**PAPER • OPEN ACCESS**

# Deploying the Big Data Science Center at the Shanghai Synchrotron Radiation Facility: the first superfacility platform in China

To cite this article: Chunpeng Wang *et al* 2021 *Mach. Learn.: Sci. Technol.* **2** 035003

View the article online for updates and enhancements.

You may also like

- Physics design of SSRF synchrotron radiation security
  Xu Yi, Dai Zhi-Min and Liu Gui-Min

- Excitation curve calibration for the SSRF magnet system
  Zhou Xue-Mei, Liu Gui-Min, Hou Jie et al.

- Application of real-time digitization techniques in beam measurement for accelerators
  Lei Zhao, , Lin-Song Zhan et al.

MACHINE
LEARNING
Science and Technology

**PAPER**

# Deploying the Big Data Science Center at the Shanghai Synchrotron Radiation Facility: the first superfacility platform in China

Chunpeng Wang[1,6] ⓘ, Feng Yu[2,6] ⓘ, Yiyang Liu[3,4] ⓘ, Xiaoyun Li[1] ⓘ, Jige Chen[1] ⓘ, Jeyan Thiyagalingam[5] ⓘ and Alessandro Sepe[1,7] ⓘ

[1] Big Data Science Center, Shanghai Synchrotron Radiation Facility, Shanghai Advanced Research Institute, Chinese Academy of Sciences, 239 Zhangheng Road, Shanghai, People's Republic of China
[2] Division of Life Science, Shanghai Synchrotron Radiation Facility, Shanghai Advanced Research Institute, Chinese Academy of Sciences, 239 Zhangheng Road, Shanghai, People's Republic of China
[3] Shanghai StartComputing Cloud Technology Co. Ltd, 58 JinQiao Road, Shanghai, People's Republic of China
[4] Shanghai Supercomputer Center, 585 Guoshoujing Road, Zhangjiang, Shanghai, People's Republic of China
[5] Science and Technology Facilities Council, Rutherford Appleton Laboratory, Harwell Campus, Didcot OX11 0QX, United Kingdom
[6] Contributed equally to this work.
[7] Author to whom any correspondence should be addressed.

**E-mail:** alessandro.sepe@zjlab.org.cn

## Abstract

With recent technological advances, large-scale experimental facilities generate huge datasets, into the petabyte range, every year, thereby creating the Big Data deluge effect. Data management, including the collection, management, and curation of these large datasets, is a significantly intensive precursor step in relation to the data analysis that underpins scientific investigations. The rise of artificial intelligence (AI), machine learning (ML), and robotic automation has changed the landscape for experimental facilities, producing a paradigm shift in how different datasets are leveraged for improved intelligence, operation, and data analysis. Therefore, such facilities, known as superfacilities, which fully enable user science while addressing the challenges of the Big Data deluge, are critical for the scientific community. In this work, we discuss the process of setting up the Big Data Science Center within the Shanghai Synchrotron Radiation Facility (SSRF), China's first superfacility. We provide details of our initiatives for enabling user science at SSRF, with particular consideration given to recent developments in AI, ML, and robotic automation.

## 1. Introduction

Large-scale experimental facilities, such as synchrotron, neutron, and x-ray free-electron laser (XFEL) facilities, are central to the study of fundamental science, particularly in terms of understanding the structural properties of materials, down to their atomic level, and linking them to their corresponding material functionalities. With recent technological advances, these large-scale experimental facilities are now generating massive datasets, which is a significantly rigorous process, and serves as a precursor for the data analysis underpinning scientific investigations [1, 2].

Although the notion of creating Big Data centers and deploying data management policies is an important aspect of modern experimental facilities, the overall set of challenges extend beyond the issues of simple data management [3–5].

More specifically, the rise of artificial intelligence (AI), machine learning (ML), and robotic automation has changed the landscape of experimental facilities. The deep learning revolution [6] has been fundamental to this paradigm shift, in terms of how different datasets are leveraged for improved intelligence, operation, and data analysis.

The following are some of the key challenges common to experimental facilities, particularly in light of current technological developments:

(a) Complex multimodal datasets: Modern scientific facilities can generate very high-resolution datasets, both in the temporal and spatial domains, and often in multiple modalities. The challenges stemming from complex, multi-format datasets are not just about data management, but in fact cover the whole spectrum of experimental configurations for data analysis. Challenges arising from these complex datasets vary from optimally fusing these multimodal datasets to maximize information extraction, to ensuring that datasets with multiple formats can be ingested into data-processing pipelines. More specifically, ensuring interoperability, and a common, unified, universal, and taggable format [7–9] is a serious endeavor.

(b) Broad range of experimental configurations: Although it may be easier to perform experiments using a configuration that would gather as much as data as possible at the highest possible resolution, if the resulting processing is going to exploit such configurations, this may constitute a waste of significant extant resources within experimental facilities.

(c) Complex, AI-enabled processing pipelines: With respect to the potential of AI and ML, these must be leveraged to improve the capability of research facilities. As such, AI and ML are now becoming an integral part of every aspect of experimental facilities. Although this may appear to be a trivial issue of leveraging AI and ML across a range of problems, the key issue here is interoperability. For instance, ensuring that different AI/ML solutions can coexist and operate toward a common solution is a difficult problem.

(d) Demand for high throughput: AI and ML have revolutionized scientific data analysis, resulting in a quick return for scientists with regard to further experiments, or even on-the-fly adaptation of experiments. This results in a huge demand for facility usage; as such, modern data facilities must be fast enough to increase this throughput.

(e) Data-scalable data processing techniques: With increased volumes of data, experimental facilities must contain data processing algorithms that can be scaled in terms of computation power as well as data volume.

(f) Generic AI/ML challenges: In large experimental facilities, unifying AI/ML-ready solutions, such that they are general and transferrable, represents a serious challenge.

(g) Bleeding edge: The overall complexity of scientific problems often surpasses that of non-scientific cases, owing to the complex nature of scientific research and its relevant datasets [10, 11]. Usually, the overall knowledge or information to be maximized is not known *a priori*, including potential data sources. More specifically, given that even the most recent knowledge or information remains incomplete, the aforementioned challenges must be addressed continuously [12]. It is customary to use multiple models to explain the same results, thereby creating the overall complexity characteristic of Big Scientific Data [13].

These challenges indicate that the various stages constituting scientific investigations as a whole, from experiments to data management to data analysis to interpretation, can no longer be performed in isolation. For instance, knowing the available techniques for data analysis may help in tuning the experimental settings. This is a challenge that must be addressed by the scientific community so that scientists can focus on science. In other words, facilities that enable 'user science' are critical for the scientific community. In this paper, we use the term 'Superfacility' to refer to a facility that meets all of the aforementioned demands.

