



HHS Public Access

Author manuscript

Nat Methods. Author manuscript; available in PMC 2022 April 14.

Published in final edited form as:

Nat Methods. 2022 March ; 19(3): 262–267. doi:10.1038/s41592-022-01415-4.

*These authors contributed equally

**Co-Corresponding Authors

The HTAN Consortium

Daniel L. Abравane²⁸, Samuel Achilefu¹⁸, Foluso O. Ademuyiwa¹⁸, Andrew C. Adey¹¹, Rebecca Aft¹⁸, Khung Jun Ahn³⁸, Fatemeh Alikarami³⁸, Shahar Alon³³, Orr Ashenberg², Ethan Baker², Gregory J. Baker¹, Shovik Bandyopadhyay³⁶, Peter Bayguinov¹⁸, Jennifer Beane²³, Winston Becker¹⁹, Kathrin Bernt³⁸, Courtney B. Betts¹¹, Julie Bletz¹⁰, Tim Blosser³², Adrienne Boire⁹, Genevieve M. Boland²⁸, Edward S. Boyden³⁴, Elmar Bucher¹¹, Raphael Bueno²⁸, Qiuyin Cai¹⁷, Francesco Cambuli⁹, Joshua Campbell²³, Song Cao¹⁸, Wagma Caravan¹⁸, Ronan Chaligne⁹, Joseph Chan⁹, Sara Chasnoff¹⁸, Deyali Chatterjee¹⁸, Alyce A. Chen¹, Changya Chen³⁸, Chia-hui Chen³⁸, Bob Chen¹⁷, Feng Chen¹⁸, Siqi Chen¹⁸, Milan G. Chheda¹⁸, Koei Chin¹¹, Hyeyoung Cho¹¹, Jaeyoung Chun⁹, Luis Cisneros³⁰, Robert J. Coffey¹⁷, Ofir Cohen², Graham A. Colditz¹⁸, Kristina A. Cole³⁸, Natalie Collins¹², Daniel Cotter¹⁹, Lisa M. Coussens¹¹, Shannon Coy²⁴, Allison L. Creason¹¹, Yi Cui³⁴, Daniel Cui Zhou¹⁸, Christina Curtis¹⁹, Sherri R. Davies¹⁸, Inode Bruijn⁹, Toni M. Delorey², Emek Demir¹¹, David Denardo¹⁸, Dinh Diep⁴⁰, Li Ding¹⁸, John DiPersio¹⁸, Steven M. Dubinett²⁹, Timothy J. Eberlein¹⁸, James A. Eddy¹⁰, Edward D. Esplin¹⁹, Rachel E. Factor³⁰, Kayvon Fatahalian¹⁹, Heidi S. Feiler¹¹, Jose Fernandez¹⁹, Andrew Fields¹¹, Ryan C. Fields¹⁸, James A.J. Fitzpatrick¹⁸, James M. Ford¹⁹, Jeff Franklin¹⁷, Bob Fulton¹⁸, Giorgio Gallia²⁴, Luciano Galdieri¹⁸, Karuna Ganesh⁹, Jianjiong Gao⁹, Benjamin L. Gaudio¹, Gad Getz², David L. Gibbs⁷, William E. Gillanders¹⁸, Jeremy Goecks¹¹, Daniel Goodwin³⁴, Joe W. Gray¹¹, William Greenleaf¹⁹, Lars J. Grimm³⁰, Qiang Gu¹¹, Jennifer L. Guerriero¹², Tuhin Guha¹⁹, Alexander R. Guimaraes¹¹, Belen Gutierrez¹⁹, Nir Hacohen³³, Casey Ryan Hanson¹⁹, Coleman R. Harris¹⁷, William G. Hawkins¹⁸, Cody N. Heiser¹⁷, John Hoffer¹, Travis J. Hollmann⁹, James J. Hsieh¹⁸, Jeffrey Huang³⁸, Stephen P. Hunger³⁸, Eun-Sil Hwang³⁰, Christine Iacobuzio-Donahue⁹, Michael D. Iglesias¹⁸, Mohammad Islam³⁸, Benjamin Izar¹², Connor A. Jacobson¹, Samuel Janes³⁹, Reyka G. Jayasinghe¹⁸, Tiarah Jeudi¹², Bruce E. Johnson¹², Brett E. Johnson¹¹, Tao Ju¹⁸, Humam Kadara³⁵, Elias-Ramzey Karnoub⁹, Alla Karpova¹⁸, Aziz Khan¹⁹, Warren Kibbe³⁰, Albert H. Kim¹⁸, Lorraine M. King³⁰, Elyse Kozlowski¹², Praveen Krishnamoorthy¹⁸, Robert Krueger³², Anshul Kundaje¹⁹, Uri Ladabaum¹⁹, Rozelle Laquindanum¹⁹, Clarisse Lau⁷, Ken Siu Kwong Lau¹⁷, Nicole R. LeBoeuf²⁸, Hayan Lee¹⁹, Marc Lenburg²³, Ignaty Leshchiner², Rochelle Levy¹², Yize Li¹⁸, Christine G. Lian²⁴, Wen-Wen Liang¹⁸, Kian-Huat Lim¹⁸, Yiyun Lin³⁵, David Liu¹², Qi Liu¹⁷, Ruiyang Liu¹⁸, Joseph Lo³⁰, Pierrette Lo¹⁰, William J. Longabaugh⁷, Teri Longacre¹⁹, Katie Lockett⁹, Cynthia Ma¹⁸, Chris Maher¹⁸, Allison Maier³¹, Danika Makowski³⁸, Carlo Maley³⁰, Zoltan Maliga¹, Parvathy Manoj⁹, John M. Maris³⁸, Nick Markham¹⁷, Jeffrey R. Marks³⁰, Daniel Martinez³⁸, Jay Mash¹⁸, Ignas Masilionis⁹, Joan Massage⁹, Marceij A. Mazurowski³⁰, Eliot T. McKinley¹⁷, Joshua McMichael¹⁸, Matthew Meyerson¹², Gordon B. Mills¹¹, Zahi I. Mitrani¹¹, Andrew Moorman⁹, Jacqueline Mudd¹⁸, George F. Murphy²⁴, Nataly Naser AL Deen¹⁸, Nicholas E. Navin³⁵, Tal Naway⁹, Reid M. Ness¹⁷, Stephanie Nevins¹⁹, Ajit Johnson Nirmal¹², Edward Novikov¹, Stephen T. Oh¹⁸, Derek A. Oldridge³⁶, Kourou Owzar³⁰, Shishir M. Pant³¹, Wungki Park⁹, Gary J. Patti¹⁸, Kristina Paul¹⁹, Roxanne Pelletier¹, Daniel Persson¹¹, Candi Petty¹⁸, Hanspeter Pfister³², Kornelia Polyak¹², Sidharth V. Puram¹⁸, Qi Qiu³⁶, Alvaro Quintanal Villalonga⁹, Marisol Adelina Ramirez¹⁷, Rumana Rashid²⁴, Ashley N. Reeb¹⁸, Mary E. Reid³⁷, Jan Remsik⁹, Jessica L. Riesterer¹¹, Tyler Risom¹⁹, Cecily Claire Ritch²⁴, Andrea Rolong¹⁷, Charles M. Rudin⁹, Marc D. Ryser³⁰, Kazuhito Sato¹⁸, Cynthia L. Sears¹⁷, Yevgeniy R. Semenov³³, Jeanne Shen¹⁹, Koresh I. Shoghi¹⁸, Martha J. Shrubsole¹⁷, Yu Shyr¹⁷, Alexander B. Sibley³⁰, Alan J. Simmons¹⁷, Anubhav Sinha³⁴, Shamilene Sivagnanam¹¹, Sheng-Kwei Song¹⁸, Austin Southar-Smith¹⁸, Avrum E. Spira²³, Jeremy St. Cyr¹², Stephanie Stefankiewicz³⁸, Erik P. Storrs¹⁸, Elizabeth H. Stover¹², Siri H. Strand¹⁹, Cody Straub³⁰, Cherease Street¹⁸, Timothy Su¹⁷, Lea F. Surrey³⁸, Christine Suver¹⁰, Kai Tan³⁸, Nadezhda V. Terekhanova¹⁸, Luke Ternes¹¹, Anusha Thadi³⁸, George Thomas¹¹, Rob Tibshirani¹⁹, Shigeaki Umeda⁹, Yasin Uzun³⁸, Tuulia Vallius¹, Eliezer R. Van Allen¹², Simon Vandekar¹⁷, Paige N. Vega¹⁷, Deborah J. Veis¹⁸, Sujay Vennam¹⁹, Ana Verma²⁴, Sebastien Vigneau¹², Nikhil Wagle¹², Richard Wahl¹⁸, Thomas Walle⁹, Liang-Bo Wang¹⁸, Simon Warchol³², M. Kay Washington¹⁷, Cameron Watson¹¹, Annika K. Weimer¹⁹, Michael C. Wendl¹⁸, Robert B. West¹⁹, Shannon White¹⁹, Annika L. Windon¹⁷, Hao Wu³⁶, Chi-Yun Wu³⁶, Yige Wu¹⁸, Matthew A. Wyczalkowski¹⁸, Jason Xu³⁶, Lijun Yao¹⁸, Wenbao Yu³⁸, Kun Zhang⁴⁰, Xiangzhu Zhu¹⁷

