

## DATA METROLOGY FOR LIFE SCIENCES, MEDICINE AND PHARMACEUTICAL MANUFACTURING

*Paul M. Duncan*<sup>a,\*</sup>, *Nadia A. S. Smith*<sup>a</sup>, *Marina Romanchikova*<sup>a</sup>

<sup>a</sup> Data Science Department, National Physical Laboratory, United Kingdom

\* Corresponding author. E-mail address: paul.duncan@npl.co.uk

**Abstract:** In many disciplines, such as physics and engineering, the application of tools to support data metrology is encouraged and embedded in many processes and applications while in the life sciences, medicine and pharmaceutical manufacturing sectors these tools are often added as an afterthought, if considered at all. The use of data-driven decision making and the advent of machine learning in these industries has created an urgent demand for harmonised high-quality, instantly available, datasets across domains. The Findable, Accessible, Interoperable, Reproducible principles are designed to improve overall quality of research data. However, this alone does not guarantee that data is fit-for-purpose. Issues such as missing data and metadata, insufficient knowledge of measurement conditions or data provenance are well known and can be aided by applying metrological concepts to data preparation to increase confidence. This work presents the data metrology projects conducted by the National Physical Laboratory Data Science team in healthcare applications.

**Keywords:** NMI, metrology, digital pathology, medicines manufacturing, metadata standards, data quality, ontologies, FAIR principles.

### 1. INTRODUCTION

Technological advancements in medicines and pharmaceutical manufacturing have been traditionally focused on advances in drug discovery, experimental procedures and manufacture. Medicines and treatments are becoming more expensive to produce, as pricing models drive down profit margins compounded with patents expiry [1]. Therefore, a greater emphasis is being placed on maximising the efficiency of medicine development and manufacture.

For the quality and the repeatability of processes, most pharmaceutical firms operate at high variation levels in terms of accurately manufacturing materials. These variations, at levels between 3 and 4 $\sigma$ , are estimated to cost ~\$20bn annually through waste and inefficiency [2]. Therefore, companies are increasingly moving to developing controlled and flexible processes to offer digital health solutions for their customers.

The National Physical Laboratory (NPL) has been addressing the problem of aiding digitalisation in healthcare by focusing on the issue of data metrology for life sciences, medicines and pharmaceutical manufacturing. Data metrology refers to the uncertainty present in the data generated in each of these

areas, from the quality of measurements accompanying the manufacturing to the quality of the data used for decision making processes.

This paper describes the similarities and differences between data metrology challenges addressed by NPL in the context of several cross-disciplinary projects with the goal of helping users to identify their data metrology needs and delivering confidence in the effective use of data.

### 2. DATA METROLOGY PROJECTS

The NPL Data Science team has been involved in multiple data metrology projects including life sciences, healthcare, and medicines manufacturing applications to highlight the similarities and domain-specific requirements to data quality and management. These projects and data metrology challenges are described below.

#### 2.1. Pharmaceutical manufacturing

Recent developments in digital pharmaceutical manufacturing are generating a large amount of data across varying temporal resolutions and manufacturing routes. This data provides unprecedented opportunities for pharmaceutical manufacturing to derive new insights and efficiencies from experiments but imposes great challenges in data processing, management, sharing, and integration. Not only data integrity and authenticity are to be ensured, but the processes that lead to the generation of data must be traceable to enable trust.

The pharmaceutical industry introduced “Good Manufacturing Practices” (GMP) to standardise processes around quality, security and effectiveness, but did not make allowances for metrological concepts such as traceability and measurement uncertainty. Data metrology therefore becomes a critical component in understanding and controlling pharmaceutical processes and reducing the variation seen in the final product. NPL has worked with major pharmaceutical manufacturers and researchers to explore their data metrology needs and develop a set of applied research programs.

1) *Ontologies for clinical trial release.* NPL has developed techniques [3] to develop a Domain Agnostic Measurement Ontology, with a view to applying these techniques across different industries. For the past 2 years, NPL has worked with the Medicines Manufacturing Innovation Centre (MMIC) to develop an ontology to aid in the automation and digitalisation of all data required for regulators to approve drugs for consumption. Much of the approvals process is manual data processing which can be

replaced by processing of data through modern data driven techniques. The ontology developed by NPL “codifies” all the data relationships which are pertinent to the identification of an expiry date for the release of a drug, facilitating true automation, reducing human input and significantly decreasing the potential for errors in the process due to the development of an approved automated decision process.