A superfacility enables facility users to fully interpret large and complex datasets by adopting a novel approach, deploying advanced data processing and networking infrastructures at beamlines via the active cooperation of multidisciplinary groups. It aims to innovate in areas such as data curation, movement, annotation, and storage pipelines, to equip itself with novel mathematics, algorithms, and interactive visualization, and to augment data interpretation capabilities for users at large scientific facilities. A noteworthy characteristic of the superfacility involves the integration of all components of a large scientific facility within a unified and centralized platform, ranging from accelerators to instruments and beamlines, seamlessly and elastically connecting them by means of high-performance computing (HPC) facilities, equipped with the most advanced modeling, simulation, analysis, and visualization tools. Therefore, the conversion of light source facilities into superfacilities will dramatically extend the user base at large scientific facilities to the most diverse scientific disciplines, allowing users from a plethora of different disciplines to comprehensively utilize the very valuable data produced at these facilities [13–22].

The concept of the superfacility is gaining momentum globally, as demonstrated by the recent creation of a novel organization, i.e., the Center for Advanced Mathematics for Energy Research Applications (CAMERA); this is an umbrella organization, encompassing and coordinating the most diverse groups, from applied mathematicians to computer scientists, beamline scientists, materials scientists, and computational chemists, to address all the challenges associated with the creation of superfacilities [23–25]. Moreover, the

U.S. National Energy Research Scientific Computing Center Superfacility project has proposed a superfacility model detailing its mission and services, including data management, sharing, transfer, and discovery capabilities [26]. The Pacific Northwest National Laboratory further contributed to the superfacility vision by developing an automated analysis and Big Data analytics pipeline, equipped with an adaptive environment; this enables augmented experimental control based on data evaluation as it is produced by the experiments [27]. The Advanced Photon Source has also proposed several scientific data standards, forming the base of a superfacility, and is investigating their usage and requirements for different types of data analysis usage [5]. Furthermore, the superfacility is of particular relevance in terms of user interaction with synchrotrons, where the community is witnessing an increase in data rates and volume, further triggered by the development of new detectors, advancements in beamline automation, and remote access [28, 29]. The demand for superfacilities is well supported by the creation of consortiums such as PaNdata, where 14 European research infrastructures, operating hundreds of instruments, with a yearly user base in excess of 30 000 researchers, have united to address common problems with a collaborative approach, providing researchers with a joint platform, with data interaction standards and tools for their experiments, available throughout a variety of facilities [30]. This platform is echoed by the Research Data Alliance Photon and Neutron Science Interest Group, which focuses on gathering synergies between photon and neutron research facilities, promoting a cooperative model for scientific data management and analysis [31]. These initiatives are well supported by the development of techniques, hardware, and software for data acquisition, as well as real-time and offline analysis, documentation, archiving and remote data access, as proposed by the High Data Rate Processing and Analysis Initiative [32].

In this paper, we discuss the process of establishing China's first superfacility, the Big Data Science Center (BDSC), housed within the Shanghai Synchrotron Radiation Facility (SSRF). We provide complete details of our initiatives for enabling user science at SSRF, with particular consideration given to recent developments in AI, ML, and robotic automation. The BDSC delivers a unique, fully operational, facility-wide multimodal scientific framework. It fully integrates multiple beamlines at SSRF within a facility-wide dedicated computational and storage infrastructure, encompassing data acquisition, quasi-real-time data analysis, final results, and fully automatic feedback, combined with a unified and centralized data movement, annotation, archiving, management, and curation system, complemented by a metadata ingestion system, rendering the framework AI-ready. Furthermore, modules including user proposal management, user publication management, and an open data sharing system are being developed, which will complete the centralization and unification of the data lifecycle of the whole facility. In addition, the unique and fully operational facility-wide framework of the BDSC is optimized for flexibility and scalability, thereby creating an adaptive system helpful to reduce the redundant work of beamline scientists and users at different beamlines.

The remainder of this paper is organized as follows: In section 2, we provide an overview of our community, outlining the roles of the SSRF and BDSC in relation to scientific research, its user communities, and the role of the superfacility. In section 3, we discuss the various infrastructures and initiatives underpinning the community around the SSRF and BDSC, i.e., the BDSC platform, encompassing Big Data and metadata frameworks, systems, software, data, and control infrastructures. We then discuss our efforts to facilitate high throughput data science for use in scientific research in section 4. We focus on two key aspects here: high performance computing, and AI/ML. We then conclude the paper in sections 5 and 6, outlining the vision for the BDSC in the future.

## 2. Shanghai Synchrotron Radiation Facility (SSRF) and Big Data Science Center (BDSC)

### 2.1. Shanghai Synchrotron Radiation Facility (SSRF)

The SSRF is located in Pudong New District, Shanghai, China (see figure 1); it is a third-generation synchrotron radiation facility, where the electromagnetic radiation is produced by relativistic particles moving along a curved orbit, under the effect of electromagnetic fields. It is equipped with a 150 MeV linear accelerator, a 3.5 GeV booster, a 3.5 GeV storage ring, beamlines, experimental stations, and supporting facilities [33]. Currently, there are 15 beamlines and 19 experimental end stations in operation at the SSRF [34–46]. The SSRF Phase-II Beamline Project upgrade is currently underway; by end of 2022, there will be nearly 40 beamlines in operation (figure 1).

