


METHODOLOGY ARTICLE

Open Access



# A data management infrastructure for the integration of imaging and omics data in life sciences

Luis Kuhn Cuellar<sup>1</sup>, Andreas Friedrich<sup>1</sup>, Gisela Gabernet<sup>1</sup>, Luis de la Garza<sup>1</sup>, Sven Fillinger<sup>1</sup>, Adrian Seyboldt<sup>1</sup>, Tobias Koch<sup>1</sup>, Sven zur Oven-Krockhaus<sup>2</sup>, Friederike Wanke<sup>2</sup>, Sandra Richter<sup>2</sup>, Wolfgang M. Thaiss<sup>3</sup>, Marius Horger<sup>4</sup>, Nisar Malek<sup>4</sup>, Klaus Harter<sup>2</sup>, Michael Bitzer<sup>4</sup> and Sven Nahnsen<sup>1,5\*</sup> 

\*Correspondence:  
sven.nahnsen@uni-tuebingen.de

<sup>1</sup> Quantitative Biology Center (QBiC), University of Tübingen, Tübingen, Germany  
Full list of author information is available at the end of the article

## Abstract

**Background:** As technical developments in omics and biomedical imaging increase the throughput of data generation in life sciences, the need for information systems capable of managing heterogeneous digital assets is increasing. In particular, systems supporting the findability, accessibility, interoperability, and reusability (FAIR) principles of scientific data management.

**Results:** We propose a Service Oriented Architecture approach for integrated management and analysis of multi-omics and biomedical imaging data. Our architecture introduces an image management system into a FAIR-supporting, web-based platform for omics data management. Interoperable metadata models and middleware components implement the required data management operations. The resulting architecture allows for FAIR management of omics and imaging data, facilitating metadata queries from software applications. The applicability of the proposed architecture is demonstrated using two technical proofs of concept and a use case, aimed at molecular plant biology and clinical liver cancer research, which integrate various imaging and omics modalities.

**Conclusions:** We describe a data management architecture for integrated, FAIR-supporting management of omics and biomedical imaging data, and exemplify its applicability for basic biology research and clinical studies. We anticipate that FAIR data management systems for multi-modal data repositories will play a pivotal role in data-driven research, including studies which leverage advanced machine learning methods, as the joint analysis of omics and imaging data, in conjunction with phenotypic metadata, becomes not only desirable but necessary to derive novel insights into biological processes.

**Keywords:** Data integration, Imaging, Omics, Metadata models, Distributed systems, Service oriented architecture, Data management infrastructure



## Background

Current technical advances in the life sciences allow researchers to collect large amounts of data that open fundamentally new routes for a multitude of research questions. Typically, multi-omics and biomedical imaging data from biological samples are building the most important data basis in such studies. The large volume of data produced by various omics disciplines (e.g. genomics, transcriptomics, proteomics, metabolomics), biomedical imaging techniques (e.g. X-ray CT and PET), and the unprecedented increase in spatial resolution achieved by conventional confocal and super-resolution light microscopy [1] and modern electron microscopes [2], present a challenge for long-term storage and management of these high-dimensional digital assets, especially with regard to the requirements imposed by the FAIR data management principles [3]. Moreover, it is of particular importance to employ rich metadata models that allow researchers to relate data from different disciplines, and design experiments using an integrative approach to handle both, multilayer omics, as well as biomedical imaging data. The increasing data volume and experimental design complexity call for data management systems allowing for the integration of multi-omics and imaging data.

In recent years, a variety of data management and analysis systems for life science have been introduced. Among many, two examples are the Galaxy platform, a web-based workflow system for reproducible genomic analysis [4, 5], and the cBio Portal, which focuses on exploration and visualization of large datasets of cancer genomics data [6]. For proteomics studies, the Swiss Grid Proteomics Portal (iPortal) provides web-based analysis tools [7]. However, the above-mentioned platforms manage only omics data. In contrast, the Open Microscopy Environment Remote Objects (OMERO) is a sophisticated image management system for biology and medicine. The OMERO server is capable of managing imaging data a large variety of microscopy and medical imaging modalities (e.g. fluorescence and electron microscopy, histological imaging and medical CT), since it is able to handle multi-channel, 2D or 3D imaging data, with time series support. Given that the abovementioned imaging characteristics and amount of data produced by different modalities may change drastically, management of biomedical imaging is a challenging task. Nevertheless, despite its advanced image-related capabilities, it offers limited support for managing omics data [8]. While several systems focus on the management of omics or imaging data [8, 9], or are dedicated to bioinformatics workflows [4, 5, 10], to our knowledge no approach has been suggested to provide a solution for the integrated management of multi-omics and imaging data in parallel. Additionally, most of the available platforms focus on a limited set of omics disciplines, thus leading to metadata models that are not well-suited to describe complex experimental designs with multiple omics and imaging modalities.

The lack of comprehensive and validated metadata storage is one of the pitfalls of reproducible research in the genomic age [11]. Accordingly, many researchers are starting to embrace proposed standards, such as the FAIR data principles, for omics data management. To address the aforementioned issues, we recently introduced qPortal, a web-based solution, which provides a basis to support the FAIR guidelines, for the management of the entire added-value chain of omics-based biomedical data [12]. qPortal is a web-based portal for scientific data management and analysis, which uses the Open Biological Information System (openBIS) platform [9] as a backend. While FAIR data

management for genomics and proteomics studies has been well supported using qPortal, it lacks the capability of managing imaging data (e.g. from light microscopy or MRI).

Here, we propose a solution based on principles of Service Oriented Architecture (SOA) design [13], to integrate an interoperable image management system for the accommodation of data from microscopy and medical imaging studies in qPortal. Our suggested infrastructure enables the integrated analysis of imaging and omics data, which is an increasingly deployed strategy in biology and medical studies [14–16]. We suggest using an OMERO server as a backend component for all imaging data. In architectural conjunction with openBIS, building the conventional backbone for qPortal. The OMERO server already provides an infrastructure that supports FAIR principles for imaging data itself, while providing large-scale data storage, management, and compatibility with common imaging file formats [17]. On the other hand, the openBIS server assists in the implementation of the FAIR principles on the metadata describing the experimental design of research projects and the resulting omics data, while qPortal provides a web-based interface to present users with a set of applications for scientific data management.