- ²⁷ Arizona State University, Tempe, AZ, USA
- ²⁸ Brigham and Women's Hospital, Boston, MA, USA
- ²⁹ David Geffen School of Medicine at UCLA, Los Angeles, CA, USA
- ³⁰ Duke University Medical Center, Durham, NC, USA
- ³¹ Harvard Medical School, Boston, MA, USA
- ³² Harvard University, Cambridge, MA, USA
- ³³ Massachusetts General Hospital, Boston, MA, USA
- ³⁴ Massachusetts Institute of Technology, Cambridge, MA, USA
- ³⁵ MD Anderson Cancer Center, Houston, TX, USA
- ³⁶ Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA
- ³⁷ Roswell Park Comprehensive Cancer Center, Buffalo, NY, USA
- ³⁸ The Children's Hospital of Philadelphia, Philadelphia, PA, USA
- ³⁹ Division of Medicine, University College London, London, UK
- ⁴⁰ University of California San Diego, San Diego, CA, USA

Author contributions

D.S., C.Y and A.S. initiated and implemented the MITI guidelines with extensive guidance from other authors and direct supervision by P.K.S. and S.S.. All authors contributed to and reviewed the final MITI guidelines. D.S., C.Y, A.S., P.K.S and S.S. wrote the manuscript with input from all authors.

MITI Minimum Information guidelines for highly multiplexed tissue images

Denis Schapiro^{1,2,3,4,*}, **Clarence Yapp**^{1,5,*}, **Artem Sokolov**^{1,6,*}, **Sheila M. Reynolds**⁷, **Yu-An Chen**¹, **Damir Sudar**⁸, **Yubin Xie**⁹, **Jeremy Muhlich**¹, **Raquel Arias-Camison**¹, **Sarah Arena**¹, **Adam J. Taylor**¹⁰, **Milen Nikolov**¹⁰, **Madison Tyler**¹, **Jia-Ren Lin**¹, **Erik A. Burlingame**^{11,26}, **Human Tumor Atlas Network**, **Young H. Chang**¹¹, **Samouil L Farhi**², **Vésteinn Thorsson**⁷, **Nithya Venkatamohan**¹², **Julia L. Drewes**¹³, **Dana Pe'er**⁹, **David A. Gutman**¹⁴, **Markus D. Herrmann**¹⁵, **Nils Gehlenborg**⁶, **Peter Bankhead**¹⁶, **Joseph T. Roland**¹⁷, **John M. Herndon**¹⁸, **Michael P. Snyder**¹⁹, **Michael Angelo**¹⁹, **Garry Nolan**¹⁹, **Jason R. Swedlow**²⁰, **Nikolaus Schultz**²¹, **Daniel T. Merrick**²², **Sarah A. Mazzili**²³, **Ethan Cerami**¹², **Scott J. Rodig**²⁴, **Sandro Santagata**^{1,24,**}, **Peter K. Sorger**^{1,25,**}

¹Laboratory of Systems Pharmacology, Ludwig Center for Cancer Research at Harvard, Harvard Medical School, Boston, MA, USA

²Klarman Cell Observatory, Broad Institute of MIT and Harvard, Cambridge, MA, USA

³Institute for Computational Biomedicine, Faculty of Medicine, Heidelberg University Hospital and Heidelberg University, Heidelberg, Germany.

