew, N. et al. `napari/napari`: 0.4.7rc1 <https://doi.org/10.5281/zenodo.3555620> (2021).
8. Vergara, H. M. et al. Whole-body integration of gene expression and single-cell morphology. *Cell* <https://doi.org/10.1016/j.cell.2021.07.017> (2021).
9. Manz, T. et al. Viv: multiscale visualization of high-resolution multiplexed bioimaging data on the web. *OSF Preprints* <https://doi.org/10.31219/osf.io/wd2gu> (2020).
10. Williams, E. et al. The Image Data Resource: a bioimage data integration and publication platform. *Nat. Methods* **14**, 775–781 (2017).
11. Lin, J.-R. et al. Highly multiplexed immunofluorescence imaging of human tissues and tumors using t-CyCIF and conventional optical microscopes. *eLife* **7**, e31657 (2018).
12. Wan, Y., McDole, K. & Keller, P. J. Light-sheet microscopy and its potential for understanding developmental processes. *Annu. Rev. Cell Dev. Biol.* **35**, 655–681 (2019).
13. Kang, D., Rübner, O., Byna, S. & Blanas, S. Predicting and comparing the performance of array management libraries. In *2020 IEEE International Parallel and Distributed Processing Symposium* <https://doi.org/10.1109/IPDPS47924.2020.00097> (2020).
14. Abernathy, R. et al. Cloud-native repositories for big scientific data. *Comput. Sci. Eng.* <https://doi.org/10.1109/MCSE.2021.3059437> (2021).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021

Methods

Bioimage latency benchmark: synthetic data generation. *Imaging modality and dataset sizes.* Synthetic datasets were generated for two established imaging modalities: a large multi-channel two-dimensional image typical of CycIF¹¹ of *xyzct* dimensions $64,000 \times 64,000 \times 1 \times 8 \times 1$ and a time-lapse isotropic volume typical of LSM¹² of *xyzct* dimensions $1,024 \times 1,024 \times 1,024 \times 1 \times 100$.

For each modality, the chunk size of the benchmark dataset was chosen as a compromise between the size of individual chunks and the total number of chunks in the Zarr dataset. To make this decision, the individual chunk size was computed against the total number of chunks for typical sizes ranging from 16 up to 1,024 (chunks.py¹³ and Extended Data Fig. 2). Based on this, we chose a 2D chunk size of 256×256 for the CycIF-like dataset and a 3D chunk size of $32 \times 32 \times 32$ for the LSM-like dataset. Note that owing to the planar limitation of TIFF, the LSM dataset was stored as 2D TIFF tiles of size 32×32 but the benchmark loaded 32 tiles to measure the total access time. All data was stored uncompressed to keep chunk sizes consistent for the random generated data. Note that with the default `aws s3 cp` command, data upload decreased from over 100 MiB s^{-1} for the single HDF5 file to under 20 MiB s^{-1} for the Zarr dataset.

Dataset generation. The HDF5 version of each synthetic dataset was first generated by using the ImapisWriter library¹⁶ (v.2021-04-07) with a version of the ImapisWriterTest example^{16,17} modified to allow setting the desired chunk size and generate gradient images rather than random data. This HDF5-based Imapis file was converted into Zarr using a modified version of `bioformats2raw v.0.2.6` with support for chunks using a / dimension separator¹⁸. Finally, the Zarr was converted into TIFF with a modified version of `raw2ometiff v.0.2.6`, allowing it to consume Zarr filesets with a / dimension separator¹⁹. Both modifications have been released since in `bioformats2raw v.0.3.0` and `raw2ometiff v.0.3.0`.

For the CycIF-like dataset, this conversion generated a single 86 GB TIFF file, a single 86 GB HDF5 file and a Zarr dataset composed of 700,000 files of 86 GB in total. For the LSM-like dataset, the conversion generated a single 300 GB TIFF file, a single 229 GB HDF5 file and a Zarr dataset of 4.3 million files of 264 GB in total.