## Results

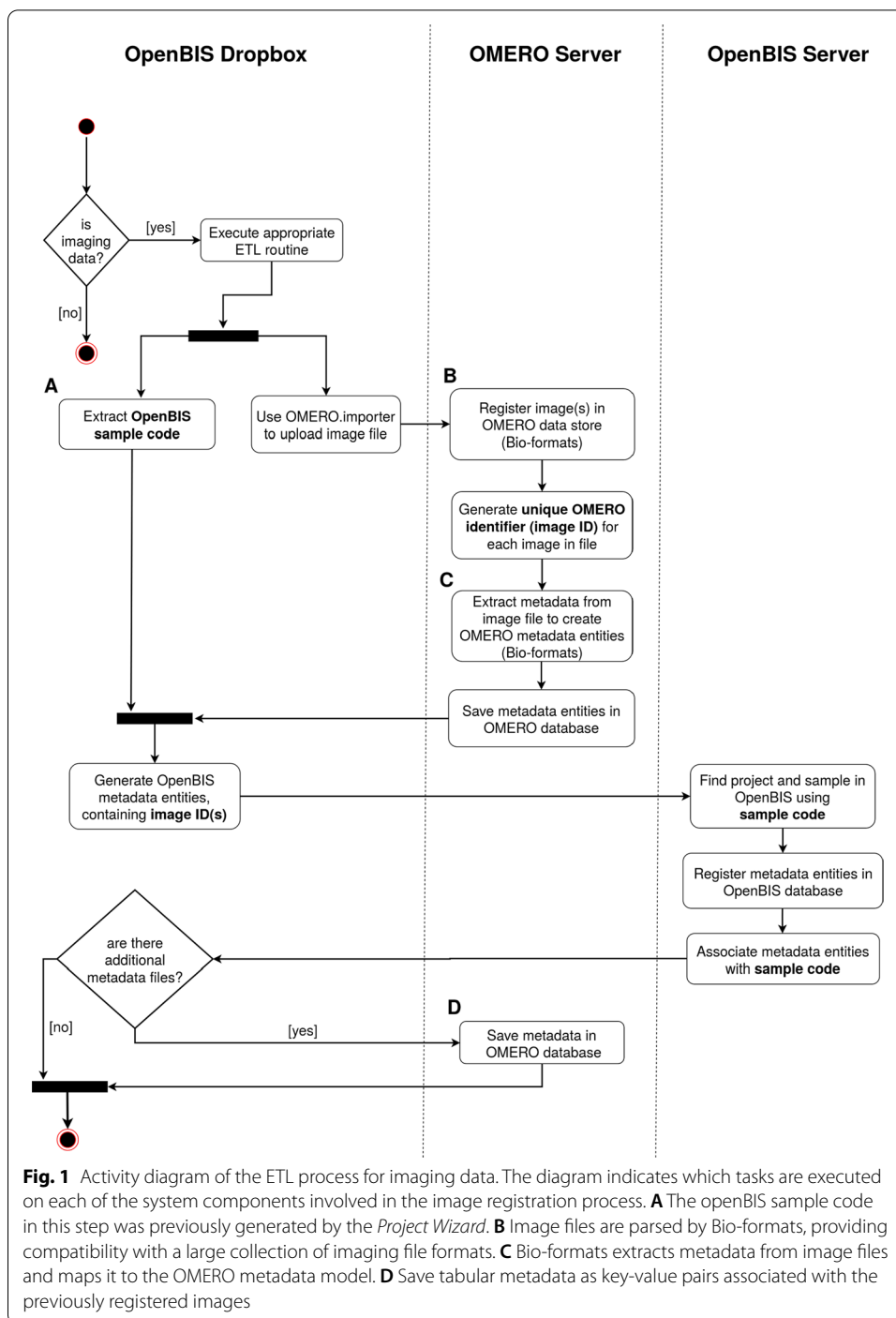
### Integration of OMERO into the qPortal platform

The integration of an OMERO server as an additional backend component of qPortal, requires substantial changes in the data handling and managing processes, considering that metadata information between the openBIS and OMERO servers has to be connected and kept synchronized in a structured manner (Additional file 1: Fig. S1). First, an integrative model was established, which defines detailed metadata boundaries between the project and omics domains, as managed by the openBIS server, and the imaging domain stored in the OMERO server. The technical implementation of this model requires extending server-side qPortal applications and processes, and the development of modules to facilitate communication between these applications and the OMERO server.

In order to allow communication between the OMERO server and qPortal, we developed a Java library wrapping the OMERO server API. This *OMERO client* library is analog to the *openBIS Client* component of qPortal [12], and is used by the *Project Wizard* application to implement the metadata model connection during project creation, i.e. creating synchronized metadata entities and structures in both the openBIS and OMERO platforms (see [Metadata Models](#)).

To effectively provide web-based image visualization in qPortal, we developed an *Image Viewer* portlet. This application provides easy access to imaging metadata and image visualization, using the *OMERO Client* component to query the OMERO server. The *Image Viewer* application also accesses OMERO.web functionality, in particular the 5D image viewers.

Finally, a suitable ETL routine for imaging data was created. This routine uploads image data into the OMERO platform, and creates a record of the available imaging data in the openBIS server. This can be achieved by creating metadata entities in the openBIS database that contain symbolic links to the images stored in OMERO binary repository (see Fig. 1).



### Metadata models

The full data management potential of both OMERO and qPortal can be achieved within a unified system, when their metadata models are centralized, or at least connected in a systematic way. This allows for full access to the metadata management capacities of openBIS, the metadata structure for project and multi-factorial experiments implemented by qPortal, including its domain-specific and feature-rich metadata support,

while using the metadata model and storage facilities of OMERO to describe microscopy and medical imaging data and its acquisition process. It is imperative that both models are linked so that any metadata transaction preserves the correct execution of CRUD operations (create, read, update, delete) on the semantic level, considering the databases of both platforms. Metadata redundancies, and ambiguous, or orphaned metadata relations have to be avoided. Since the metadata model of both platforms has a hierarchical organization, with conceptually similar metadata entities in the main structure, it is straightforward to align and connect metadata entities across models using cardinality relationships (e.g. one-to-many, many-to-many). The general openBIS model is composed of four main hierarchical levels, with project entities on top, followed by experiments, samples, and finally datasets. Similarly, the OMERO model uses a hierarchical structure, with project entities on top, followed by datasets, images, pixels, and finally features. In order to connect both metadata structures, we propose a loose-coupling approach, following SOA principles [13]. Here, two major one-to-one correspondences are established, one between “project” entities and another one relating qPortal sample entities with OMERO dataset entities. Additional file 1: Fig. S1 depicts the unified metadata model.

While both models can handle extensible XML schemas, we aim to distribute metadata information in a manner that benefits from the most useful features of each platform. Therefore, we allow the qPortal model to continue describing the main experimental design of the project and capture the biology of the samples, while using the OMERO model to store image related metadata, e.g. pixel size, magnification, image acquisition parameters, technical specifications of the microscope or medical imaging device. With this approach we can benefit from both Bio-Formats and the project planning functionality offered by qPortal.

#### **Technical implementation**

The *OMERO client* module was implemented as a java-based library (<https://github.com/qbicsoftware/omero-client-lib>). It provides an interface to read and write metadata on the OMERO server from portlet applications, the interface is tailored to implement the necessary operations needed to maintain the proposed coupling of metadata models. This library communicates with the OMERO server via the OMERO Java API and encapsulates connectivity logic with the backend, allowing portlet applications to focus on user interface and metadata synchronization between the openBIS and OMERO servers.

The *Project Wizard* application (<https://github.com/qbicsoftware/projectwizard-portlet>) was extended to provide the user with an option to enable imaging support to their omics projects (i.e. OMERO functionality). The *Project Wizard* uses the *OMERO client* library to create the appropriate metadata structure in the OMERO server, respecting the correspondences with metadata entities in the openBIS server. In technical terms, entity correspondence is achieved by providing the OMERO entities with the unique identifiers of the corresponding entities in openBIS.

An additional portlet application was developed to allow easy web-based access to imaging data from qPortal. The *Image Viewer* application (<https://github.com/qbicsoftware/omero-client-portlet>) uses the *OMERO client* component to query metadata

information (e.g. image identifiers, spatial size, image time points and channels) and image thumbnails from the OMERO server, and display this information in a structured manner that follows the proposed metadata model coupling. That is, a user must select a previously created project and sample (as created by the *Project Wizard*) to access the associated imaging data (Fig. 3A). Through the *Image Viewer* application, users can directly access a full, 5D view of any image, which is provided by the OMERO.web server (Fig. 3B). To allow direct access to the 5D image viewer, the *Image Viewer* application creates an HTTP session with the OMERO.web server via the OMERO json API, providing the user with a consistent single-sign-on (SSO) experience. Similarly, the OMERO.iviewer application (an OMERO.web plugin) can be employed to facilitate web-based annotation of ROIs (Fig. 3C).

Imaging data registration is achieved by creating a specialized openBIS *dropbox* and respective ETL procedure, which can connect to the OMERO server through the Python API, and uses the OMERO.importer application as an external tool. In short, the OMERO.importer uploads raw imaging data and uses Bio-formats to extract all metadata from open and proprietary file formats before mapping it to the OMERO metadata model and registering it accordingly. Subsequently, registration of the uploaded images in the openBIS server is achieved by creating specialized metadata entities that contain the unique OMERO identifiers of each uploaded image. Once an image with a correct sample identifier has been uploaded using the ETL routine, the OMERO server will generate a thumbnail of the image and allow access to all associated metadata and raw data via the *Image Viewer* application. Figure 1 depicts the activity diagram of the aforementioned ETL routine.