The ongoing SSRF Phase-II Project is a key national scientific infrastructure project that includes the construction of 16 state-of-the-art beamlines, new supporting laboratories and buildings, and an upgrade to the accelerator [34, 47]. It aims to significantly improve the performance and capabilities of the SSRF, as well as the final SSRF user experience. Among other capabilities, it includes the following: 10 nm ultra-high spatial resolution; 100 ps ultra-fast time resolution; part-per-billion ultra-high sensitivity elemental analysis; *in situ* analysis for multi-elements and multi-components across multiple energy regions; multi-scale structural analysis, ranging from nanometers to centimeters; multi-level dynamic analysis, ranging from
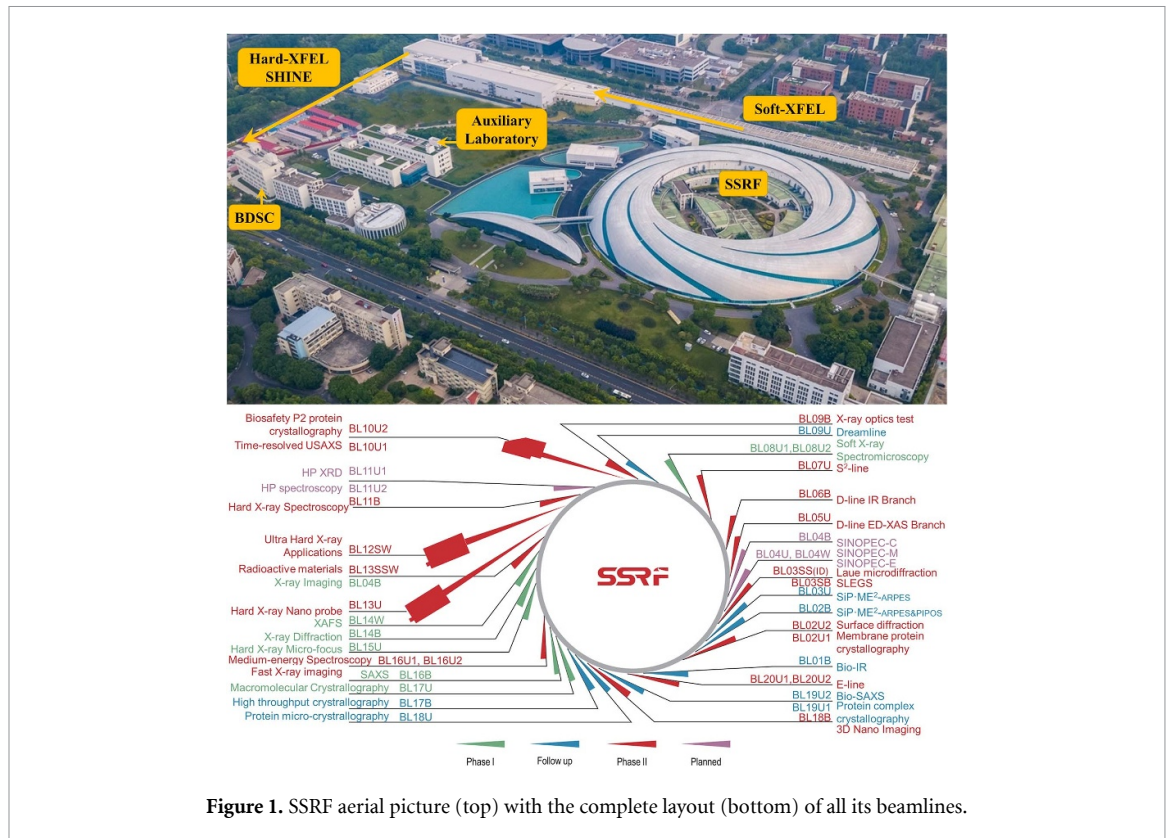
**Figure 1.** SSRF aerial picture (top) with the complete layout (bottom) of all its beamlines.

100 ps to 1000 s; experimental handling of samples requiring special conditions, such as biohazardous and radioactive materials; *in situ* experiments under extreme conditions (high pressure, strong magnetic field, etc); and comprehensive user support, including massive data storage and Big Data analysis frameworks.

In addition to the 16 state-of-the-art beamlines [34, 48–53], an x-ray test beamline [54] (see table 1), and a resonant inelastic x-ray scattering station, the SSRF Phase-II Beamline Project includes supporting laboratories, specifically materials, chemical, and biomedical laboratories, which are located inside the new auxiliary laboratory building. Moreover, as part of the SSRF Phase-II Beamline Project, a new five-floor user data center building has been constructed, which hosts the BDSC. The accelerator upgrade includes a super-bend-based lattice configuration, which aims to increase the number of straight sections, the development of various undulators, and a superconductive wiggler [55–57] to extend the photon energy range, and control of the electron bunch length to achieve a high single bunch current in the storage ring.
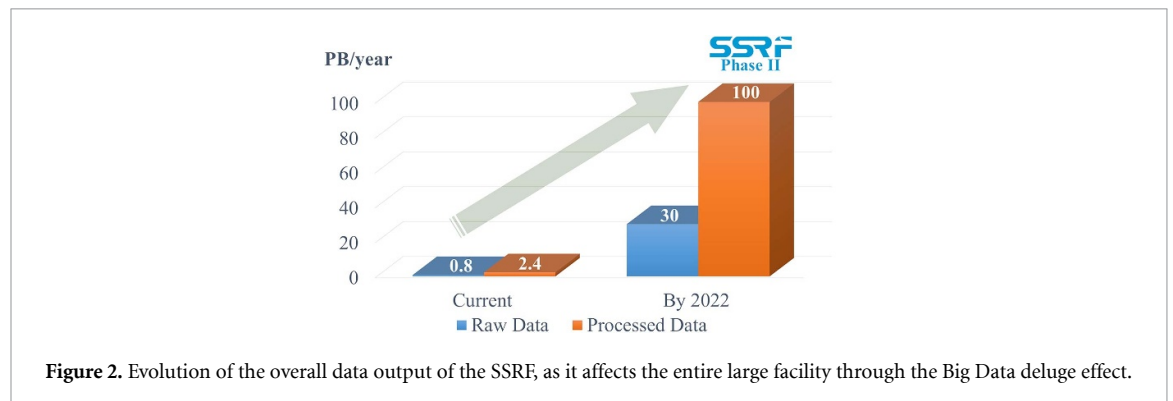
By 2022, upon completion of all the SSRF Phase-II beamlines currently under construction, the average raw data output of the SSRF will amount to approximately 160 TB d$^{-1}$, with a maximum burst peak rate of approximately 600 TB d$^{-1}$, which will produce more than 30 PB of raw data and 100 PB of processed data per year (table 1 and figure 2). Currently the BDSC requires 3 min on average to process a dataset of approximately 10 GB; therefore to process 160 TB of raw data per day, equivalent to 16 384 datasets, the BDSC would require 819 h to complete its daily task, which is equivalent to 30–50 times more computing resources than those currently available. Therefore, to achieve real-time Big Data processing capabilities, a total of 10 PFlops and a data throughput of 50 GB s$^{-1}$ are required.