Bioimage latency benchmark: measurements and results. *Measurements.*

All three datasets along with a 1 byte dummy file for measuring overhead were placed in three types of storage: local disk, a remote server and object storage. We measured the reading time of individual chunks for all four file types across the three storage systems.

A random sequence of 100 chunk locations was chosen for the benchmark. All 100 chunks were loaded from each file in the same order. The time taken to retrieve the chunk, independent of the time taken to open a file or prepare the remote connection, was recorded.

Raincloud plots. Raincloud plots²⁰ combine three representations (split-half violin plots, box plots, raw data points) so that the true distribution and the statistical parameters can be compared. Split-half violin plots show a smoothed version of a histogram with a kernel density estimate. This type of plot is useful to determine, at a glance, if the mean is lower or higher than the median depending on the skewness of the curve. Box plots show the median and the boundaries of quartiles on either side of the median of the distribution to determine statistical differences at a glance. Below each box plot, the raw data points are additionally plotted with slight vertical jittering to avoid overlaps.

All code for reproducing the plots and the runs both locally with Docker or Amazon EC2 instances are available under a BSD-2 license on Zenodo¹⁵.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

The synthetic data generated for the benchmark are 1.05 TB. All code necessary to regenerate the data, including at different sizes, is available on Zenodo under a BSD-2 license¹⁵. The SARS-CoV-2 EM dataset from Extended Data Fig. 1, originally from Lamers et al.²¹ and published in IDR²², was converted into OME-NGFF and is available at Zenodo²³ under a CC-BY 4.0 license.

Code availability

Data generation and analysis code for file format benchmarking is available on Zenodo under a BSD-2 license¹⁵.

References

- Moore, J. et al. [ome/bioimage-latency-benchmark: 2021-10-05](https://doi.org/10.5281/zenodo.4668606). <https://doi.org/10.5281/zenodo.4668606> (2021).
- Beati, I., Andreica, E. & Majer, P. ImapisWriter: open source software for storage of large images in blockwise multi-resolution format. Preprint at *arXiv* <https://arxiv.org/abs/2008.10311> (2020).
- Gault, D. [ome/ImarisWriterTest](https://doi.org/10.5281/zenodo.5547849). <https://doi.org/10.5281/zenodo.5547849> (2021).
- Allan, C. et al. [ome/bioformats2raw](https://doi.org/10.5281/zenodo.5548102). <https://doi.org/10.5281/zenodo.5548102> (2021).
- Allan, C., Linkert, M. & Moore, J. [ome/raw2ometiff](https://doi.org/10.5281/zenodo.5548109). <https://doi.org/10.5281/zenodo.5548109> (2021).
- Allen, M. et al. Raincloud plots: a multi-platform tool for robust data visualization. *Wellcome Open Res.* **4**, 63 (2021).
- Lamers, M. M. et al. SARS-CoV-2 productively infects human gut enterocytes. *Science* **369**, 50–54 (2020).
- Lamers, M. M. et al. SARS-CoV-2 productively infects human gut enterocytes. *Image Data Resource* <https://doi.org/10.17867/10000135> (2020).
- Moore, J. & Besson, S. OME-NGFF: EM image of SARS-CoV-2. <https://doi.org/10.5281/zenodo.4668606> (2020).