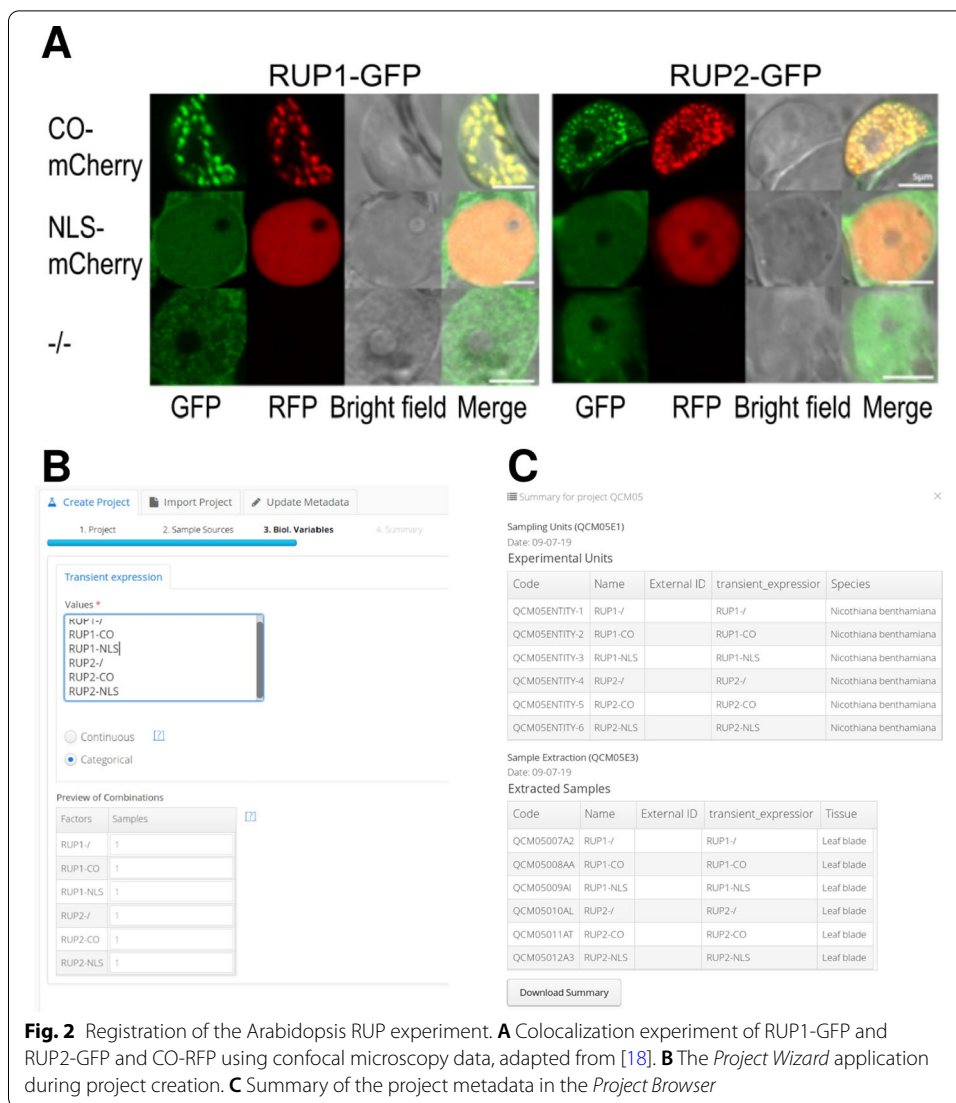
### **Proof of concept and use case**

We tested our newly developed infrastructure using two technical proofs of concept and a use case, with projects ranging from basic single layer experiments to large multi-scale clinical studies. Here we illustrate the broad applicability of an architecture for integrated management of omics and imaging data using two proof of concept projects, one focuses on fluorescence microscopy data from plant biology research, while the other is synthetic clinical project, containing a large dataset of medical imaging and gene sequencing data. Additionally, we present a use case with a 3-layer omics dataset and multi-modal imaging of liver cancer.

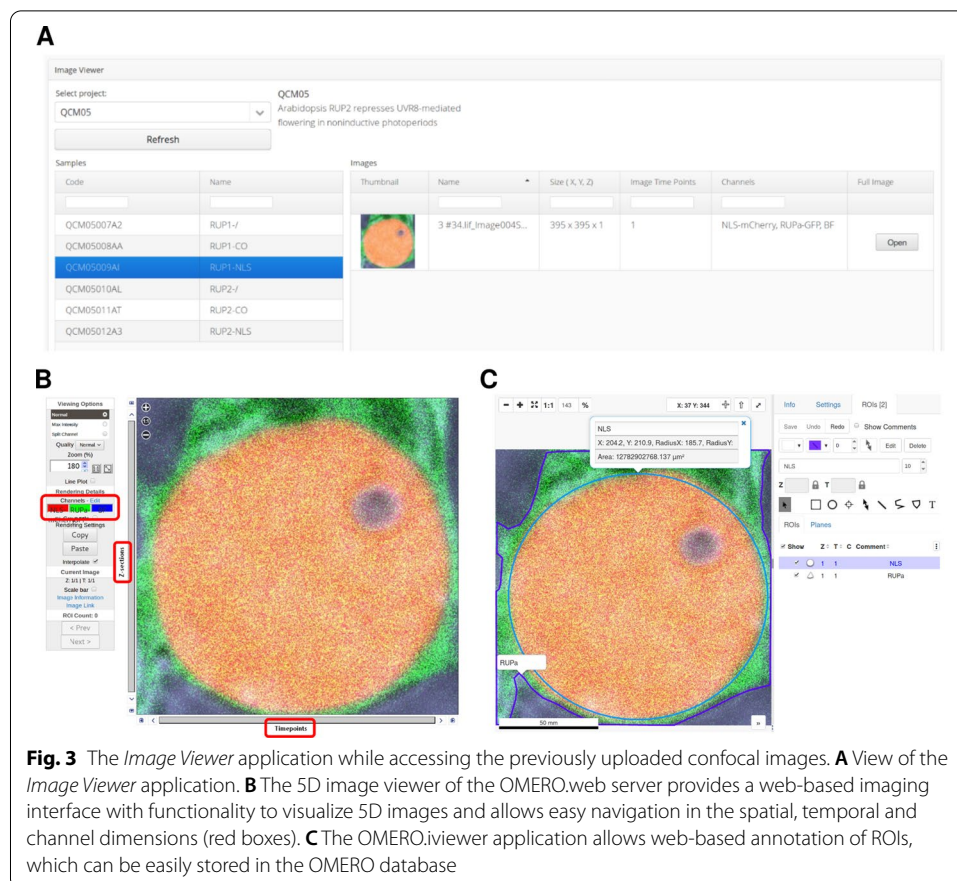
#### ***Proof of concept 1: Arabidopsis RUP2***

This study aimed to improve our molecular understanding of photoperiodic flowering, a light-controlled mechanism in plants, which allows them to optimize the timing of the transition from vegetative growth to flowering, by responding to seasonal changes in day length [18]. An experiment was designed to test the interaction between RUP1, RUP2 (repression of UV-B photomorphogenesis 1 and 2, respectively) and CO (constans transcription factor), two proteins involved in photoperiodic flowering control, using epidermal leaf cells of *Nicotiana benthamiana* that were transiently transformed. The localization of GFP-tagged RUP1 or RUP2 when coexpressed with either CO-mCherry or NLS-mCherry (nuclear localization signal) was investigated with confocal laser scanning microscopy (Fig. 2A).