2) *Controlled vocabularies for pharmaceutical data exchange.* The development of the ontology for clinical trials release exposed an issue in how data from different companies and manufacturers can differ semantically when describing similar terms. For example, separate companies may use the terms “pill” and “tablet” to describe the same concept. This inconsistency decreases the quality of the information used in digitalised or automated systems. NPL has been developing a controlled vocabulary for the clinical trials process to ensure that any automated system can understand the terminology used by each party. This *controlled* vocabulary can provide a traceable link to quality processes for each company which can aid in automating the verification of the processes used. NPL has also been exploring the idea of working with industry to create an industry-wide standard to create a unified approach to solving this problem and reducing the uncertainty of the information.

3) *Mapping of measurement uncertainty propagation in manufacturing.* NPL has been working on an approach to understand the measurement uncertainty generated at each stage of a continuous manufacturing process. Currently, uncertainty generated at each node is not propagated, so to ensure greater traceability of variation present in the final product, we are currently developing methods to “map” out the uncertainty present and propagate this to each stage.

The goal of these use cases under development is to truly understand the uncertainty of the information produced during pharmaceutical manufacturing and to provide industry with frameworks to understand their data metrology.

## 2.2. Minimum metadata for biological imaging

Biological imaging (bioimaging) includes a vast array of techniques that includes optical microscopy, spectroscopy, multispectral imaging, among others. In the pharmaceutical industry, these techniques are used both in R&D and in clinical studies that evaluate drug resistance, efficacy, targeting mechanisms and pharmacodynamics. Complexity, diversity, and volume of data generated by high-resolution imaging techniques drive the need for advanced analysis and data management methods. While work to improve data interoperability is ongoing [4], [5], achieving reproducibility of results and re-usability of data is challenging without agreed metadata requirements and data formats [6].

Working with two major industry partners, NPL has identified three bioimaging case studies characterised by high data volumes and need for re-use: 1) mass spectrometry imaging (MSI); 2) high content screening, and 3) light sheet microscopy (LSM). Minimum metadata requirements were collated and shared between project partners using the metadata categories illustrated in Figure 1.

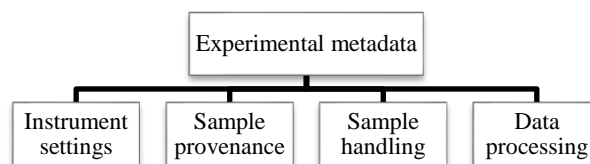


Figure 1: Metadata categories in bioimaging.

The minimum metadata requirements for LSM and MSI were used to develop frameworks for bioimaging data capture and annotation at NPL [7]. While efforts have been undertaken to define minimum reporting standards and metadata [8], [9], stronger engagement of equipment vendors, researchers and funding authorities is needed to create future-proof re-usable and reproducible data repositories.

## 2.3. Digital pathology

Clinical histopathology describes a study of stained tissue sections on glass slides under a microscope, whereby pathologists manually change the brightness, focus depth, and the region of interest. In digital pathology (DP), tissue samples are digitised using a whole slide imaging (WSI) scanner. The resulting high-resolution images (1-4 GB) can be studied *in silico* by image analysis software or on-display by a trained pathologist. The DP workflow poses multiple metrological challenges: reproducibility and repeatability of tissue processing, calibration and traceability of WSI, as well as uncertainty analysis to support diagnosis.