⁴Institute of Pathology, Heidelberg University Hospital, Heidelberg, Germany.

⁵Image and Data Analysis Core, Harvard Medical School, Boston, MA, USA

⁶Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

⁷Institute for Systems Biology, Seattle, WA, USA

⁸Quantitative Imaging Systems LLC, Portland, OR, USA

⁹Program in Computational and Systems Biology, Memorial Sloan Kettering Cancer Center, New York, NY, USA

¹⁰Sage Bionetworks, Seattle, WA, USA

¹¹Oregon Health and Science University, Portland, OR, USA

¹²Dana-Farber Cancer Institute, Boston, MA, USA.

¹³Johns Hopkins University School of Medicine, Baltimore, MD, USA

¹⁴School of Medicine, Emory University, Atlanta, GA, USA

¹⁵Department of Pathology, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA

¹⁶Edinburgh Pathology, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh, UK.

¹⁷Vanderbilt University School of Medicine, Nashville TN, USA

¹⁸Department of Surgery, Washington University School of Medicine, St. Louis, MO USA.

¹⁹School of Medicine, Stanford University, Stanford, CA, USA

²⁰Division of Computational Biology and Centre for Gene Regulation and Expression, University of Dundee, Dundee, UK

²¹Department of Epidemiology & Biostatistics at Memorial Sloan Kettering Cancer Center, New York, NY, USA

²²Pathology, University of Colorado, Aurora, CO, USA

²³Boston University, Boston, MA, USA

²⁴Department of Pathology, Brigham and Women's Hospital, Boston, MA, USA

²⁵Department of Systems Biology, Harvard Medical School, Boston, MA, USA

²⁶Current Affiliation: Indica Labs, Albuquerque, NM, USA

Abstract

The imminent release of tissue atlases combining multi-channel microscopy with single cell sequencing and other omics data from normal and diseased specimens creates an urgent need for data and metadata standards that guide data deposition, curation and release. We describe a Minimum Information about highly multiplexed Tissue Imaging (MITI) standard that applies best practices developed for genomics and other microscopy data to highly multiplexed tissue images and traditional histology.

Highly multiplexed tissue imaging using any of a variety of optical and mass-spectrometry based methods (Supplemental Table 1) combines deep molecular insight into the biology of single cells with spatial information traditionally acquired using histological methods, such as hematoxylin and eosin (H&E) staining and immunohistochemistry (IHC)¹. As currently practiced, multiplexed tissue imaging of proteins involves 20–60 channels of 2D data, with each channel corresponding to a different antibody or colorimetric stain (Figure 1). Multiple inter-institutional and international projects, such as the Human Tumor Atlas Network (HTAN)², the Human BioMolecular Atlas Program (HuBMAP)³, and the LifeTime Initiative⁴ aim to combine such highly multiplexed tissue images with single cell sequencing and other types of omics data to create publicly accessible “atlases” of normal and diseased tissues. Easy public access to primary and derived data is an explicit goal of these atlases and is expected to encompass native-resolution images, segmented single-cell data, anonymized clinical metadata and treatment history (for human specimens), genetic information (particularly for animal models), and specification of the protocols used to acquire and process the data. Given the imminent release of the first atlases, an urgent need exists for data and metadata standards consistent with emerging FAIR (Findable, Accessible, Interoperable, and Reusable) standards⁵. In this commentary, we establish the MITI (Minimum Information about highly multiplexed Tissue Imaging) standard and associated data level definitions; we also discuss the relationship of MITI to existing standards, practical implementations, and future developments.

Scope and target audiences

MITI covers biospecimen, reagent, data acquisition and data analysis metadata, as well as data levels for imaging with antibodies, aptamers, peptides, dyes and similar detection reagents (Supplemental Table 1). The standard is also compatible with images based on H&E staining, low-plex immunofluorescence (IF) and IHC. A working group is currently extending MITI to cover subcellular resolution imaging of nucleic acids using methods such as MERFISH⁶. While conceived with today's two-dimensional (2D) images in mind (these typically involve 5 – 10 µm thick sections of fixed or frozen specimens), MITI accommodates three-dimensional (3D) datasets acquired using confocal, deconvolution and light sheet microscopes⁷. MITI has been established as its own organization with its own GitHub repository, governing structure, and procedures for proposing and incorporating revisions. The definition of MITI is available in machine readable YAML format (<https://github.com/miti-consortium/MITI>) along with other relevant information. MITI has also been implemented in practice (<https://github.com/ncihtan/data-models>) and used to structure metadata available via the HTAN data portal (<https://htan-portal-nextjs.vercel.app>). However, MITI is independent of HTAN or any single research consortium.

Highly multiplexed imaging is derived from methods such as IHC and IF that are in widespread use in pre-clinical research using cultured cells and model organisms, and in clinical practice with human tissue specimens. Many standards and best practices have been established for these types of data (Supplemental Table 2), but high-plex imaging presents unique challenges: images are expensive to collect and can be very large (up to 1TB in size), specimens are often difficult to acquire and may have data use restrictions, and accurate clinical and genomic annotation is a necessity. Recent interest in highly multiplexed tissue imaging has been driven by applications in oncology, largely due to the importance of the tumor microenvironment in immuno-editing and responsiveness to immunotherapy, but the approach is broadly applicable to studying normal development, infectious disease, immunology and other topics. HuBMAP³, for example, is using high-plex imaging to study a range of normal human tissues. MITI is also relevant to studies with model organisms and data tables have already been created to store data from genetically engineered mouse models (GEMMs) in a standardized manner.

Multiplexed imaging also promises to impact the pathological diagnosis of diseases, which is rapidly switching to digital approaches⁸. For over a century, histological analysis of anatomic specimens (from biopsies and surgical resection) has been the primary method of diagnosing diseases such as cancer⁹, and this remains true today, despite the impact of gene sequencing. Multiplexed tissue imaging promises to augment conventional pathological diagnosis with the detailed molecular information needed to specify use of contemporary precision therapies. This is therefore an opportune time to seek alignment of research and diagnostic approaches by establishing public standards able to take full advantage of the detailed molecular information revealed by emerging imaging methods.