The fluorescence microscopy data acquired in the abovementioned experiment was registered into our prototype version of qPortal with OMERO support. The initial step in registering this experiment was to create a project using the *Project Wizard* application, and input the experimental design alongside sample information (e.g. species, tissue, experimental condition). The *Project Wizard* will then create the metadata structures in the backend and generate unique identifiers (codes) of each of the biological samples (Fig. 2B). Once a project containing the colocalization experiment was created, it could be accessed using the *Project Browser*, an application that can provide a summarized view of the project or display an in-depth diagram representing the metadata structure and entities (Fig. 2C). Using the sample codes generated by the *Project Wizard*, confocal images of each of the experimental samples in the colocalization analysis were uploaded using the dropbox mechanism (see implementation). Once all data was uploaded, images could be accessed through the *Image Viewer* application, they were



**Fig. 3** The *Image Viewer* application while accessing the previously uploaded confocal images. **A** View of the *Image Viewer* application. **B** The 5D image viewer of the OMERO.web server provides a web-based imaging interface with functionality to visualize 5D images and allows easy navigation in the spatial, temporal and channel dimensions (red boxes). **C** The OMERO.iviewer application allows web-based annotation of ROIs, which can be easily stored in the OMERO database

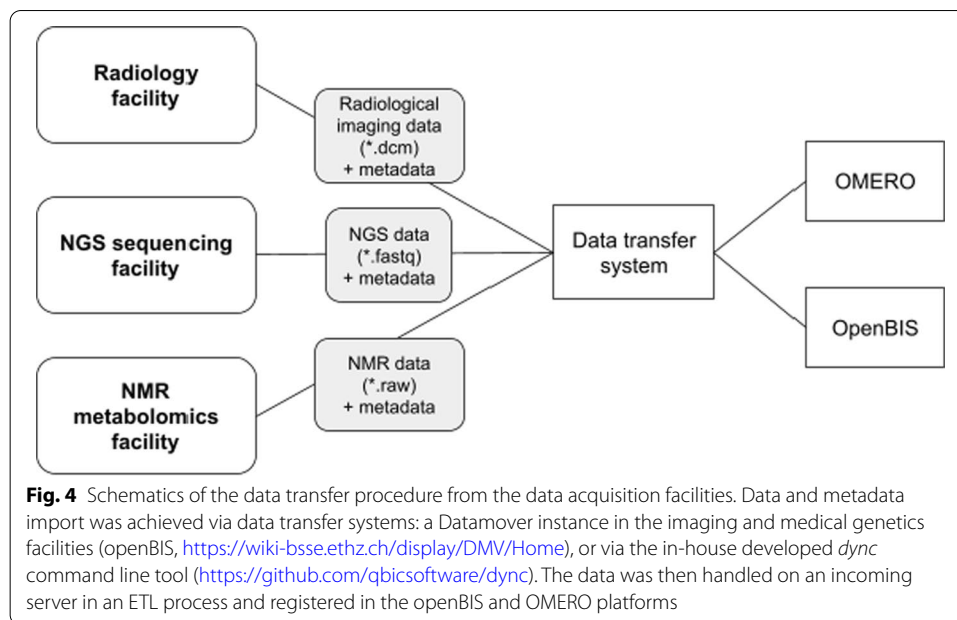
easily findable using sample codes and metadata filters (Fig. 3A), and could be visualized using the 5D image viewer of OMERO.web server (Fig. 3B).

### **Proof of concept 2: synthetic clinical project**

To test the functionality of the proposed architecture with a significantly large and heterogeneous dataset, we created a synthetic clinical project, and registered publicly available medical imaging and genomics sequencing data. Using the *Project Wizard*, an experimental design was defined, comprising 130 patients with three measured biological samples each: abdomen, healthy tissue, and tumor tissue. For each abdomen sample, one X-ray CT from the LiTs dataset [19] was sampled randomly and registered. Similarly, for each liver and tumor tissue sample, 10 randomly selected Hematoxylin and Eosin (H&E) stained tissue images from the MoNuSeg dataset [20] were registered. The publicly available 1000 genomes dataset from the International Genome Sample Resource [21] was employed as a showcase for the registration of genomics data. For each liver and tumor biological sample, a DNA test sample was created and a paired-end illumina whole-genome sequencing dataset was registered. The total size of the project is 34.7 GB, the distribution of file sizes per modality is shown in Additional file 1: Fig. S2A.

We conducted a benchmark test of data registration for the abovementioned synthetic dataset. To evaluate registration of heterogeneous data modalities, we measure the time



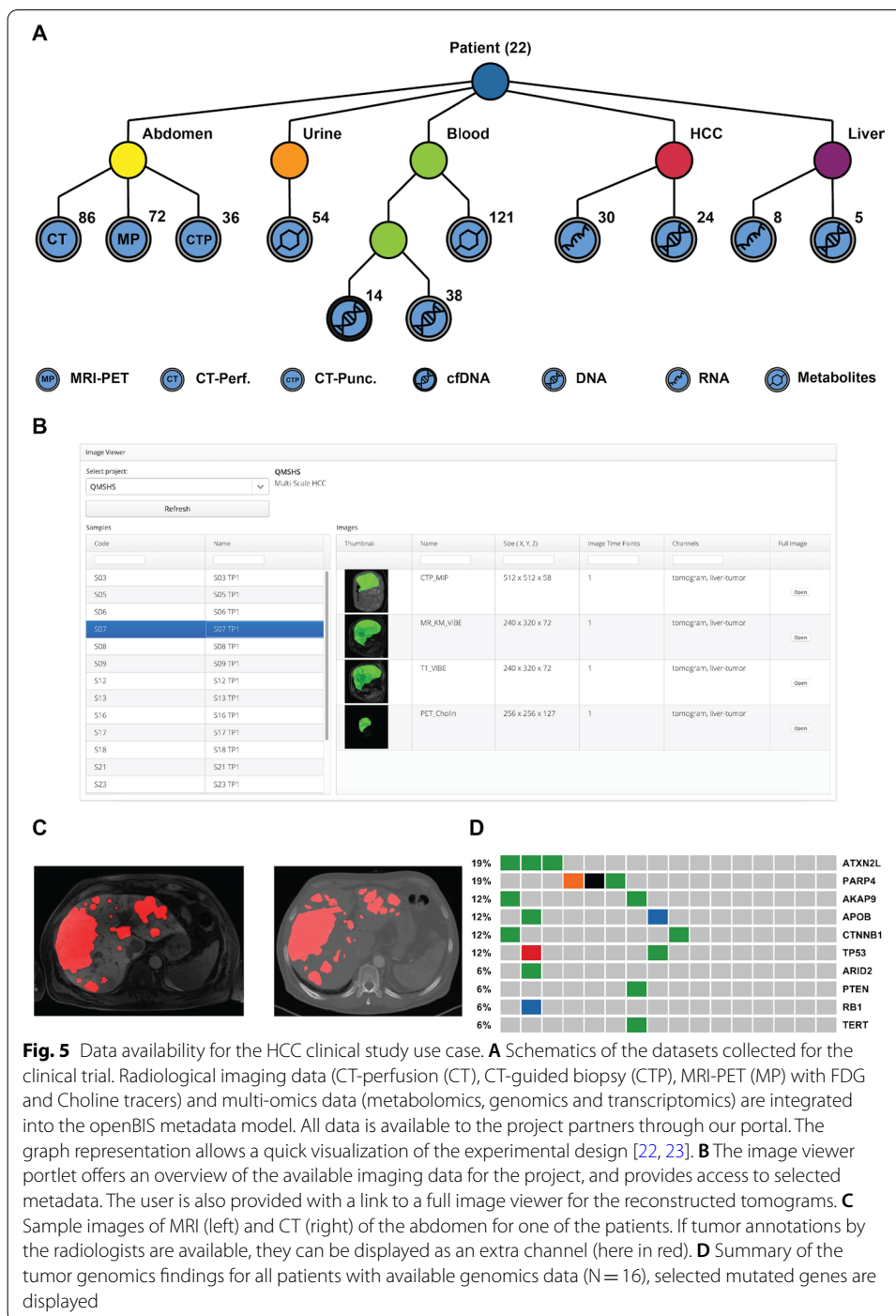


required by the ETL process to sequentially register data for an increasing number of patients, where the data for a single patient is approximately 600 MB in size, consisting of 1 X-ray CT (~250 MB), 20 histology images (~1.75 MB each), and 2 paired-end whole genome sequencing datasets (~600 MB each). The results of the benchmark test show that the time required to register both imaging and sequencing data increased linearly. Registration times are shown in Additional file 1: Fig. S2B.