## 2.2. Big Data Science Center (BDSC)

The major source of Big Data at large scientific facilities can be attributed to the beamline detectors. Beamlines generate diverse and complex data from the most diverse unrelated multidisciplinary beamlines, thereby creating the Big Data deluge [13], which must then be then centralized and standardized at Big Data centers. An increasing number of synchrotron beamlines, which employ the most diverse experimental methods, use area detectors that are characterized by large areas and high spatial resolutions, which, combined with the fast and ultra-fast experimental methods adopted, rapidly increase the data acquisition rate, resulting in the Big Data deluge phenomenon [1, 10, 13]. Furthermore, at large scientific facilities, Big Data centers must unify data in a centralized pipeline, before rendering it machine-readable, which constitutes a secondary source of data, as it is generated by the beamline control system. This is crucial for future applications that envisage the use of AI-assisted robotic automation, because the beamline control system controls all the hardware at the beamlines, including the encoders, engines, mirrors, and mechanical

**Table 1.** Detailed overview of the 16 state-of-the art beamlines and the x-ray test beamline, components of the SSRF Phase-II Beamline Project.
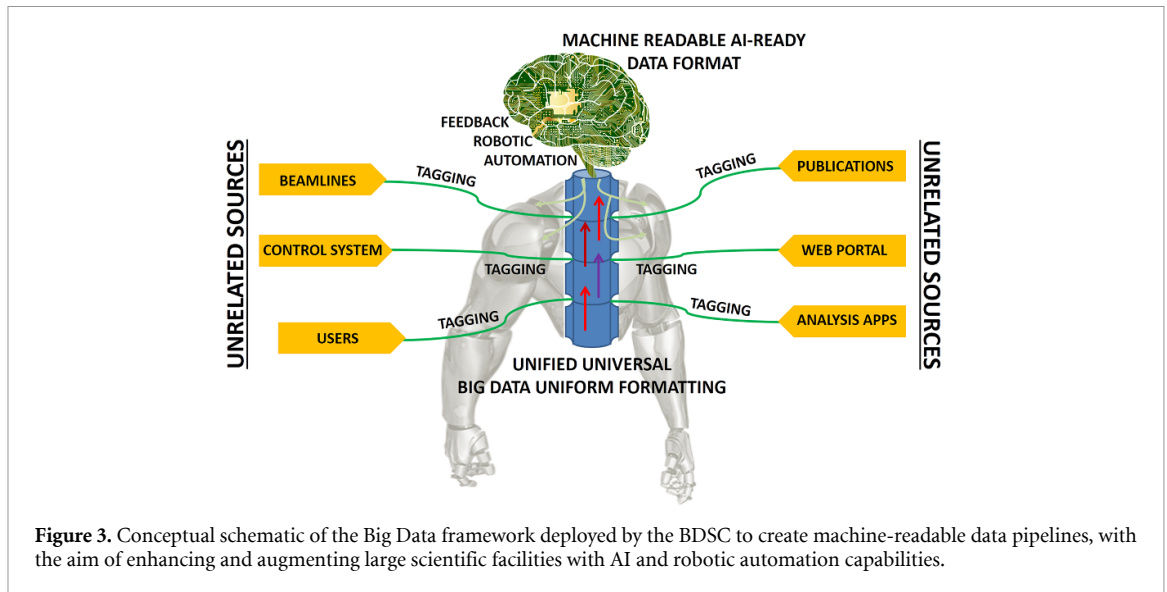
| ID | Name | Source | Energy range | Average output (byte d$^{-1}$) | Max burst (byte d$^{-1}$) |
|---|---|---|---|---|---|
| BL20U1 BL20U2 | Energy material beamline (E-line) | IVU + EPU | 130 eV–18 keV | 300 M | 500 M |
| BL11B | Hard x-ray spectroscopy beamline | BM | 5–30 keV | 300 M | 500 M |
| BL16U1 | Medium-energy spectroscopy beamline | U | 2.1–16 keV | 300 M | 500 M |
| BL07U | Spatial-resolved and spin-resolved ARPES and magnetism beamline (S$^2$-line) | Twin EPU | 50 eV–2 keV | 300 M | 500 M |
| BL02U1 | Membrane protein crystallography beamline | IVU | 7–15 keV | 15 T | 20 T |
| BL02U2 | Surface diffraction beamline | CPMU | 4.8–28 keV | 15 T | 20 T |
| BL03SB | Laue micro-diffraction beamline | SB | 7–30 keV | 10 T | 15 T |
| BL13U | Hard x-ray nanoprobe beamline | IVU | 5–25 keV | 5 T | 10 T |
| BL18B | 3D nano imaging beamline | BM | 5–14 keV | 5 T | 10 T |
| BL05U BL06B | Dynamics beamline (D-line) | IVU + BM | 5–25 keV 10–10 000 cm$^{-1}$ | 800 G | 2 T |
| BL10U1 | Time-resolved USAXS beamline | IVU | 8–15 keV | 4 T | 10 T |
| BL16U2 | Fast x-ray imaging beamline | CPMU | 8.7–30 keV | 80 T | 480 T |
| BL10U2 | Biosafety P2 protein crystallography beamline | IVU | 7–18 keV | 15 T | 20 T |
| BL13SSW | Radioactive materials beamline | W | 5–50 keV | 300 M | 500 M |
| BL12SW | Ultra-hard x-ray applications beamline | SCW | 30–150 keV | 5 T | 10 T |
| BL03SS (ID) | Laser electron gama source beamline (SLEGS) | ID | 0.4–20 MeV | — | — |
| BL09B | X-ray test beamline | BM | 4–30 keV | — | — |



**Figure 2.** Evolution of the overall data output of the SSRF, as it affects the entire large facility through the Big Data deluge effect.

and mechatronic elements [58, 59]; therefore, it is important to remotely feed back the AI cloud into Big Data centers, via robotic Internet of Things (IoT) edge technologies, locally deployed at the beamlines.

The BDSC at SSRF is therefore tasked with centralizing, unifying, managing, curating, and formatting all data. These data should be in a universal, taggable format, to facilitate easy analysis, visualization, and incorporation of data from all data sources and pipelines, encompassing the entire SSRF by means of a Big Data framework. This approach enables the creation of machine-readable pipelines that feed AI and robotics for science, training, and operations (figure 3) [13, 60, 61].

Augmenting the SSRF with the BDSC will evidently benefit the SSRF user community since, conventionally, users are overloaded with beamline operations during their experiments at SSRF, along with demanding data post-processing tasks subsequent to their experiments at SSRF (or indeed at any other large scientific facility worldwide). Currently, users can only access their final results after months of data analysis, and only then can they begin the process of interpreting the core science resulting from their experiments at SSRF. Moreover, it is unfortunate that, owing to the Big Data produced at large scientific facilities, and the limited resources available to users at their home institutions (both in terms of human expertise and scientific computational capabilities), only a small fraction of the data collected at the beamlines of these

**Figure 3.** Conceptual schematic of the Big Data framework deployed by the BDSC to create machine-readable data pipelines, with the aim of enhancing and augmenting large scientific facilities with AI and robotic automation capabilities.

large scientific facilities worldwide (approximately 40%) can then be effectively analyzed [24]. This constrains access to a significantly larger number of scientific discoveries and technological advancements, which remain hidden in the larger portion of the data collected at the beamlines, which are then archived and forgotten at the user's home institutions, and never analyzed. Thus the Big Data deluge in science accounts for considerable losses in terms of scientific discoveries and their corresponding technological implementations. As a result, only a small fraction of the Big Data produced at large scientific facilities is utilized in publications [62], and in a few years, owing to the accelerated growth of novel detector and computational technologies within the research and development domain, large scientific facilities may simply run out of capacity to produce meaningful scientific results [14]. Big Data at large scientific facilities must be appropriately analyzed, rather than being archived. Therefore, it is necessary for large scientific facilities to support users with their beamline experiments along with their scientific data interpretation. In fact, a significant increase has been observed in the number of Big Data centers for science [63]. The BDSC aims to provide the users at SSRF with an unprecedented tool, combining data science with user science, thereby effectively moving the user focus from data analysis to the interpenetration of scientific results, because the latter, ultimately, impacts on users' publications . Therefore, the objective of the BDSC is to considerably increase the number of scientific discoveries and technological advancements based on the researchtaking place at large scientific facilities, and to increase the number of scientific publications, making scientific results widely available to the international community that can then go on to be utilized by society as a whole.