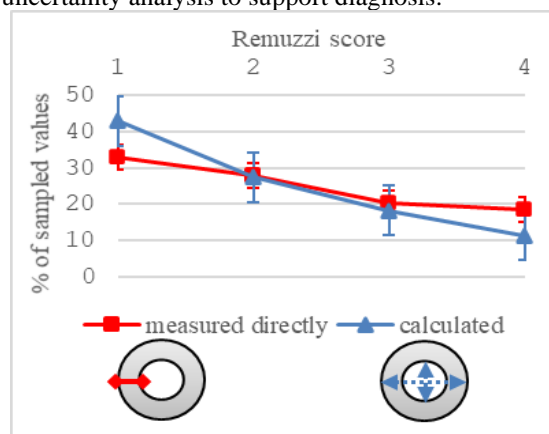


Figure 2: Impact of measurement method on clinical assessment (Remuzzi score). Red line: wall thickness is measured directly. Blue line: wall thickness is calculated from vessel outer diameter and lumen diameter. Image courtesy of Tobi Ayori.

The NPL Digital Pathology inter-disciplinary project, launched in 2020, comprised a landscape exercise during which DP experts and stakeholders identified priority areas for metrology support [10]. The outcomes of the landscape exercise were used to shape demonstrator studies with real-world data. Within the collaboration with PITHIA trial (<http://www.pithia.org.uk>, Grant Reference Number PB-PG-1215-20033), we are studying the uncertainties in the diagnosis based on kidney biopsy images, aiming to 1) locate the sources of uncertainty in decision making and find tools to reduce it, and 2) find image features that correlate with clinical outcomes to increase reproducibility and explainability in WSI evaluation

The preliminary findings (Figure 2) show how on-display assessment method influences the diagnostic results: when a blood vessel wall thickness is measured directly (red line,

lower left diagram), the assessors show preference for more uniform score assignment than if the wall thickness is calculated as a difference  $(\text{outer diameter} - \text{lumen diameter}) / 2$  (blue line, lower right diagram).

Further case studies will include analysis of measurable image features and their association with diagnostic predictions, as well as impact assessment of intra- and inter-WSI device variability on image features and diagnosis.

Future work will include engaging with standards bodies to include metrology-enabling contextual data such as calibration results, device settings etc. into clinical DP standards such as Digital Communications for Medical Imaging (DICOM) and Fast Healthcare Interoperability Resources (FHIR). These standards have high maturity levels and provide mechanisms to include metrological metadata and requirements such as units of measure, clinical terminologies, ontologies, unique identifiers etc.

#### 2.4. Medical sensors case study

While WSI data and associated measurement information can be captured using the existing DICOM standard, novel medical devices require modification of existing standards to capture new data types and provide integration into the healthcare infrastructure. NPL worked with a UK-based medical device developer to create clinically interoperable data structures to store and manage the data from a novel surgical sensor. This opportunity facilitated the capture of valuable metrological information including traceability and calibration *ab initio*, creating a metrologically sound data model at the early stage of device development.

An example of how custom measurement-related information can be included into DICOM metadata is presented in Table 1. A custom value (patient tilting angle in degrees) is enclosed in a Concept Name Code Sequence that refers to the coding document and provides the value inclusive of its format (value representation/VR). Note that the value description includes the unit of measure and the reference terminology (Measurement Units Code Sequence).

Table 1. Including custom measurement value, units of measure and reference to ontology in DICOM metadata.

Tag description	Tag	VR	Value
<b>Concept Name Code Sequence</b>	(0040, A043)	SQ	-
Code Value	(0008, 0100)	SH	'1.2.2-1'
Coding Scheme Designator	(0008, 0102)	SH	'ASCODE'
Coding Scheme Version	(0008, 0103)	SH	'1.0'
Code Meaning	(0008, 0104)	LO	'Patient tilting angle'
Numeric Value	(0040, A30A)	DS	'-19.05'
<b>Measurement Units</b>	(0040, 08EA)	SQ	-

Code Sequence				
	Code Value	(0008, 0100)	SH	'deg'
	Coding Scheme Designator	(0008, 0102)	SH	'UCUM'
	Coding Scheme Version	(0008, 0103)	SH	'1.4'
	Code Meaning	(0008, 0104)	LO	'degrees'

#### 2.5. Digital health

New measurement modalities within healthcare are creating vast amounts of high-dimensional data from disparate sources and of varying quality, including genomic, imaging, biomarkers, electronic healthcare records and data from wearable devices. The current and future healthcare practices across the world are increasingly reliant on the integration of these diverse, complex, and large datasets as well as trusted and robust analysis methods [11]. The data curation process in healthcare includes extraction, de-identification, and annotation of datasets with metadata, as well as data fusion and linkage. Therefore, future-proof secure scalable curation methods that handle rapidly growing data volumes are needed.