#### **Use case: hepatocellular carcinoma clinical study**

We recently applied our infrastructure to a clinical study involving patients with Hepatocellular Carcinoma (HCC, NCT02372162, “Fingerprint characterization of advanced HCC”). The study aimed at predicting the patient’s response to the drug sorafenib and identifying the molecular or image determinants for therapy response. The disease progression was followed by medical imaging analysis before therapy, after 4 weeks of therapy and 4–6 months after therapy in case of progression. The imaging data was acquired at the Tübingen University Hospital and included magnetic resonance imaging (MRI) combined with Positron-Emission Tomography with two different tracers (PET), and Computerized Tomography (CT) scans for tumor and liver perfusion analysis. CT-guided liver biopsies were additionally acquired. Molecular genetics and gene expression data profiles were obtained from the tumor biopsies and from the healthy liver tissue at the same time points. Next-generation sequencing (NGS) was performed at the medical genetics facility in Tübingen. Metabolite markers in blood and urine of the patients were measured using Nuclear Magnetic Resonance (NMR) at the Werner Siemens Imaging Center, to follow progression markers of the disease and the drug’s pharmacokinetics. The data generated at the three facilities was integrated into our infrastructure (Fig. 4).

The study generated a total of 3.1 Terabytes of data, comprising NGS (transcriptomics and genomics), radiological imaging, including radiomics analysis, and NMR



metabolomics datasets (Fig. 5A). The different project partners were given access to the complete study data through the qPortal. The Image Viewer portlet provides an overview of the acquired radiomics data, and enables the display of the radiologist annotations corresponding to the tumor regions (Fig. 5B, C). By following the OMERO link, the reconstructed 3D images can be inspected (Fig. 5C). Finally, the radiomics findings were contrasted with molecular information about the tumor, provided by the other omics

modalities stored in our portal. The genetic variations of the patient's tumors were analyzed, in order to identify possible driver genes responsible for tumor development and regression (Fig. 5D).

## Discussion

We presented a data management architecture for life science research capable of handling projects with complex experimental designs, and managing multi-omics data in conjunction with imaging data from the majority of microscopy and medical imaging disciplines (e.g. PET, electron and fluorescence microscopy). The integration of these datasets was achieved by integrating an OMERO server into qPortal, a web-based multi-omics platform built on openBIS. An important feature of this approach is the integrative metadata model, which results from the coupling of the OMERO and openBIS models. This model defines clear metadata boundaries, assigning project, experimental design and omics metadata to the openBIS server, while using OMERO to store metadata that describes available imaging data and its acquisition process. The objective of these metadata boundaries is to take advantage of the strengths of each platform. However, it is important to carefully assign metadata elements to at least one platform, and avoid replicating metadata on both servers, since this may lead to significant complexity in future metadata curation.

A salient feature of the proposed architecture is the SOA-based approach used to implement a loose-coupling between openBIS and OMERO, as modular software components of the qPortal backend. In particular, for defining the metadata connections between the qPortal and OMERO models, which benefited from similarities in the main structure of the models (see Additional file 1: Fig. S1), and allowed their complementary characteristics to be leveraged with few restrictions. Specifically, the OMERO server was treated as a modular component, responsible for storing raw imaging data and metadata related only to imaging datasets and their acquisition process, thereby leveraging OMERO's core functionality while allowing qPortal and openBIS to continue to support in FAIR data management, and to store omics data and metadata related to the experimental design of projects.

In SOA, a variety of specialized and modular components provide services via APIs, allowing them to operate in concert within a distributed system. In this context, qPortal can serve as a suitable distributed framework for data management of multi-omics and imaging data, since it provides additional abstraction layers where the middleware needed to implement the communication protocol between the OMERO and openBIS components can be developed. Moreover, qPortal operates as a web-based platform with a set of applications for scientific data management, and capable of deploying custom Java applications for a particular omics or imaging modality (e.g. the *Image Viewer* application).

An alternative approach to the proposed architecture would be to extend a single platform (e.g. openBIS or OMERO), leading to a monolithic backend system with aggregated functionality. This type of information systems tend to be highly complex and difficult to maintain. On the one hand, when extending qPortal or openBIS to accommodate biomedical imaging data, it is relevant to notice that while qPortal is capable of managing complex multi-omics data, it was not designed to support raw imaging data

and metadata from a large variety of microscopy and medical imaging disciplines without a significant extension to its core functionality, that would likely lead to replication of most of the functionality and the metadata model offered by OMERO, e.g. an image rendering engine, proprietary file format parsers, and database support for ROIs (see [The omero platform](#) section).

On the other hand, while OMERO offers an extensible metadata model and allows customization to support omics data, as demonstrated by an extension to accommodate data from a Genome Wide Association Study (GWAS) of human autoimmune diseases [8, 24], it is still unclear if complex and highly heterogeneous multi-omics experiments can be stored and queried using the OMERO server alone. Extending the OMERO server to enforce FAIR-supporting, multi-factorial experimental designs, storage and annotation of multiple, high-throughput omics disciplines, would likewise require significant software development and long-term support, that would also replicate functionality already offered by openBIS and qPortal.

Therefore we opted for a SOA approach to integrate management of omics and biomedical imaging data, allowing us to leverage several technologies in a modular way. By abstracting the openBIS and OMERO servers as backend service components within a distributed system (qPortal), it was possible to selectively use the metadata and raw data storage facilities of each component during the execution of the required data management operations, while maintaining the integrity of metadata structures across all backend components.

We have also presented two technical proofs of concept and a use case where the proposed qPortal with imaging support was used to integrate imaging data and metadata into a multi-omics environment. The Arabidopsis RUP proof of concept described the registration process of a plant biology project containing data from a confocal scanning microscope, and showed how such imaging data could be easily accessed via the web interface of qPortal. While this project did not contain omics data, it was necessary proof of concept to introduce microscopy data and metadata into an otherwise omics-only platform. The second proof of concept is based on a large synthetic dataset, containing X-ray CT, histology images, and NGS data. Using this synthetic project we showed how the proposed platform can be applied to large and heterogeneous datasets. Additionally, we conducted a benchmark test of data registration, where we show that registration time for the synthetic dataset increased linearly with respect to the amount of data. The presented use case shows a hepatocellular carcinoma study with highly heterogeneous medical imaging and multi-omics data, allowing physicians to relate several disciplines (e.g. NGS, gene expression data, confocal, PET and MRI imaging data) via the metadata structure of the project, which relates biological samples entities (and all their associated datasets) with the corresponding patient and time point of the clinical study.

Our data management architecture allows the deployment of data repositories capable of linking omics and biomedical imaging modalities. Leveraging the interoperability of backend components and storing metadata in an integrative and standardized manner, allows searching across studies to derive datasets of integrated omics and biomedical imaging data, thus facilitating multi-modal data reusability. Such datasets are particularly attractive given the complementary nature of omics and biomedical imaging modalities, since omics disciplines enable the simultaneous measurement of a large

variety of molecular components and species in vitro, but lack the spatio-temporal information provided by imaging modalities. Moreover, since significant breakthrough in microscopy techniques permit the visualization of cellular structures in situ, at ever increasing resolutions (e.g. super-resolution light microscopy [1], direct electron detectors [2]), and both omics and microscopy modalities have reached single-cell sensitivity, the need to correlate such datasets has significantly increased [15].

Joint analysis of imaging and omics data can provide new insights into biological processes. Multi-modal data analysis can be accomplished by independently analyzing data from different modalities, that has been extracted from the same biological sample, ideally using reproducible bioinformatics workflows [25]. Integrated multi-modal analysis is also possible by using higher-order data representations or by employing deep learning methods [15]. Importantly, supervised machine learning methods, such as deep neural networks, rely on the availability of large training datasets annotated with rich and high quality phenotypic metadata, underscoring the importance of multi-modal and FAIR-supporting data repositories.