Acknowledgements

Work on the `bioformats2raw` and `raw2ometiff` converters was funded by awards from InnovateUK to Glencoe Software (PathLAKE, ref. 104689 and iCAIRD ref. 104690). Work on OME-NGFF by J.M., J.-M.B., S.B., D.G., D.L. and W.M. was funded by grant no. 2019-207272 from the Chan Zuckerberg Initiative DAF, an advised fund of Silicon Valley Community Foundation, the Wellcome Trust (ref. 212962/Z/18/Z) and BBSRC (ref. BB/R015384/1). Work by C.T. has been made possible in part by grant number 2020-225265 from the Chan Zuckerberg Initiative DAF, an advised fund of Silicon Valley Community Foundation. C.P. has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement no. 824087. T.M. has received funding from the National Science Foundation Graduate Research Fellowship under grant no. DGE1745303. The authors thank the originators of the Zarr and N5 formats, A. Miles and S. Saalfeld and the vibrant communities they have built for working together to unify their formats. The development of OME-NGFF has been and will continue to be a community endeavor. Everyone who has participated in the format specification and/or an implementation is invited to request software authorship (<http://credit.niso.org/contributor-roles/software/>) by contacting the corresponding author.

Author contributions

J.M., C.A., J.-M.B. and S.B. conceived the project; J.M. and S.B. wrote the specification; C.A., M.L. and E.D. wrote and tested the conversion software; D.G. wrote the data generation code; K.K. wrote the AWS-deployment scripts; J.M. performed the benchmark and created the figures; T.M., W.M., J.-M.B., C.P. and C.T. wrote and validated software tools to visualize data; D.L., S.B. and W.M. tested the formats for data publishing; J.-M.B. validated the formats using workflows with public data; J.R.S. acquired the funding; and J.M., S.B. and J.R.S. wrote the paper with input from all the authors.

Competing interests

C.A., E.D., K.K., M.L. and J.R.S. are affiliated with Glencoe Software, a commercial company that builds, delivers, supports and integrates image data management systems across academic, biotech and pharmaceutical industries. The remaining authors declare no competing interests.

Additional information

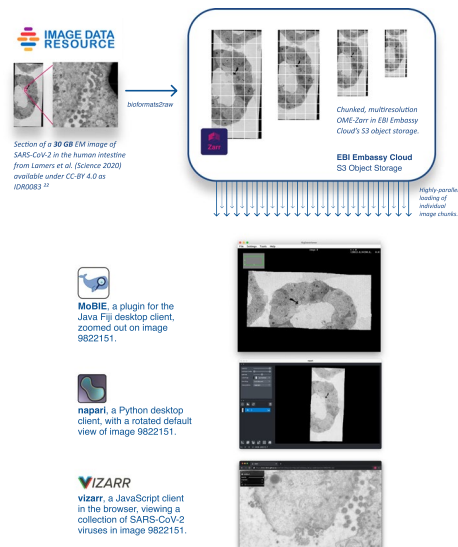
Extended data is available for this paper at <https://doi.org/10.1038/s41592-021-01326-w>.

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41592-021-01326-w>.

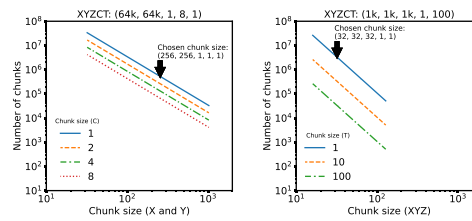
Correspondence and requests for materials should be addressed to Jason R. Swedlow.

Peer review information *Nature Methods* thanks Albert Cardona and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Rita Strack was the primary editor on this article and managed its editorial process and peer review in collaboration with the rest of the editorial team

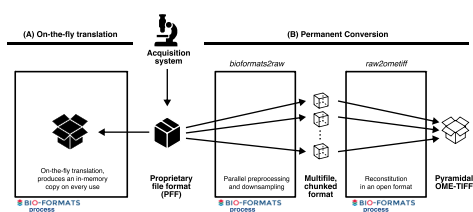
Reprints and permissions information is available at www.nature.com/reprints.