NPL runs an ongoing inter-disciplinary Digital Health programme<sup>1</sup> aimed to use data metrology tools to help solve some of the important and emerging challenges of utilising healthcare data. The project includes several case studies detailed in the 2021 report [12].

One of the case studies investigates whether it is possible to improve the data quality and comparability by linking patient images with imaging device calibration data. The study set out to link megavoltage computed tomography (MVCT) images used for image-guided radiotherapy with MVCT device calibration data from the routine monthly quality assurance tests that check whether the scanner is fit-for-purpose. MVCT images are routinely used for patient positioning, radiation dosimetry, and in-treatment therapy effect assessment. Like other medical imaging modalities, MVCT images are subject to temporal and inter-device variations that are known to have negative influence on the accuracy of subsequent radiation dose calculation and image segmentation. We implemented a procedure that includes the device calibration information into the DICOM header information of the patient scan. We expect that the MVCT calibration data can be used to remove the device-related variability and make the patient images more inter-comparable, reduce the variations in the image quality, improving the accuracy of analysis, safety, and efficiency of data-driven clinical interventions.

A further case study in the Digital Health programme evaluates how data linkage can be used to improve the quality of life and long-term treatment outcomes for prostate cancer patients by using the patient care data acquired outside of clinical trials [13]. We developed an ontology-based data curation framework to identify and collate information about diagnosis, symptoms, and treatment side effects from routine primary care electronic health records. This work is a first step

<sup>1</sup> <https://www.npl.co.uk/data-science/digital-health> (accessed

on 11/04/2022)

to increase the utility of primary care data for oncology by a) creating a knowledge base of data sources, b) mapping out the required integration efforts, and c) developing a practical ontology-based method for systematic and reproducible prostate cancer case identification and validating this method on real-world datasets. The developed ontology can be used to standardise the identification and retrieval of prostate cancer cases from primary care data.

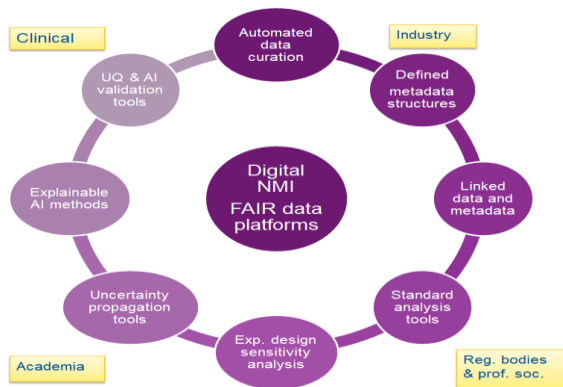


Figure 3: FAIR data platforms for clinically relevant research

NPL's most recent endeavours to increase availability and reliability of medical data include developing a curated data platform. The platform will provide mechanisms for curation, storage, metadata annotation, linkage, and analysis of clinically relevant imaging, audit, and calibration data (Figure 3). Such a platform would provide a much-needed foundation to enable access to a much richer and larger dataset than what is currently available, rendering the data FAIR-er, and thus increasing its value and utility.

### 3. CONCLUSIONS

While the domains for data metrology applications vary significantly within pharmaceutical manufacturing, healthcare and bioimaging, similarities have emerged from the projects outlined above. Firstly, the need for FAIR-ness and data reusability calls for a systematic approach to data curation and metadata annotation. Ontologies and controlled vocabularies help to address this gap but striving towards standards and minimum data quality requirements is recommended to increase data re-usability and impact across different sectors/companies. Second, there is a variety of existing open standards and formats that can and should be used to manage data from new medical devices and imaging modalities. These standards can be adapted to incorporate information pertaining to metrological traceability and uncertainty. Third, while the use of, and need for, metrology methods is widely recognised in physics and engineering, in life sciences, medicine and pharmaceutical manufacturing these tools are often added as an afterthought, if considered at all. Therefore, work is required to demonstrate the impact of data metrology via case studies in the respective domains.