Furthermore, the BDSC intends to provide real-time analysis capabilities to the users, so that they can make judgments directly during beamline experiments in terms of data quality and the corresponding scientific results; if these do not satisfy their expectations, the experimental procedure can immediately be reconfigured in real-time at the beamline. This real-time capability provided by the BDSC will address a critical issue currently affecting measurement procedures at large scientific facilities worldwide, where users can access their scientific results only after the conclusion of the experiment at the beamline, not to mention a subsequent lengthy data post-processing procedure, which is often conducted only when the users return to their home institutions [64]. Therefore, before the BDSC was in operation at SSRF, the users were often unable to verify whether or not their experiments at the beamlines were progressing well. Considering the high cost of beamtime preparation and operation in terms of human resources, preparation time (often amounting to months prior to the actual beamtime), funding, and the extremely limited amount of beamtime available to each user (only few days per year), combined with the recent Big Data deluge at large scientific facilities [13], any approach not including the type of Big Data infrastructure capable of properly supporting user science at large facilities is no longer sustainable.

The role of the BDSC is therefore central to enabling users to successfully accomplish their overall scientific objectives. Moreover, the dramatic simplification introduced by the BDSC with respect to the data and operation pipelines for users at large scientific facilities has now broadened access to these facilities to a plethora of researchers from the most diverse disciplines, who were previously intimidated by the cumbersome experimental framework and subsequent data interpretation process typical of large scientific facilities. In fact, many users who might benefit most from these large scientific facilities, while producing

**Figure 4.** The superfacility framework.

scientifically and socially impactful results, bridging science with technology and industrial applications, were left out, because they required experts to prepare their experiments and evaluate their data; these are not always available in every field of study. The BDSC aims to reduce the need for large facility experts within user research groups, thereby dramatically expanding the user base of these large scientific facilities, and sharing the most advanced large facility science to benefit all scientific research fields.

Another objective of the BDSC includes contributing to the international scientific community in addressing the aforementioned scientific challenges, which are not limited to the SSRF; these are common issues associated with any large scientific facility worldwide. Furthermore, these issues are not limited to synchrotrons; they extend to other large research infrastructures, including neutron reactors and XFELs. Therefore, the BDSC is working on a horizontal solution to address Big Data issues in science via a centralized infrastructure, extendable to international large scientific facilities, creating a widespread impact across the heterogeneous international scientific community.

Therefore, by combining Big Data science, AI, and robotic automation, the BDSC aims to create an augmented user experience at the SSRF, to increase the number of scientific discoveries, publications and industrial technologies resulting from experiments at the SSRF. To achieve this, the BDSC is engaged in supporting users at the SSRF to return to their home institutions with publication-ready scientific results, rather than simply raw data, and to provide them with local expertise to accelerate the scientific interpretation of their data. Accordingly, real-time data analysis and final result visualization are at the core of the BDSC infrastructure; this requires fast modeling, reconstructions and simulations, while visualizing the final results, which should be publication-ready, and can be directly comprehended by non-experts.

The BDSC is therefore augmenting the SSRF with a superfacility [13], which is the first of its kind in China.

### 2.3. Superfacility

To implement a Big Data framework capable of supporting AI and robotic automation at SSRF, the BDSC is deploying a superfacility infrastructure (see figure 4). The Superfacility is a new concept for large scientific facilities, which aims to create a virtual facility [13, 14, 24, 63]. A superfacility essentially comprises an interconnected network, utilizing the most relevant components necessary to accelerate users' scientific productivity, facilitating a connection between experiments at the beamlines and eventual scientific publications. It should be noted that user experiments at the beamlines constitute the starting point of a significantly longer scientific journey. A superfacility is intended to accompany users through this journey (figure 4). User experiments at the beamlines result in Big Complex Data, which, within the Superfacility framework, are analyzed by the most advanced algorithms, supported by the latest theoretical frameworks, encompassing several scientific domains, to produce simulations, modeling and reconstructions based on these Big Data; these algorithms must be implemented locally, using the most recent generation of HPC clusters, which are then elastically scaled-out to national supercomputers [64, 65] when the large facility workload exceeds the Big Data center's cluster capabilities. Next, the processed Big Data must travel thorough Big Data networks [66–68], supporting the real-time requirements of the demanding experiments taking