Existing standards and approaches

The Human Genome Project, the Cancer Genome Atlas (TCGA)¹⁰ and similar large-scale genomic programs have developed several approaches to data management of immediate relevance to tissue atlases. The first is the concept of “minimum information” metadata, which has been employed in microarrays (the MIAME standard)¹¹, genome sequences (MIGS)¹², and biological investigation in general (MIBBI)¹³. The second is the idea of “data levels” (<https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/data-levels>), which specify the extent of data processing (raw, normalized, aggregated or region of interest, corresponding to data levels 1–4) and access control. Access control is required because even anonymized DNA sequencing data pose a re-identification risk¹⁴. As a result, the database of Genotypes and Phenotypes (dbGaP), the NCI Genome Data Commons (GDC)¹⁵, and the US Federal Register (79 FR 51345) control access to primary sequencing data (so-called level 1 & 2 sequencing data) based on policies set by a data access committee. Higher level genomic data, which are generally more consolidated, involve information aggregated from many patients, and pose little or no re-identification risk can be freely shared¹⁶ (Figure 2). When datasets are combined, they acquire the most stringent restriction applied to any constituent element. While we are not aware of any policies addressing the anonymity of histological images, consultation with our Institutional Review Boards (IRBs, ethics committees) has led us to conclude that public release of tissue images does not constitute a risk to patient privacy. MITI data levels are nonetheless consistent with the existing GDC and dbGaP practice that data intended for unrestricted distribution are classified as level 3 and up. In the case of images adhering to the MITI standard, level 3 data have been subjected to quality control and some degree of human annotation, making them more useful in a shared environment than raw images. We anticipate that IRBs and government agencies will in the future provide further guidance on sharing of datasets that combine clinical history, sequence information, and tissue images; MITI will be adapted to accommodate such guidance.

The MITI standard also draws extensively on image formats developed for cultured cells and model organisms and on a wide variety of open-source software tools (Supplemental Table 3). Noteworthy among these are the Open Microscopy Environment (OME) TIFF standard¹⁷ and the BioFormats¹⁸ approach to standardization of microscopy data. MITI field definitions are harmonized with the QUality Assessment and REProducibility for Instruments and Images in Light Microscopy (QUAREP-LiMi)¹⁹ effort, the Resource Identification Initiative²⁰, and antibody standardization efforts by the Human Protein Atlas²¹ and are also compliant with the recently developed Recommended Metadata for Biological Images initiative²². Metadata on model organisms (particularly GEMMs - and patient derived xenografts - PDXs) are aligned with existing standards, many developed for genomic information (see Supplemental Table 2 for a full list of antecedent resources). Well-curated clinical information is essential for the interpretation of data from human specimens but standardizing such information has proven to be a major challenge in the past, for example in TCGA²³. Thus, HTAN and other current NCI projects focused on human specimens are emphasizing standardization of clinical metadata, and the MITI standard is

designed to closely align with the Genomic Data Commons (GDC) Data Model²⁴ in this regard (Supplemental Tables 5–6).

All imaging methods generate data that comprises a sequence of intensity values on a raster; multi-spectral imaging simply adds new dimensions to the raster. The cameras that collect H&E and IHC images from bright-field microscopes or high-plex images from fluorescence microscopes generate a raster; ablation-based mass-spectrometry imaging (e.g. MIBI and IMC) is also raster based. As currently defined, MITI specifies that raster images should be stored in the OME-TIFF 6 standard, but OME formats are currently being migrated to a set of next generation file formats (collectively OME-NGFF)²⁵ to improve scalability and performance on the cloud. MITI will be updated to align with these new formats as they come into general use. Another area of translational and clinical research in which imaging is commonly encountered is radiology, which is almost entirely digital, and uses data interchange standards governed by DICOM (<https://www.dicomstandard.org/>). DICOM has recently been extended to accommodate both radiology data and OME-TIFF standards²⁶. The NCI's ongoing program to create an Imaging Data Commons²⁷ is expected to be based on this dual standard, or on a successor using OME-NGFF. MITI is, or will be, compatible with these foundational data standards.

In highly multiplexed tissue imaging antibodies are either conjugated to fluorophores directly or via oligonucleotides, or are bound to secondary antibodies (Figure 1, Supplemental Table 4). Images are then acquired serially, one to six channels at a time, to assemble data from 20–60 antibodies. In ablation-based methods, antibodies are labelled with metals and vaporized with lasers or ion beams after which they are detected by atomic mass spectrometry (Supplementary Table 4). In all cases, the raw output of data acquisition instruments comprises Level 1 MITI data (Figure 2), analogous to the Level 1 FASTq files in genomics.

Whole slide imaging is required for clinical applications²⁸ and also necessary to ensure adequate power in pre-clinical studies²⁹. However, resolution and field of view have a reciprocal relationship – both with respect to optical physics and the practical process of mapping image fields onto the fixed raster of a camera (or ablating beam). Whole slide images of histological specimens⁸ must therefore be acquired by dividing a large specimen into contiguous tiles. This usually involves acquisition of ~100 to 1,000 tiles by moving the microscope stage in both X and Y, with each tile being a multi-dimensional, subcellular resolution, TIFF image. Tiles are combined at sub-pixel accuracy into a mosaic image in a process known as stitching. When high-plex images are assembled from multiple rounds of lower-plex imaging, it is also necessary to register channels to each other across imaging cycles and to correct for any unevenness in illumination (so-called flat fielding)³⁰. Stitched and registered mosaics can be as large as 50,000 × 50,000 pixels × 100 channels and require ~500 GB of disk space. They correspond to Level 2 MITI data and represent full-resolution primary images that have undergone automated stitching, registration, illumination correction, background subtraction, intensity normalization and have been stored in a standardized OME format. The level of processing is analogous to BAM files, a common type of Level 2 data in genomics.