NPL is building a range of use cases and demonstrators to highlight the need for data metrology across the healthcare industries, including applications across digital pathology, bioimaging, pharmaceutical and bio-manufacturing through interacting with industry and researchers in these fields.

### ACKNOWLEDGMENTS

This work was funded by the UK Government Department for Business, Energy & Industrial Strategy through the UK's National Measurement System. We would also like to thank our partners at the MMIC; CPI, University of Strathclyde, UKRI, Scottish Enterprise, AstraZeneca and GSK as well as ArtioSense Ltd and the PITHIA trial investigators.

### REFERENCES

- [1] D. Taylor, 'The Pharmaceutical Industry and the Future of Drug Development', in *Pharmaceuticals in the Environment*, 2015, pp. 1–33. doi: 10.1039/9781782622345-00001.
- [2] J. S. Srai, C. Badman, M. Krumme, M. Futran, and C. Johnston, 'Future Supply Chains Enabled by Continuous Processing—Opportunities and Challenges. May 20–21, 2014 Continuous Manufacturing Symposium', *J. Pharm. Sci.*, vol. 104, no. 3, pp. 840–849, 2015, doi: 10.1002/jps.24343.
- [3] J.-L. Hippolyte, M. Chrubasik, F. Brochu, and M. Bevilacqua, 'A domain-agnostic ontology for unified metrology data management', *Meas. Sens.*, vol. 18, p. 100263, Dec. 2021, doi: 10.1016/j.measen.2021.100263.
- [4] U. Sarkans *et al.*, 'REMBI: Recommended Metadata for Biological Images—enabling reuse of microscopy data in biology', *Nat. Methods*, pp. 1–5, May 2021, doi: 10.1038/s41592-021-01166-8.
- [5] C. Allan *et al.*, 'OME Remote Objects (OMERO): a flexible, model-driven data management system for experimental biology', *Nat. Methods*, vol. 9, no. 3, pp. 245–253, Feb. 2012, doi: 10.1038/nmeth.1896.
- [6] A. Stuppel, D. Singerman, and L. A. Celi, 'The reproducibility crisis in the age of digital medicine', *Npj Digit. Med.*, vol. 2, no. 1, pp. 1–3, Jan. 2019, doi: 10.1038/s41746-019-0079-z.
- [7] F. Brochu *et al.*, 'Federation of Imaging Data for Life sciences: current status of metadata collection for high content screening, mass spectrometry imaging and light sheet microscopy of AstraZeneca, GlaxoSmithKline and NPL', May 2020. <https://doi.org/10.47120/npl.MS24> (accessed Jun. 09, 2021).
- [8] O. J. R. Gustafsson *et al.*, 'Balancing sufficiency and impact in reporting standards for mass spectrometry imaging experiments', *GigaScience*, vol. 7, no. 10, Oct. 2018, doi: 10.1093/gigascience/giy102.
- [9] M. Huisman *et al.*, 'Minimum Information guidelines for fluorescence microscopy: increasing the value, quality, and fidelity of image data', *ArXiv191011370 Cs Q-Bio*, Jan. 2020, Accessed: Mar. 09, 2020. [Online]. Available: <http://arxiv.org/abs/1910.11370>
- [10] M. Adeogun *et al.*, 'Metrology for Digital Pathology. Digital pathology cross-theme project report', Mar. 2021. <https://doi.org/10.47120/npl.AS102> (accessed Mar. 03, 2022).
- [11] 'The Topol Review — NHS Health Education England', *The Topol Review — NHS Health Education England*. <https://topol.hee.nhs.uk/> (accessed Mar. 31, 2022).
- [12] N. A. S. Smith, D. Sinden, S. A. Thomas, M. Romanchikova, J. E. Talbott, and M. Adeogun, 'Building confidence in digital health through metrology', *Br. J. Radiol.*, vol. 93, no. 1109, p. 20190574, May 2020, doi: 10.1259/bjr.20190574.
- [13] N. Smith, M. Romanchikova, I. Partarrieu, E. Cooke, A. Lemanska, and S. Thomas, 'NMS 2018-2021 Life-sciences and healthcare project "Digital health: curation of healthcare data" - final report', National Physical Laboratory, Nov. 2021. doi: 10.47120/npl.MS31.