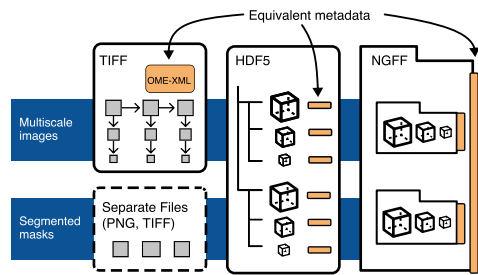
Extended Data Fig. 1 | Maximizing re-use by allowing popular tools to access bioimaging data in the cloud. An example of using NGFFs for promoting the distribution of public image datasets. Selection of current tools streaming different portions of the same SARS-CoV-2 virus image at various resolutions directly from S3 storage at the European Bioinformatics Institute (EBI). Original data from Lamers et al. is available in IDR while the converted data is available on Zenodo.21-23.



Extended Data Fig. 2 | Effect of Chunk Size on Chunk Number. For each modality, the chunk size of the benchmark dataset was chosen as a compromise between the size of individual chunks and the total number of chunks in the Zarr dataset. The plots above show typical power of 2 chunk sizes: between 32 and 1024 for the 2D data and between 16 and 128 for the 3D data. We chose a 2D chunk size of 256×256 for the CyIF-like dataset and a 3D chunk size of 32×32×32 for the LSM-like dataset. Note that due to the planar limitation of TIFF, the LSM dataset was stored as 2D TIFF tiles of size 32×32 but the benchmark looped over 32 tiles to measure the access time of the same chunk size.



Extended Data Fig. 3 | Conversion tools provide an alternative to continual, on-the-fly translation of PFFs. Figure shows workflows for file format access. **(a)** The classical approach to access images produced by an acquisition system is to use a library like Bio-Formats to translate the proprietary file format (PFF) and produce an in-memory copy of the imaging data on-the-fly. This translation needs to be repeated on every use. **(b)** With the existence of open, community-supported formats, converting PFFs becomes the most cost-efficient method for long-term storage and sharing of microscopy data. `bioformats2raw` and `raw2ometiff` (Supplementary Note) parallelize the creation of an open format, OME-TIFF, by using an intermediate NGFF format consisting of many, individual files each with one chunk of the original image data.



Extended Data Fig. 4 | Unification of metadata specifications will allow interoperability between TIFF, HDF5, and Zarr. Each proposed container (TIFF, Zarr, HDF5) can be used interchangeably to store pixel data, but trade-offs described in this manuscript can be used to determine what is the best target. TIFF is ideal for interoperability in digital pathology and other 2-dimensional domains since the format is widely accessible by established open source and proprietary software. In higher-dimensional domains, HDF5 and Zarr are better suited. HDF5 will likely be preferred for local access. If data is intended for sharing in the cloud, Zarr will likely be preferred. High throughput image analysis will benefit from the lower-latency access to data in HDF5 and Zarr. If original image data is paired with derived representations like pixel or object classification, a shared structure in HDF5 or Zarr is likely the best choice.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- | | |
|-----|-----------|
| n/a | Confirmed |
|-----|-----------|
- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
 - A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
 - The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
 - A description of all covariates tested
 - A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
 - A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
 - For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
 - For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
 - For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
 - Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The synthetic data generated for the benchmark is 1.05 TB. All code necessary to regenerate the data, including at different sizes, is available in <https://github.com/ome/bioimage-latency-benchmark> under a BSD-2 license.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

- Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	In Extended Figure 2, The line plots show how the number of chunks needed for an image of the specified dimensions varies with size of the chunk. Chunk sizes were chosen here to balance chunk size with number of files. In Figure 1, the access latency was sampled 100 times for each condition.
Data exclusions	No data were excluded
Replication	N/A
Randomization	Chunks to access in file were chosen at random.
Blinding	Blinding of data was not performed as no human intervention in measurement or interpretation was part of the study

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

- | n/a | Involved in the study |
|-------------------------------------|--|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Antibodies |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Eukaryotic cell lines |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Palaeontology and archaeology |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Animals and other organisms |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Human research participants |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Clinical data |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Dual use research of concern |

Methods

- | n/a | Involved in the study |
|-------------------------------------|---|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> ChIP-seq |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Flow cytometry |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> MRI-based neuroimaging |