Level 3 data represent images that have been processed with some interpretive intent, which may include (i) full-resolution images following quality control or artifact removal, (ii) segmentation masks computed from such images, (iii) machine-generated spatial models, and (iv) images with human or machine-generated annotations. Level 3 MITI data is roughly analogous to Level 3 mRNA expression data in genomics. However, whereas many users of genomic data only require access to processed level 3 and 4 data, which are usually quite compact, quantitative analysis of tissue images adds a requirement for full-resolution primary images so that images and computed features can be examined in parallel³¹. Level 3 MITI data is intended to be the primary type of image data distributed by tissue atlases and similar projects.

Assembled level 3 images are typically segmented to identify single cells³¹, which are quantified to produce a “spatial feature table” that describes marker intensities, cell coordinates and other single-cell features. The Level 4 data in spatial feature tables are a natural complement to count tables in single cell sequencing data (e.g. scRNA-seq, scATAC-seq, scDNA-seq) and can be analyzed using many of the same dimensionality reduction methods (e.g. PCA, t-SNE and U-MAP)³² and on-line browsers such as cellxgene (Supplemental Table 3)³³. These types of tabular data are all examples of “Feature Observation Matrixes” which are themselves being standardized across domains of biology to improve their utility and inter-compatibility. Level 5 MITI data comprise results computed from spatial feature tables or primary images. Because access to TB-size full-resolution image data is impractically burdensome when reading a manuscript or browsing a large dataset, a specialized type of Level 5 image data has been developed to enable panning and zooming across images using a standard web browser. In the case of Level 5 images viewed with MINERVA software, the aim is to exploit similar functionality and concepts as those in Google Maps or electronic museum guides³⁴. The inclusion of digital docents with images makes it possible to combine pan and zoom with guided narratives that greatly facilitate comprehension of complex datasets and promote new hypothesis generation³⁵.

For any metadata standard to be used, a balance must be struck between ease of data entry, which minimizes non-compliance by data generators, and level of detail, which must be sufficient for data retrieval, analysis, and publication in a reproducible manner. Moreover, specifying a metadata standard is separate from the essential task of developing a practical and reliable means for capturing information needed to ensure adherence to the standard. Two approaches have proven most effective in addressing this requirement. One, exemplified by OMeta³⁶, involves a relational database and web interface that data generators use to input necessary information in a controlled manner. Another approach, exemplified by MAGE-TAB³⁷, involves a standardized format for collecting metadata via a series of structured documents, which are then used to populate web pages and databases³⁸. As a practical test of MITI we have implemented the latter approach in a JSON schema (<https://github.com/ncihtan/data-models>) that also conforms to the design principles of [SCHEMA.org](https://schema.org). These principles focus on the creation, maintenance and promotion of schemas for structured data that is supported by major web search engines, thereby enhancing discoverability. In this TAB-like approach the MITI standard is exposed to data collectors as Google Sheets with dropdowns representing controlled vocabularies and highlighting required or optional elements; many fields are automatically validated upon

entry. These documents are ingested using SCHEMATIC (Schema Engine for Manifest Ingress and Curation; <https://github.com/Sage-Bionetworks/schematic>), automatically linked to primary imaging data, and stored as cloud assets. These implementations continue to evolve, and entirely different approaches are possible: nothing in a MITI-type standard constrains how data are collected.

Whereas many research agencies and countries have made a major investment in curating, storing, and distributing genomic data, fewer repositories exist for primary image data. The Image Data Resource³⁹ maintained by the European Bioinformatics Institute (EBI) is an exception, but as the volume of image data grows, other means of data distribution will almost certainly be required. In the U.S., in the absence of a major public investment in data storage, the development of “requester pays”⁴⁰ access to datasets is a promising development. The primary cost associated with creation and maintenance of a dataset on a commercial cloud service involves data download, not data ingress and storage. In a “requester pays” model, a user seeking access to a dataset pays the cost of data egress directly to the cloud provider making access both secure and anonymous (moreover, the cost of egress into another account on the same commercial cloud is low). Although the “requester pays” approach might appear to create an impediment to research, the actual cost of egress is quite low (currently, about hundred US dollars per TB) compared to any form of data acquisition and a key goal is to avoid a tragedy of the commons in which frequent, duplicate downloads overwhelm the system. A combination of a MITI implementation on a cloud service (as described above) with “requester pays” cloud access will also make it possible for individuals to distribute very large FAIR image datasets at relatively low cost. Such an approach does not obviate the need for public investments, such as those being made but EBI, but does represent a practical way forward to democratize release of standardized data – some of which can then be incorporated into publicly supported resources. Regardless, the MITI standard described here is available for immediate use, without being impacted by how access to the primary data is provisioned.

Public data and metadata standards have been essential for the success of genomics and other fields of biomedicine, but the creation of a new standard is no guarantee of successful adoption. An outpouring of effort 10–20 years ago led to the development of widely adopted and well maintained standards such as MIAME¹¹, MIGS¹² and MIBBI¹³, and these have been consolidated and further documented by the Digital Curation Center (<https://www.dcc.ac.uk/>), [FairSharing.org](https://www.fairsharing.org/), and similar projects. However many other minimum information projects have been left unattended⁴¹, and it remains unclear whether existing metadata adequately conform to user needs⁴². The development of MITI and of the initial HTAN implementation enjoys NCI support and is expected to become part of the NCI Cancer Research Data Commons²⁷, helping ensure its viability. However, individuals and organizations are invited to join in the further development of MITI and should make contact via the [image.sc](https://www.image.sc/) forum or submit pull requests (i.e. requests for inclusion in the MITI “code base” at <https://github.com/miti-consortium/MITI>). Because high high-plex tissue imaging is in its infancy and MITI has attracted the great majority of developers of existing high-plex tissue image acquisition methods, it represents a solid beginning for what will need to be an evolving standard. By having its own repository and governance structure,

independent of any particular research program or constituency, MITI also conforms with other requirements of successful open standards⁴³.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

This work is supported by the HTAN Consortium and the Cancer Systems Biology Consortium (CSBC). A list of all current Consortium members can be found at <https://humantumoratlas.org/>.