## Conclusion

Here we present a SOA-based method to integrate an image management system, the OMERO server, as a modular, backend component of qPortal, to allow the integrated management and analysis of multi-omics and biomedical imaging data. The implementation of a structural coupling between the openBIS and OMERO metadata models, the development of software components to facilitate the communication with the OMERO server, and an extension to the data management operations of qPortal, facilitated storage and analysis of raw data and metadata from various omics, microscopy and biomedical imaging modalities in an integrative manner, with the ability of accepting metadata queries from web-based, scientific applications. The applicability of the proposed architecture was demonstrated with two proofs of concept and a use case, a plant biology study using confocal scanning microscopy, a synthetic project with a large and heterogeneous dataset (medical imaging and sequencing data), and a clinical study on hepatocellular carcinoma, with data from heterogeneous medical imaging and omics modalities.

As emerging developments in omics and biomedical imaging drive the increase in resolution, modality, and throughput of data generation in life science studies, following SOA principles in extending, integrating, or building the required information systems for long-term storage and FAIR-driven management of these valuable digital assets, offers a feasible approach to manage a variety of multi-omics and imaging data repositories. Moreover, given the success of SOA-oriented information systems in other domains, we strongly believe that similar architectures will play an important role in implementing management systems for heterogeneous scientific datasets.

Such FAIR data management infrastructures for data repositories, capable of tackling the ever increasing volume and complexity of omics and biomedical will not only enable high-throughput, multi-modal data management in life science, but also allow the generation of large and highly multi-dimensional datasets, via metadata queries across stored scientific projects, enabling data-driven analysis and training of machine learning models in future predictive applications.

## Methods

Our implementation leverages the openBIS, qPortal, and OMERO platforms to provide data management for multi-modal data in life science research. Specifically, it allows management and analysis of omics data (e.g. gene sequencing and expression data, mass spectrometry) in conjunction with data from microscopy and medical imaging disciplines. Integration of omics and imaging data modalities at the metadata level is achieved by a combination of well-established metadata models (i.e. the qPortal and OMERO models), and implemented through middleware software that provides the logic, and connective tissue between the application program interface (API) of the aforementioned platforms, allowing them to work in concert as components of a unified system.

## Architecture

Our data management architecture integrates the OMERO platform into qPortal, it leverages openBIS and OMERO functionality to facilitate the integration of both imaging and omics data within qPortal. openBIS serves the purpose of storing and managing raw data and metadata from omics sources, including the general experimental design of research projects and basic information regarding the biology of the samples. On the other hand, OMERO is best suited for modeling, storing and managing medical imaging and microscopy data, employing a metadata model focused on describing multi-channel spatio-temporal data, the technical specification of imaging instruments, and the experimental parameters used during data acquisition.

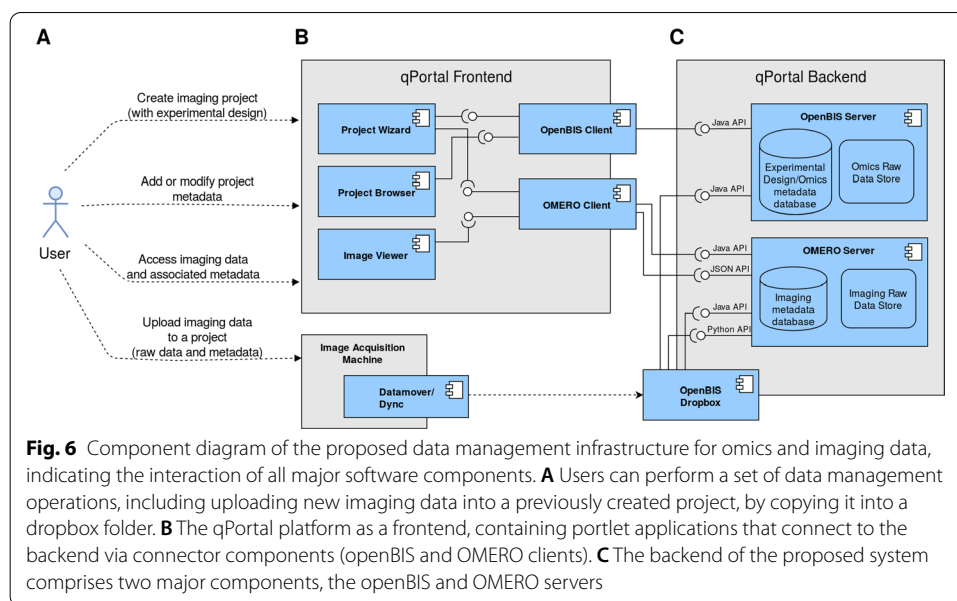
While the openBIS and OMERO servers function as the backend of the system, qPortal provides a web-based user interface with a set of applications for scientific data management. These applications support various project management operations and orchestrate the complex data management tasks by using utility components designed to facilitate a connection to the openBIS and OMERO servers for metadata transfer. Figure 6 depicts the component diagram of the proposed architecture.

Raw imaging data and accompanying metadata can be uploaded into the system using an openBIS-based *dropbox* mechanism. In short, data is securely transferred to an *incoming server* within the backend of the proposed system, a process is then triggered to allow for structured ingestion of raw data and metadata. This process uploads images to the OMERO server and stores the resulting unique identifiers in the openBIS server as metadata items, in order to maintain a consistent record in both servers.

## The OMERO platform

The OME Remote Objects (OMERO) is a data management platform for biological imaging that is based on the OME Data Model [26] and the Bio-Formats project [17]. It is composed of databases, middleware and remote client applications. The main component is the OMERO.server, a Java application that acts as middleware connecting various databases that store different data types, and provides access to stored data via a single API. The OMERO.server not only provides access to the underlying storage facilities, it must also process data before delivering it to the client applications.





OMERO provides a Java-based client for imaging data upload, the OMERO.importer. This tool is a command line interface that uses Bio-Formats to read image data and metadata from supported file formats. During data import, Bio-Formats reads metadata from the raw data file and maps it to the OME Data Model [26]. While a binary repository is used to store binary data such as images and thumbnails, the OMERO relational database archives all metadata (see [Metadata models](#)) associated with binary images, user information, and simple data annotations.

Since modern biological images may consist of various image frames recorded at different positions, channels, or timepoints, the OMERO.server contains a multi-threaded image rendering engine that can rapidly display image planes, and transfer them via the API to clients. This engine reads images from the binary repository and can apply transforms, according to the parameters provided by an OMERO client or the OMERO relational database. Supported operations include image compression, overlay and projection. Users can therefore create multiple views of the data without modifying the originally acquired data.

While data can be imported using the provided OMERO.importer or OMERO.dropbox tools (a filesystem monitoring tool), third-parties can also develop specialized tools to access metadata and raw imaging data by querying the Java or Python APIs of the OMERO.server. Remote data access between client applications and the OMERO.server is achieved via ZeroC's Internet Communication Engine [27].

OMERO also includes a customizable web-based client for image visualization and annotation. OMERO.web is a Django-based application that uses the APIs of the OMERO.server to provide a web interface for metadata access, metadata annotation, and full 5D (i.e. space, time, and channels) image visualization.

Regions of interest (ROI) describing the spatial boundaries of detected objects in an image, which are often the product of manual annotation or segmentation algorithms,

can be stored in the OMERO database as geometric objects (e.g. points, circles, polygons), thus supporting ROI storage for 5D images.

### The OpenBIS platform

The openBIS platform is an open source management system for data acquired in biological experiments [9]. The main components of this platform are a data store, a flexible and extendable metadata model supporting complex XML-based annotation, a relational database to store metadata, and an application server to browse, access and manage both data and metadata.

In order to upload data into openBIS, a *dropbox* is usually employed. In this *context*, a *dropbox* is a filesystem monitoring mechanism that can transfer files containing raw data and metadata, from an acquisition machine or adjacent computer within a source laboratory, to a server in the backend of the data management system, which temporarily stores *incoming* data. Usually, this data transfer operation is carried out securely using the openBIS Datamover or the Dync application [28].

Subsequently, this *incoming* backend server executes the *extract, transform, load* (ETL) routine associated with the specific source laboratory and raw data type, the objective of this routine is to process raw data before uploading it into the openBIS data store. ETL routines can execute additional external scripts and often extract metadata from raw data files. During the ETL process all the necessary metadata entities (e.g. experiments, samples, datasets) are created, and all available metadata is stored as properties of the newly created entities.