place at large facilities. This must be presented seamlessly to the users via a well-packaged, middleware, user-friendly interface, capable of performing real-time data interpretation directly during the users' beamtimes. The final results are expected to leverage Big Data visualization on thin client terminals (including laptops, tablets, mobiles, etc), both locally at the beamlines, and remotely, to produce a final outcome, which closely resembles the final results, which can then be directly included by the users in their scientific publications. This whole process must be transparent to the users. A superfacility, with all due and evident differences, can thus, conceptually, transform a large scientific facility into a super-microscope: easy to use, easy to understand, easy to share among non-experts, and easy to publish. Accordingly, the superfacility aims to shift the focus of Big Data centers at large scientific facilities from being data-centric to scientific-publication-centric.
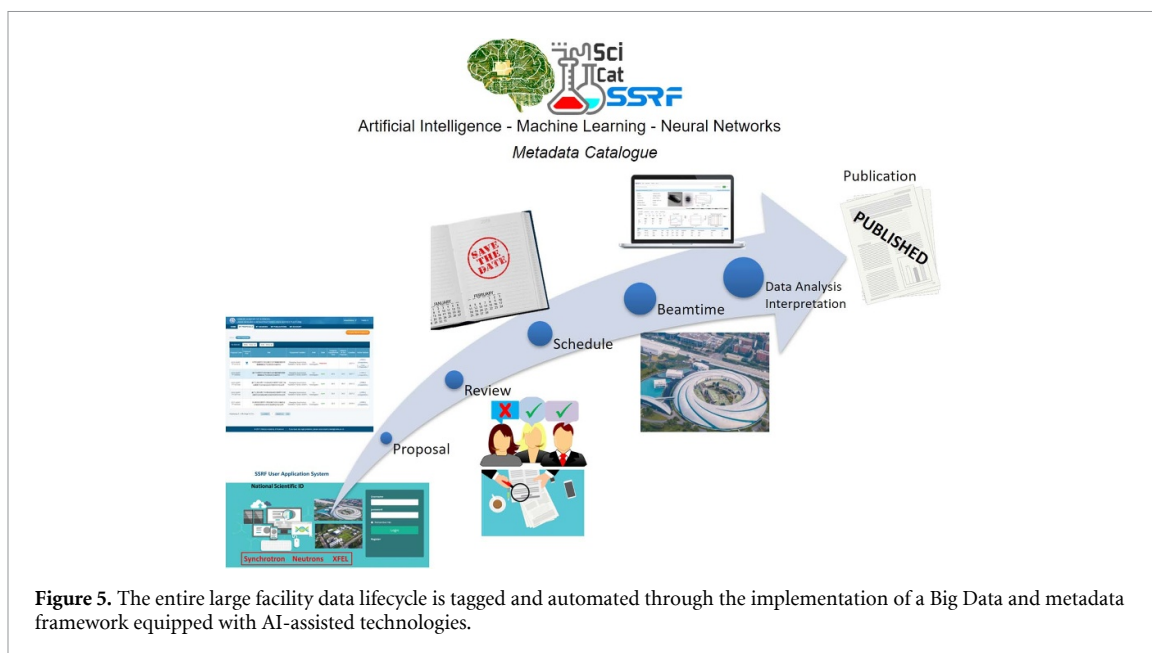
### 2.4. User base

Since its inception, the SSRF has supported more than 23 000 scientists per year, from more than 2500 research groups, from more than 500 institutions in China [34, 69]. Owing to this support, a series of important scientific discoveries, involving frontier and applied research in structural biology, condensed matter physics, chemistry, material science, the environment, and medicine have been made [70, 71]. More than 5000 peer-reviewed publications have been published by SSRF users, with an annual average beamtime availability of approximately 5500 h [34]. By the end of 2022, when all SSRF Phase-II beamlines are expected to be operational, nearly 10 000 additional researchers will be supported per year.

Before the BDSC was in operation at SSRF, the beamlines were uncoordinated, and each had their own processing pipelines and operational aspects. Most of the beamlines did not even have the notion of metadata, and largely relied on manual data acquisition and analysis, which lacked any data security, resulting in missing authentication, permission, and access control management capabilities.

However, this has changed recently, owing to the primary focus of the BDSC, i.e., the user science focus. The users of the BDSC, and thus of the SSRF, now rely on real-time data analysis, remote data collection, visualization and interpretation, enabling the rapid evaluation of their data onsite during their beamtimes, as well as remotely from their home institutions, thereby increasing the number of scientifically valuable datasets produced, and enabling more flexible interactions with the users in terms of data transfer and remote analysis.

More specifically, the user community relies on the Big Data framework, metadata and data tagging capabilities of the BDSC, which are crucial for training AI models [72, 73]. The overall user experience within SSRF is fully followed through. In addition to user experience, the overall information concerning the experiments and their resulting datasets is maintained with quality control in place. This includes, but is not limited to, sample specifications, beamline techniques, experimental setups, beamline configurations, parameters for user data analysis, simulation, modeling, reconstruction, and visualization of publication-ready results. Thisinformation is routinely utilized for training AI and ML systems. In fact, the provision of real-time analysis capabilities to the users can further facilitate AI training for data analysis and experimental configurations. An AI can then be trained by the most experienced users; subsequently, this knowledge can be transferred and shared between other SSRF users. This will be particularly advantageous for users with less experience. Furthermore, the BDSC, once fully interfaced with the SSRF beamline control system, can collect all relevant information from the beamlines concerning configurations, alignment, and setups [74–76]. The BDSC can therefore monitor all the configuration parameters used by a specific beamline to perfect experiments based on user requests, then feed them to the AI via ML. As such, AIs can be trained to achieve the optimal beamline configuration, based on user requirements, and to provide feedback to the robotics and IoT at the beamlines to implement a fully automated and unmanned beamline station. As the work of specific beamline teams is, at the moment, highly involved, and does not allow for any focus on beamline improvements, given that beamline staff tend to be fully engaged in configuring their beamlines for the users, our Big Data framework can also alleviate beamline responsibilities, shifting their focus to fully support the user science. In addition, the BDSC Big Data framework intends to implement real-time data pipelines at the SSRF, where the users do not require assistance from beamline scientists once the AI-assisted technologies are fully connected to the robotics, which, in many past cases, had forced users to suspend their experiments for long periods. This capability should also benefit users during night-shifts, when limited beamline assistance is available. The use of AI-assisted technology, once fully implemented at the SSRF beamlines, will therefore assist users to focus on their experiment's optimization, based on scientific results that they are able to visualize in real-time at the beamline. It should be noted that each beamline has its own interface, even within the same facility; as a result, a significant amount of time is lost in learning different configurations not strictly related to the experiment. Once the AI deployed by the BDSC at the SSRF is fully operational, and has been trained using the experience of expert beamline operators, a user-friendly universal interface can then be presented to the users. The AI will automatically configure unmanned

**Figure 5.** The entire large facility data lifecycle is tagged and automated through the implementation of a Big Data and metadata framework equipped with AI-assisted technologies.

beamlines for users' experiments, based on their experimental requirements. Furthermore, users cannot currently change their samples once a beamline is in operation, and is irradiated with synchrotron radiation, owing to the intense x-ray exposure within the experimental hatch, or when using experimental setups intended to operate under extreme conditions (in terms of pressure, temperature etc); therefore, when a sample needs to be changed, the experiment must be suspended. Consequently, real-time experiments are not achievable under these conditions. However, the use of AI-assisted technologies enables this, by implementing ML to automatically instruct unmanned robotics and IoT inside the beamline hatch to change the samples on behalf of the users [77].

## 3. Infrastructure

### 3.1. Overview

The BDSC is deploying a superfacility with a particular focus on Big Data and metadata frameworks, which constitute the cornerstones of the superfacility.

The infrastructure of the BDSC covers:

- An automated Facility Data Management System.
- A metadata system.
- An on-premise computing platform.
- Robotic systems for remote experiments.
- AI/ML data analysis systems.

We discuss these in the following sub-sections.