This work was supported by the following grants from the National Cancer Institute under the Human Tumor Atlas Network (HTAN) U2C CA233262 (Harvard Medical School), U2C CA233280 (OHSU), U2C CA233195 (Boston DFCI Broad), U2C CA233291 (Vanderbilt University Medical Center), U2C CA233311 (Stanford University), U2C CA233238 (Boston University Medical Campus), U2C CA233285 (Children's Hospital of Philadelphia), U2C CA233303 (Washington University St. Louis), U2C CA233280 (Oregon Health and Science University), U2C CA233284 (Memorial Sloan Kettering Cancer Center), U2C CA233254 (Duke University Medical Center) and by other public support including U54 CA225088 (SS, PKS) and U24 CA233243 (Dana-Farber Cancer Institute, Emory University, Institute for Systems Biology, Memorial Sloan Kettering Cancer Center, Sage Bionetworks). DS was funded by an Early Postdoc Mobility fellowship (no. P2ZZHP3_181475) from the Swiss National Science Foundation and was a Damon Runyon Fellow supported by the Damon Runyon Cancer Research Foundation (DRQ-03-20); DS is currently supported by the BMBF (01ZZ2004). NG was funded by the NIH Human BioMolecular Atlas Program (HuBMAP) OT2 OD026677 and MDH by NCI/NIH Task Order No. HHSN26110071 under Contract No. HHSN2612015000031.

Competing Interests Statement

PKS is a member of the SAB or BOD of Applied Biomath, RareCyte Inc., and Glencoe Software, which distributes a commercial version of the OMERO data management platform; PKS is also a member of the NanoString SAB and a consultant to Merck and Montai Health. In the last five years the Sorger lab has received research funding from Novartis and Merck. Sorger declares that none of these relationships have influenced the content of this manuscript. SS is a consultant for RareCyte Inc. NG is a co-founder and equity owner of Datavisyn. DS is a consultant for Roche Glycart AG. JRS is Founder and CEO of Glencoe Software, which distributes a commercial version of the OMERO data management platform. SR receives research funding from Bristol-Myers-Squibb, Merck, Affimed, and Kite/Gilead. SR is on the Scientific Advisory Board for Immunitas Therapeutics. DSU is employed by Quantitative Imaging Systems LLC. EAB is an employee of Indica Labs.

Data and Code Availability Statement

The detailed specification of the guidelines outlined in this manuscript are available at <https://github.com/miti-consortium/MITI> and <https://www.miti-consortium.org/>

References

1. Bodenmiller B Multiplexed Epitope-Based Tissue Imaging for Discovery and Healthcare Applications. *Cell Syst* 2, 225–238 (2016). [PubMed: 27135535]
2. Rozenblatt-Rosen O et al. The Human Tumor Atlas Network: Charting Tumor Transitions across Space and Time at Single-Cell Resolution. *Cell* 181, 236–249 (2020). [PubMed: 32302568]
3. HuBMAP Consortium. The human body at cellular resolution: the NIH Human Biomolecular Atlas Program. *Nature* 574, 187–192 (2019). [PubMed: 31597973]
4. Rajewsky N et al. LifeTime and improving European healthcare through cell-based interceptive medicine. *Nature* 587, 377–386 (2020). [PubMed: 32894860]
5. Wilkinson MD et al. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* 3, 160018 (2016). [PubMed: 26978244]
6. Chen KH, Boettiger AN, Moffitt JR, Wang S & Zhuang X RNA imaging. Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* 348, aaa6090 (2015). [PubMed: 25858977]

7. Fischer RS, Wu Y, Kanchanawong P, Shroff H & Waterman CM Microscopy in 3D: a biologist's toolbox. *Trends in Cell Biology* 21, 682–691 (2011). [PubMed: 22047760]
8. Ghaznavi F, Evans A, Madabhushi A & Feldman M Digital imaging in pathology: whole-slide imaging and beyond. *Annu Rev Pathol* 8, 331–359 (2013). [PubMed: 23157334]
9. Amin MB et al. The Eighth Edition AJCC Cancer Staging Manual: Continuing to build a bridge from a population-based to a more 'personalized' approach to cancer staging. *CA Cancer J Clin* 67, 93–99 (2017). [PubMed: 28094848]
10. Cancer Genome Atlas Research Network et al. The Cancer Genome Atlas Pan-Cancer analysis project. *Nat Genet* 45, 1113–1120 (2013). [PubMed: 24071849]
11. Brazma A et al. Minimum information about a microarray experiment (MIAME)-toward standards for microarray data. *Nat. Genet* 29, 365–371 (2001). [PubMed: 11726920]
12. Field D et al. The minimum information about a genome sequence (MIGS) specification. *Nat Biotechnol* 26, 541–547 (2008). [PubMed: 18464787]
13. Taylor CF et al. Promoting coherent minimum reporting guidelines for biological and biomedical investigations: the MIBBI project. *Nat Biotechnol* 26, 889–896 (2008). [PubMed: 18688244]
14. Benitez K & Malin B Evaluating re-identification risks with respect to the HIPAA privacy rule. *J Am Med Inform Assoc* 17, 169–177 (2010). [PubMed: 20190059]
15. Grossman RL et al. Toward a Shared Vision for Cancer Genomic Data. *N Engl J Med* 375, 1109–1112 (2016). [PubMed: 27653561]
16. Byrd JB, Greene AC, Prasad DV, Jiang X & Greene CS Responsible, practical genomic data sharing that accelerates research. *Nat Rev Genet* 21, 615–629 (2020). [PubMed: 32694666]
17. Swedlow JR, Goldberg I, Brauner E & Sorger PK Informatics and quantitative analysis in biological imaging. *Science* 300, 100–102 (2003). [PubMed: 12677061]
18. Li S et al. Metadata management for high content screening in OMERO. *Methods* 96, 27–32 (2016). [PubMed: 26476368]
19. Nelson G et al. QUAREP-LiMi: A community-driven initiative to establish guidelines for quality assessment and reproducibility for instruments and images in light microscopy. *arXiv:2101.09153 [physics, q-bio]* (2021).
20. Bandrowski A et al. The Resource Identification Initiative: A cultural shift in publishing. *F1000Res* 4, (2015).
21. Edfors F et al. Enhanced validation of antibodies for research applications. *Nat Commun* 9, 4130 (2018). [PubMed: 30297845]
22. Sarkans U et al. REMBI: Recommended Metadata for Biological Images—enabling reuse of microscopy data in biology. *Nat Methods* 1–5 (2021) doi:10.1038/s41592-021-01166-8. [PubMed: 33408396]
23. Liu J et al. An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-Quality Survival Outcome Analytics. *Cell* 173, 400–416.e11 (2018). [PubMed: 29625055]
24. Zhang Z et al. Uniform genomic data analysis in the NCI Genomic Data Commons. *Nat Commun* 12, 1226 (2021). [PubMed: 33619257]
25. Moore J et al. OME-NGFF: scalable format strategies for interoperable bioimaging data. 2021.03.31.437929 <https://www.biorxiv.org/content/10.1101/2021.03.31.437929v4> (2021) doi:10.1101/2021.03.31.437929.
26. Clunie DA Dual-Personality DICOM-TIFF for whole slide images: A migration technique for legacy software. *Journal of Pathology Informatics* 10, 12 (2019). [PubMed: 31057981]
27. Fedorov A et al. NCI Imaging Data Commons. *Cancer Res canres.0950.2021* (2021) doi:10.1158/0008-5472.CAN-21-0950.
28. Health, C. for D. and R. Technical Performance Assessment of Digital Pathology Whole Slide Imaging Devices. U.S. Food and Drug Administration <http://www.fda.gov/regulatory-information/search-fda-guidance-documents/technical-performance-assessment-digital-pathology-whole-slide-imaging-devices> (2019).
29. Lin J-R et al. Multiplexed 3D atlas of state transitions and immune interactions in colorectal cancer. *bioRxiv* 2021.03.31.437984 (2021) doi:10.1101/2021.03.31.437984.