### The web-based qPortal

qPortal is a web-based platform for FAIR-driven data management of biomedical data. It enforces a well-structured metadata model to capture the experimental design of projects and the biology of the samples. This platform is built on a Liferay instance and contains a large collection of loosely coupled portlets (i.e. java-based, web applications) that use the open-source framework VAADIN [29]. qPortal also offers integrated workflow support on stored data. The main portlet applications, the *Project Wizard* and *Project Browser*, facilitate the creation of projects with complex experimental designs, and provide data access and management, respectively.

The registration of an experimental design is facilitated with the *Project Wizard*. The application helps users to register a new project by taking them through a series of defined steps to describe the experimental design and metadata of the project. The first step describes the biological entities under study (e.g. patients, model organism) on species level, the second step captures the extraction of cells or tissues from biological entities. Finally, the third step describes the process used to prepare the biological sample for data acquisition. Experimental factors such as treatment or genotype can also be defined precisely and in order to maximize statistical power a full-factorial design is proposed by default. Once all this information has been collected, the *Project Wizard* creates and stores all the metadata entities and relations in the openBIS database, according to a well-defined metadata model (see [Metadata models](#) section).

The main interaction with project-associated data is facilitated by the *Project Browser*. This application allows search and access to project metadata, as well as

providing direct access to raw data and analysis results. This portlet shows all projects the logged-in user has access to, in a searchable table that provides general project information, e.g. project code and description. By clicking on a project, the user can access all project metadata, which the application visualizes as a tree-like structure of metadata entities (Additional file 1: Fig. S1A and metadata models). Metadata concerning the biology of the samples can be accessed, and the associated raw data can be directly downloaded. Data resulting from workflow executions can also be accessed from a dedicated section within the application.

As part of the qPortal backend, an openBIS server instance provides the means of storing and managing raw data and metadata. While the *Project Wizard* application only registers the main information of the experiments and project, the openBIS-based ETL routines register the remainder of the experimental metadata concerning the file (e.g. sample preparation and data acquisition parameters) during data upload. Effectively, the *Project Wizard* and ETL scripts enforce an homogenous metadata structure for all projects.

qPortal provides an interface to allow execution of computations in cluster infrastructures using the data stored in openBIS, facilitating the automation of common bioinformatics analyses. Once a workflow has been executed, qPortal can register the resulting data into openBIS with a metadata reference to the original input data and the parameters used during analysis, to allow reproducibility of workflow results. Moreover, workflow systems such as gUSE [30] and Nextflow [31] have been successfully connected to the qPortal platform. In particular, qPortal is compatible with the nf-core framework [25], which provides access to standardized and reproducible bioinformatics pipelines, along with the necessary pipeline development tools.

### Abbreviations

API: Application programming interface; CO: Constants transcription factor; CRUD: Create, read, update, delete; CT: Computer tomography; ETL: Extract, transform, load; FAIR: Findability, accessibility, interoperability, and reusability; FDG: Fluorodeoxyglucose; GFP: Green fluorescent protein; GWAS: Genome wide association study; H&E: Hematoxylin and eosin; HCC: Hepatocellular carcinoma; HTTP: Hypertext transfer protocol; MRI: Magnetic resonance imaging; NGS: Next-generation sequencing; NLS: Nuclear localization sequence; NMR: Nuclear magnetic resonance; OME: Open Microscopy Environment; OMERO: Open Microscopy Environment Remote Objects; openBIS: Open Biological Information System; PET: Positron-Emission tomography; ROI: Region of interest; RUP: Repression of UV-B photomorphogenesis; RUP1: Repression of UV-B photomorphogenesis 1; RUP2: Repression of UV-B photomorphogenesis 2; SOA: Service oriented architecture; SSO: Single-sign-on; UV-B: Ultraviolet B radiation; XML: Extensible markup language; 3D: 3 Dimensional; 5D: 5 Dimensional.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-022-04584-3>.

**Additional file 1: Fig. S1.** Diagram of the unified metadata model, depicting the coupling between the underlying models and the cardinality relationship between metadata entities. **(A)** The hierarchical metadata model used by qPortal to describe the experimental design of a research project, containing general information of sample biology. **(B)** The hierarchical metadata model used by OMERO, focused on describing imaging data. **Fig. S2.** The file size distribution of the synthetic project and benchmark test. **(A)** The distribution of file sizes for the synthetic project per modality. The X-ray CT data was sampled from the LiTS [19] training dataset (130 tomograms), the H&E stained histology images were sampled from the training dataset of MoNuSeg [20] (30 images), and the genomics dataset was obtained from the 1000 genomes project [21] (260 paired-end illumina whole genome sequencing datasets). **(B)** Results of the benchmarking test for data registration, showing sequential registration times for imaging data (x-ray CT and histology images) and genomics sequencing data (WGS) for an increasing number of patients, where data for a single patient consists of 1 x-ray CT, 20 histology images, and 4 fastq files (paired-end data for cancer and healthy tissue, respectively). The curves follow the mean values of 4 registration runs (black dots).

### Acknowledgements

We acknowledge support from the Open Access Publishing Fund of the University of Tübingen.

### Authors' contributions

LKC and AF did the research, implemented the software, metadata model, and wrote the manuscript. SF supervised the software implementation and contributed in writing the manuscript. LG, TK and AS contributed with software development and writing of the manuscript. GG and WT analyzed the clinical data and contributed in writing the manuscript. SZOK, FW and SR conducted the plant biology study, analyzed the plant biology data, tested the proposed platform, and contributed in writing the manuscript. NM and MB supervised the clinical study and reviewed the manuscript. MH conducted the clinical study and reviewed the manuscript. KH supervised the plant biology study and reviewed the manuscript. SN suggested the study, supervised the overall work and wrote the manuscript. All authors have read and approved the manuscript.

### Funding

Open Access funding enabled and organized by Projekt DEAL. Funding from all sources was pivotal to conduct this study. Here we outline the individual funding sources that contributed to the study. We acknowledge funding from BMBF Multiscale-HCC, DFG SFB/TR 209, DFG SFB 1101, DFG SFB/TR 261, Excellence cluster microbiology. S.N. and L.K.C. acknowledge funding from SFB 261 and FW, SzO-K. L.K.C. and K.H. acknowledge funding from SFB 1101 (projects D02 and Z02). G.G. acknowledges funding from the German Ministry of Research and Education (BMBF, Grant No. 01ZX1301F). S.N. acknowledges funding from Deutsche Forschungsgemeinschaft (core facilities initiative, KO-2313/6-1 and KO-2313-2, Institutional Strategy of the University of Tübingen, ZUK 63). Furthermore, S.N. acknowledges funding by the Sonderforschungsbereich SFB/TR 209 "Liver cancer" of the Deutsche Forschungsgemeinschaft (DFG), as well as from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)- Project-ID 398967434—TRR 261. We acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy—EXC 2124—390838134 and by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy—EXC 2180—390900677. The clinical study (NCT02372162) was funded by the German Ministry for Education and Research (BMBF, eMed/Multiscale HCC, FKZ 01ZX1301A und 01ZX1601G, N.P.M., M.B.). Funding bodies did not play any role in the design of this study and in writing of the manuscript.

### Availability of data and materials

Project name: OMERO Client Module.

Project home page: <https://github.com/qbicsoftware/omero-client-lib>.

Archived version: <https://doi.org/10.5281/zenodo.4067716>.

Operating system(s): Platform independent.

Programming language: Java.

Other requirements: Java 1.8 or higher, Tomcat 4.0 or higher.

License: MIT License.

Project name: Project Wizard Application.

Project home page: <https://github.com/qbicsoftware/projectwizard-portlet>.

Archived version: <https://doi.org/10.5281/zenodo.3908302>.

Operating system(s): Platform independent.

Programming language: Java.

Other requirements: Java 1.8 or higher, Tomcat 4.0 or higher.

License: MIT License.