### 3.2. Automated facility data management system

To build a superfacility, it is necessary to maximally automate the entire data lifecycle of a large facility, to tag all the Big Data produced during the experiments with the metadata produced in each step of the lifecycle itself (figure 5) [78, 79]. This will guarantee the creation of a robust basis for AI training through ML, and robotic automation through the IoT.

Therefore, the metadata are collected using a metadata catalogue, to retain all the necessary information linked to any user's interaction with the large facility [13, 80]. To devise a comprehensive tagging system for all the Big Data produced at large scientific facilities, it is necessary to collect and link all the user information, from user ID to final publications (figure 5), rather than only that strictly connected to the measurements at the beamlines. In particular, it is necessary to collect information regarding user credentials, so as to unequivocally identify all the data produced by the users. This must be then associated to the proposals submitted by the users, in order for their experiments to be considered for beamtime. Furthermore, all text entered by the users into their proposals must be machine-readable and recognizable, to facilitate text-based AI training. This will be crucial in preparing AI-assisted technologies capable of

supporting less experienced users in preparing a solid proposal; the AI will suggest the best experiment and beamline configuration to the users, along with the best beamline and technique for their experiments, based on its knowledge of all the previous proposals written by experienced users. Specifically, besides the machine-readable text, the metadata catalogue will store information concerning the proposal ID, sample ID, experimental topic, and beamline requirements. After the proposal has been assessed by the peer-review panel at the large scientific facility, the results, comments, and final evaluation are presented to the users and stored within the metadata catalogue. This automating facility data lifecycle information can be used for AI training, to classify successful and rejected proposals. In this way, AI-assisted technologies can be trained to suggest the optimal proposal framework to users for their experiments. As such, from the earliest stage of proposal submission, the users can focus on their science, rather than having to concentrate on beamline configurations and experimental setups details. This should substantially assist numerous inexperienced users in approaching a large scientific facility. Moreover, this is a necessary feature if a large scientific facility intends to include a wider, higher impacting scientific community, rather than limiting itself to experts with in-depth knowledge of the science, technology, and operations of large scientific facilities. Following a successful beamtime application (figure 5), the beamtime is then scheduled, and the corresponding available dates are presented to the users; all information concerning the beamtime calendar for each user, including chronological information relating to their previous beamtimes, both accepted and rejected, will be stored within the metadata catalogue. At the time of the actual experiment at the beamline, the metadata catalogue will be further enriched with all the parameters used by the beamline operators to configure and set up the beamlines, as well as all the parameters required by the users to analyze and interpret their data and final results. These metadata can then be used to train the AI-assisted technologies at large scientific facilities. Once fully trained, they will provide a fully automatic robotic beamline, capable of self-configuration, based on knowledge transfer from experienced beamline operators to the AI, while implementing real-time unmanned, or quasi-unmanned, data analysis. This will make the beamlines self-configurable, based on the real-time results from actual measurements. This is achievable only if the AI is trained using metadata produced by the data analysis and interpretation of the most experienced users at the beamline. To facilitate the process of knowledge transfer from experienced users to AI, all the software used for data reduction, evaluation, analysis, and interpretation by the users will be uploaded to the Big Data center cluster, fully integrated and optimized for the HPC architecture, and made available to the users for their real-time analysis at the beamline. This will simplify metadata collection on the users data analysis pipeline, because the whole analysis will be performed within the Superfacility infrastructure, which will then facilitate metadata ingestion for the purpose of AI-training (see figure 5). The final aim is to allow both users and AI to quickly evaluate the quality of data in real-time at the beamline, and, if necessary, adjust the experimental configuration at the beamline directly during the beamtime. This will also allow users to return to their home institutions with a dramatically increased quantity of higher-quality data and publication-ready results. Finally, once the users publish their results, all information connected to their publications, including the corresponding digital object identifier (DOI), will be stored within the large scientific facility metadata catalogue, linked with all other metadata produced by users in the subject area on which the specific publication is based. This will enable the comprehensive tracking of all user operations, interactions, and productivity throughout the whole data lifecycle of a large scientific facility. All the metadata collected will be linked to the Big Data produced by the user experiments, and employed as a tagging system for AI training through several ML approaches, thereby dramatically increasing users' scientific productivity by virtue of a fully automated AI-assisted robotic facility.

### 3.3. Metadata system

Metadata constitutes the core of the entire Superfacility framework (figure 5) [81]; therefore, it is important to select the most appropriate metadata catalogue to support the Big Data tagging system for AI training, which can then provide feedback to the IoT robotics at the beamlines.

Therefore, the selection of the most appropriate metadata cataloguing system is fundamental to allowing the metadata to be fully centralized, accessible, linked, uniformly formatted, organized, and machine-readable, based on all unrelated data sources present at the large scientific facility. The metadata catalogue, in turn, will support AI training, which uses the most updated ML approaches and scientific multidisciplinary theoretical frameworks. The process of centralization is relevant to a large scientific facility, because it will enable the seamless update of AI with the most recent algorithms and theories. This is particularly important within those scientific fields that use large scientific facilities where both the scientific theories and AI frameworks are still incomplete, and are constantly updated. A flexible centralized upgrading framework is crucial to avoid obsolescence, while its proper integration within the centralized infrastructure ensures the absence of any disruption within the operations and data pipelines of the large facility as a whole. A centralized and uniformly formatted machine-readable metadata structure will therfore permit rapid

large-scale reconfigurations of the whole superfacility framework, while remaining transparent to both users and the operations of the large facility.