30. Peng T et al. A BaSiC tool for background and shading correction of optical microscopy images. *Nat Commun* 8, (2017).
31. MCMICRO: A scalable, modular image-processing pipeline for multiplexed tissue imaging | bioRxiv. <https://www.biorxiv.org/content/10.1101/2021.03.15.435473v1>.
32. Heiser CN & Lau KS A Quantitative Framework for Evaluating Single-Cell Data Structure Preservation by Dimensionality Reduction Techniques. *Cell Rep* 31, 107576 (2020). [PubMed: 32375029]
33. Megill C et al. cellxgene: a performant, scalable exploration platform for high dimensional sparse matrices. bioRxiv 2021.04.05.438318 (2021) doi:10.1101/2021.04.05.438318.
34. Gutman DA et al. The Digital Slide Archive: A Software Platform for Management, Integration, and Analysis of Histology for Cancer Research. *Cancer Res.* 77, e75–e78 (2017). [PubMed: 29092945]
35. Rashid R et al. Interpretative guides for interacting with tissue atlas and digital pathology data using the Minerva browser. *Nat Biomed Eng.* 2020.03.27.001834 (2020) doi:10.1101/2020.03.27.001834.
36. Singh I et al. OMeta: an ontology-based, data-driven metadata tracking system. *BMC Bioinformatics* 20, 8 (2019). [PubMed: 30612540]
37. Rayner TF et al. A simple spreadsheet-based, MIAME-supportive format for microarray data: MAGE-TAB. *BMC Bioinformatics* 7, 489 (2006). [PubMed: 17087822]
38. Martínez-Romero M et al. Using association rule mining and ontologies to generate metadata recommendations from multiple biomedical databases. *Database* 2019, (2019).
39. Williams E et al. The Image Data Resource: A Bioimage Data Integration and Publication Platform. *Nat. Methods* 14, 775–781 (2017). [PubMed: 28775673]
40. Using Requester Pays buckets for storage transfers and usage - Amazon Simple Storage Service. <https://docs.aws.amazon.com/AmazonS3/latest/userguide/RequesterPaysBuckets.html>.
41. Kacprzak E et al. Characterising dataset search—An analysis of search logs and data requests. *Journal of Web Semantics* 55, 37–55 (2019).
42. Löffler F, Wesp V, König-Ries B & Klan F Dataset search in biodiversity research: Do metadata in data repositories reflect scholarly information needs? *PLOS ONE* 16, e0246099 (2021). [PubMed: 33760822]
43. Swedlow JR et al. A global view of standards for open image data formats and repositories. *Nat Methods* 1–7 (2021) doi:10.1038/s41592-021-01113-7. [PubMed: 33408396]