Project name: Image Viewer Application.

Project home page: <https://github.com/qbicsoftware/omero-client-portlet>.

Archived version: <https://doi.org/10.5281/zenodo.4068252>.

Operating system(s): Platform independent.

Programming language: Java.

Other requirements: Java 1.8 or higher, Tomcat 4.0 or higher.

License: MIT License.

Project name: ETL Scripts.

Project home page: <https://github.com/qbicsoftware/etl-scripts/>.

Archived version: <https://doi.org/10.5281/zenodo.4085722>.

Operating system(s): Platform independent.

Programming language: Jython, Python, Java.

Other requirements: qPortal instance, provided OMERO conda environment.

License: MIT License.

### Declarations

#### Ethics approval and consent to participate

For the illustrated use case on liver cancer, we received ethical approval by the ethical commission of the University Hospital of Tübingen.

#### Consent for publication

Not applicable.

**Competing interests**

The authors declare that they have no competing interests.

**Author details**

<sup>1</sup>Quantitative Biology Center (QBiC), University of Tübingen, Tübingen, Germany. <sup>2</sup>Center for Plant Molecular Biology (ZMBP), University of Tübingen, Tübingen, Germany. <sup>3</sup>Department of Radiology, Diagnostic and Interventional Radiology, University of Tübingen, Tübingen, Germany. <sup>4</sup>Department Internal Medicine I, University of Tübingen, Tübingen, Germany. <sup>5</sup>Biomedical Data Science, Department of Computer Science, University of Tübingen, Tübingen, Germany.

Received: 26 June 2020 Accepted: 21 January 2022

Published online: 07 February 2022

**References**

1. Sigal YM, Zhou R, Zhuang X. Visualizing and discovering cellular structures with super-resolution microscopy. *Science*. 2018;361(6405):880–7.
2. Cheng Y. Single-particle Cryo-EM at crystallographic resolution. *Cell*. 2015;161(3):450–7.
3. Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A, et al. The FAIR guiding principles for scientific data management and stewardship. *Sci Data*. 2016;3:160018.
4. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Čech M, et al. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res*. 2016;44:W3–10. <https://doi.org/10.1093/nar/gkw343>.
5. Goecks J, Nekrutenko A, Taylor J, Galaxy Team. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol*. 2010;11(8):R86.
6. Cerami E, Gao J, Dogrusoz U, Gross BE, Sumer SO, Aksoy BA, et al. The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer Discov*. 2012;2(5):401–4.
7. Kunszt P, Blum L, Hullár B, Schmid E, Srebniak A, Wolski W, et al. iPortal: the Swiss grid proteomics portal: requirements and new features based on experience and usability considerations. *Concurr Comput Pract Exp*. 2015;27(2):433–45.
8. Allan C, Burel J-M, Moore J, Blackburn C, Linkert M, Loynton S, et al. OMERO: flexible, model-driven data management for experimental biology. *Nat Methods*. 2012;9(3):245–53.
9. Bauch A, Adamczyk I, Buczek P, Elmer F-J, Enimanev K, Glyzowski P, et al. openBIS: a flexible framework for managing and analyzing complex data in biology research. *BMC Bioinformatics*. 2011;12:468.
10. de Chaumont F, de Chaumont F, Dallongeville S, Chenouard N, Hervé N, Pop S, et al. Icy: an open bioimage informatics platform for extended reproducible research. *Nat Methods*. 2012;9:690–6. <https://doi.org/10.1038/nmeth.2075>.
11. Baker M. 1,500 scientists lift the lid on reproducibility. *Nature*. 2016;533(7604):452–4.
12. Mohr C, Friedrich A, Wojnar D, Kenar E, Polatkan AC, Codrea MC, et al. qPortal: a platform for data-driven biomedical research. *PLoS ONE*. 2018;13(1):e0191603.
13. Josuttis NM. SOA in practice: the art of distributed system design. O'Reilly Media Inc; 2007.
14. Disselhorst JA, Krueger MA, Ud-Dean SMM, Bezrukov I, Jarbouli MA, Trautwein C, et al. Linking imaging to omics utilizing image-guided tissue extraction. *Proc Natl Acad Sci USA*. 2018;115(13):E2980–7.
15. Hériché J-K, Alexander S, Ellenberg J. Integrating imaging and omics: computational methods and challenges. *Annu Rev Biomed Data*. 2019. <https://doi.org/10.1146/annurev-biodatasci-080917-013328>.
16. Stoyanova R, Takhar M, Tschudi Y, Ford JC, Solórzano G, Erho N, et al. Prostate cancer radiomics and the promise of radiogenomics. *Transl Cancer Res*. 2016;5(4):432–47.
17. Linkert M, Rueden CT, Allan C, Burel J-M, Moore W, Patterson A, et al. Metadata matters: access to image data in the real world. *J Cell Biol*. 2010;189(5):777–82.
18. Arongaus AB, Chen S, Pireyre M, Glöckner N, Galvão VC, Albert A, et al. Arabidopsis RUP2 represses UVR8-mediated flowering in noninductive photoperiods. *Genes Dev*. 2018;32(19–20):1332–43.
19. Bilic P, Christ PF, Vorontsov E, Chlebus G, Chen H, Dou Q, Fu C-W, Han X, Heng P-A, Hesser J, Kadoury S, Konopczynski T, Le M, Li C, Li X, Lipková J, Lowengrub J, Meine H, Moltz JH, et al. The liver tumor segmentation benchmark (LiTS). In: arXiv [cs.CV]. 2019. arXiv: <http://arxiv.org/abs/1901.04056>.
20. Kumar N, Verma R, Sharma S, Bhargava S, Vahadane A, Sethi A. A dataset and a technique for generalized nuclear segmentation for computational pathology. *IEEE Trans Med Imaging*. 2017;36(7):1550–60.
21. Fairley S, Lowy-Gallego E, Perry E, Flicek P. The International Genome Sample Resource (IGSR) collection of open human genomic variation resources. *Nucleic Acids Res*. 2020;48(D1):D941–7.
22. Friedrich A, de la Garza L, Kohlbacher O, Nahnsen S. Interactive Visualization for Large-Scale Multi-factorial Research Designs. *Data Integr Life Sci* 2018;75–84.
23. Friedrich A, Kenar E, Kohlbacher O, Nahnsen S. Intuitive web-based experimental design for high-throughput biomedical data. *BioMed Res Int*. 2015;2015:958302.
24. Sanna S, Pitzalis M, Zoledziewska M, Zara I, Sidore C, Murru R, et al. Variants within the immunoregulatory CBLB gene are associated with multiple sclerosis. *Nat Genet*. 2010;42(6):495–7.
25. Ewels PA, Peltzer A, Fillinger S, Patel H, Alneberg J, Wilm A, Garcia MU, Di Tommaso P, Nahnsen S. The nf-core framework for community-curated bioinformatics pipelines. *Nat Biotechnol*. 2020;38(3):276–8.
26. Goldberg IG, Allan C, Burel J-M, Creager D, Falconi A, Hochheiser H, et al. The Open Microscopy Environment (OME) Data Model and XML file: open tools for informatics and quantitative analysis in biological imaging. *Genome Biol*. 2005;6(5):R47.
27. Henning M. A new approach to object-oriented middleware. *IEEE Internet Comput*. 2004;8:66–75. <https://doi.org/10.1109/mic.2004.1260706>.
28. Seyboldt A, Fillinger S. qbicssoftware/dync: first stable release (version 1.0.0). Zenodo; 2019. 10.5281/zenodo.3515438
29. Duarte A. Vaadin 7 UI design by example: beginner's guide. Packt Publishing Ltd; 2013.

30. Kacsuk P, Farkas Z, Kozlovsky M, Hermann G, Balasko A, Karoczkai K, Marton I. WS-PGRADE/gUSE generic DCI gateway framework for a large variety of user communities. *Int J Grid Util Comput*. 2012;10(4):601–30.
31. Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow enables reproducible computational workflows. *Nat Biotechnol*. 2017;35(4):316–9.

### **Publisher's Note**

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

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