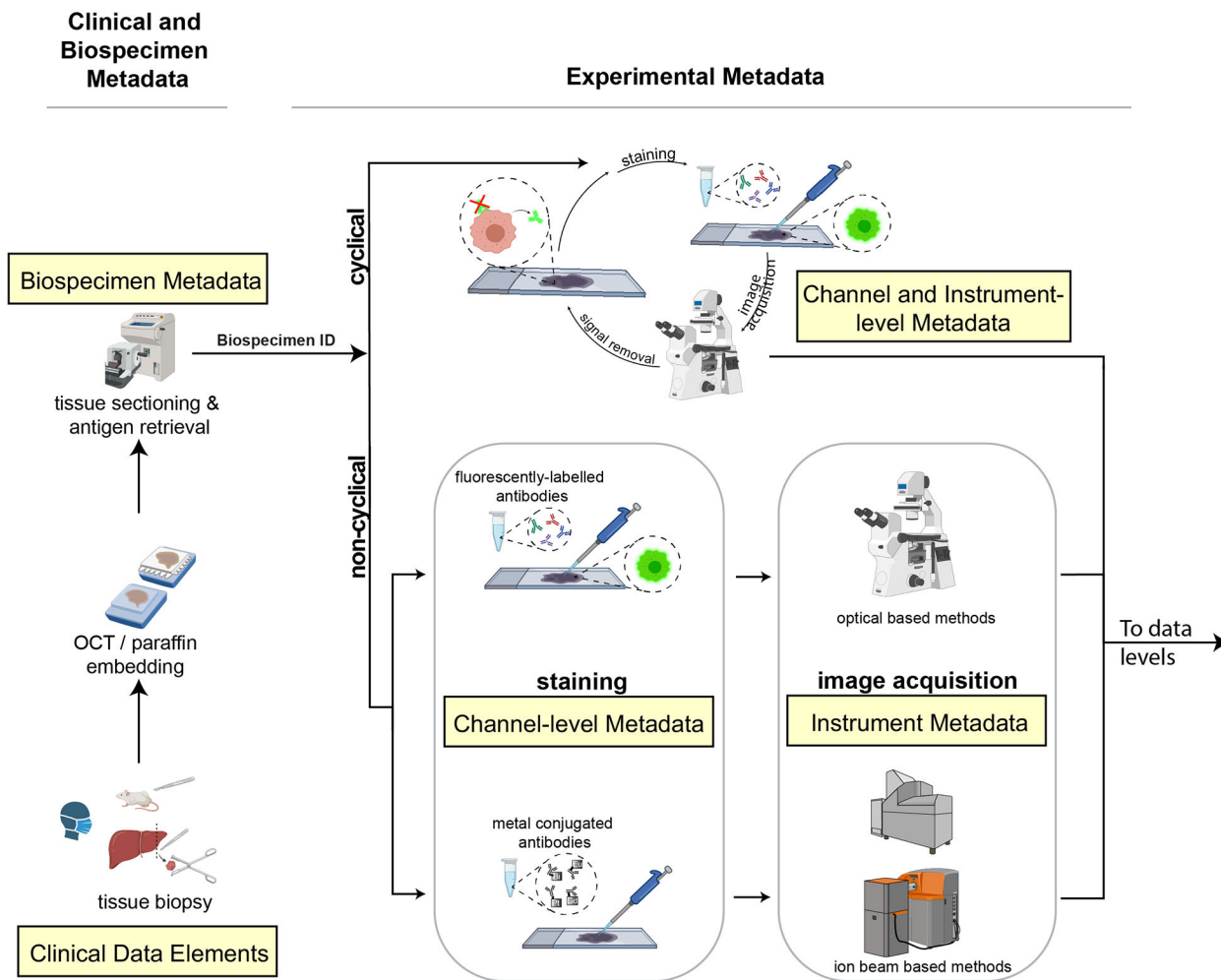


FIGURE 1: Schematic diagram of the steps in a canonical multiplexed tissue imaging experiment and the associated metadata
 In a typical workflow, samples collected from patient biopsies and resections or from animal models are formaldehyde fixed and paraffin embedded (FFPE) or frozen and then sectioned and mounted onto either a standard glass microscope slide (for CyCIF, mIHC, IMC, MELC or mxIF), fluidic chamber (for CODEX) or specialized carriers (for MIBI). Clinical and biospecimen metadata (extracted from clinical records, for example) is linked to all other levels of metadata via a unique ID (Biospecimen ID). Data is acquired using cyclical or non-cyclical staining and imaging methods and both reagent and experimental metadata collected (consisting of antibody, reagent and instrument metadata). In both cyclic and non-cyclic methods, sections undergo pre-processing, antigen retrieval, and antibody incubation and images are acquired. In cyclical imaging methods, fluorophores or chromogens are inactivated or removed and additional antibodies and/or visualization reagents are applied and data acquisition repeated. Channel and instrument metadata capture these essential details. Created with [BioRender.com](https://www.biorender.com).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

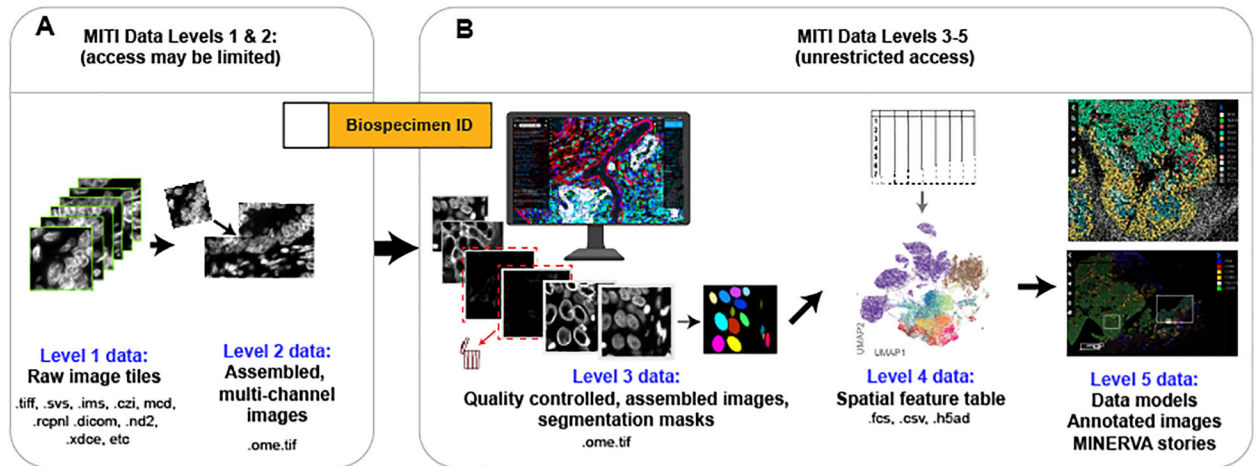


FIGURE 2: MITI data levels and formats

Data levels specify the extent of data processing and, in the case of sequencing data, whether access requires the approval of a data access committee. In common practice, data at levels 3 and up are freely shared. Primary data arising from microscopes and data acquisition instruments corresponds to level 1 data. Because the raw image data acquired from one slide usually consists of separate image fields, possibly from proprietary formats, they are processed to correct for uneven illumination and other instrumentation artifacts and assembled into a single multi-channel image in the OME-TIFF format (level 2 data). OME-TIFF image mosaics undergo quality control (including artefact removal, channel rejection, evaluation of staining quality) to generate full-resolution, assembled and curated level 3 image data; segmentation algorithms generate one or more label masks that also comprise level 3 data. The great majority of users will want to access these level 3 images. Each label mask (e.g., nuclei, cytoplasmic-regions, whole cells, organelles, etc.) is used to compute quantitative features, such as the mean signal intensity, spatial coordinates of individual cells and morphological features, which are stored as level 4 spatial feature tables (where rows represent single cells and columns the extracted cellular features); these data are suitable for analysis using the dimensionality reduction and visualization tools used for other types of single-cell data (e.g. UMAP plots). Spatial models computed from images and spatial feature tables, or by direct application of machine learning to images, as well as images annotated by humans, comprise level 5 data.