ICAT [82] is a robust metadata cataloguing solution [13, 83–86] which has emerged in recent decades; however, a novel metadata cataloguing system, SciCat [87], has recently been developed to include all the most recent advancements within the field of scientific computation, leveraging the most advanced technologies necessary to build a modern metadata framework for Big Data, ready to interface with the AI and IoT applications required for scientific research.
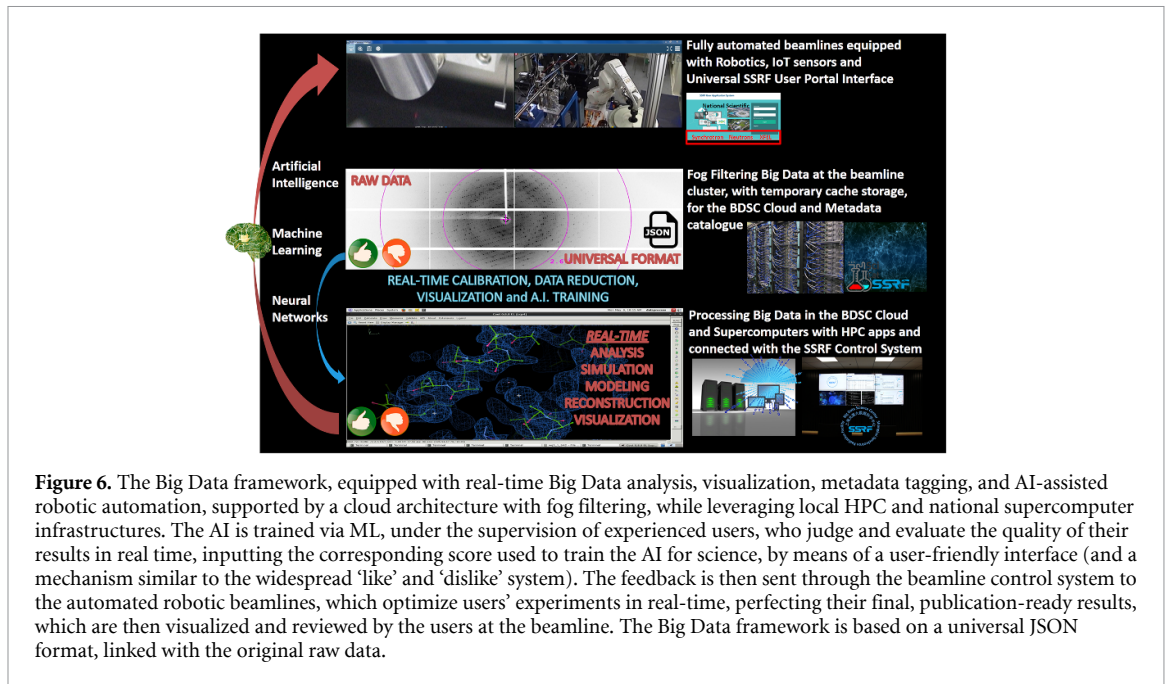
The ICAT framework, developed by the Science and Technology Facilities Council, UK Research and Innovation in U.K., has been deployed at the Diamond Light Source, the ISIS neutron facility and Central Laser Facility at the Harwell Campus in U.K., along with the Helmholtz-Zentrum Berlin für Materialien und Energie in Germany, the Optics Clustered to Output Unique Solutions facility in U.K., the Institut Laue-Langevin in France, and the China Spallation Neutron Source in China [83–86, 88–91], among others. The ICAT system aims to catalogue, archive, and manage data via the development of a metadata database, which monitors, indexes, and labels all data connected to the users, from the initial proposal, through data collection, to the final publications, allowing them to be tracked and organized, as well as conveniently accessed by the users for their data evaluation. The ICAT system is based on a study-dataset-oriented Core Scientific Metadata model (CSMD) [82, 92–94], which captures high-level information concerning scientific studies and the corresponding data that they produce; this includes core entities such as investigation, investigator, topic, keyword, publication, sample, dataset, datafiles, parameter, authorizations, etc. The ICAT schema is thus strongly based on the CSMD. ICAT adopted MySQL as a core database structure, where the data table must be customized and designed according to the unique procedure for each experiment at each beamline, commonly using different technologies and methods, with a long development cycle, poor scalability, and great difficulty in handling subsequent changes, rendering the implementation of Big Data frameworks supporting AI and IoT for science impractical.

SciCat is an open-source data cataloguing management system, jointly developed by the European Spallation Source in Sweden and Denmark, the Max IV synchrotron in Sweden, and the Paul Scherrer Institute in Switzerland. SciCat aims to manage the entire metadata lifecycle at large scientific facilities using, among other modern data technologies, MongoDB, a NoSQL database, which facilitates the deployment of Big Data architectures for AI and IoT for science. At the BDSC, SciCat has been adopted, adapted, enhanced, deployed, and integrated within a wider architecture covering the entire data lifecycle at the SSRF.

### 3.4. On-premise computing platform

This is provided via a Big Data framework, capable of combining the experiments at the beamlines with the metadata database, augmented with AI, and supported by HPC clusters, while equipped with remote and local analysis services for academic and industrial users. The BDSC provides access to its Big Data framework via a web portal and software applications, which run on thin clients, e.g., local computers at the beamlines, and user devices such as laptops, tablets, or mobiles.

The BDSC has therefore implemented a data pipeline that runs from the beamline detectors to the visualization of the final results on thin user terminals, via a user-friendly interface. The detectors generate Big Data that are temporarily cached within the clusters of beamlines, effectively acting as a fog computing layer. The Big Data framework at the BDSC has been built without any restriction on the data formats generated at the beamlines; as such, it can effectively handle any data format produced by any beamline during the experiments, and fully and seamlessly integrate them within the BDSC Big Data and metadata architecture. All the Big Data are methodically tagged with the user ID, which matches with the credentials stored within the SSRF user authentication system. The BDSC Big Data framework restricts users and SSRF operations only to those with a valid user ID, thereby allowing it to effectively tag all data produced from any source at SSRF with the corresponding user ID. This enables the reliable and consistent tracing of user operations and production across the entire superfacility lifecycle. The Big Data are then remotely and automatically transferred from the beamline clusters to the BDSC cluster via a high-speed network. Furthermore, the beamline clusters are also in charge of the basic beamline alignments and data calibrations. As these are interfaced with our metadata catalogue, they also serve as a source for Big Data tagging and subsequent AI training. The users are then offered the option to visualize their reduced data in a user-friendly manner, and select the data they wish to submit to the BDSC for further processing (figure 6). As this step is fully tracked by the BDSC metadata system, all the Big Data processed at the BDSC is fully tagged with all the relevant user information. In this way, the BDSC effectively uses a fog filter, supervised by the users. After the BDSC Big Data framework has visualized the reduced data for the users, and prior to submission to the BDSC cluster for data processing, the Big Data is tagged with user decisions made while selecting only those datasets judged valuable for their scientific research. This can effectively provide a tagging system, supervised by the users, which feeds into the AI training. Only the data selected and judged valuable by the users are then

**Figure 6.** The Big Data framework, equipped with real-time Big Data analysis, visualization, metadata tagging, and AI-assisted robotic automation, supported by a cloud architecture with fog filtering, while leveraging local HPC and national supercomputer infrastructures. The AI is trained via ML, under the supervision of experienced users, who judge and evaluate the quality of their results in real time, inputting the corresponding score used to train the AI for science, by means of a user-friendly interface (and a mechanism similar to the widespread 'like' and 'dislike' system). The feedback is then sent through the beamline control system to the automated robotic beamlines, which optimize users' experiments in real-time, perfecting their final, publication-ready results, which are then visualized and reviewed by the users at the beamline. The Big Data framework is based on a universal JSON format, linked with the original raw data.

sent to the BDSC, along with all the user input parameters, which are then stored within the BDSC metadata system. These input parameters are necessary for running the users' simulations, modeling, analysis, reconstructions, and data interpretation in real time, using the software uploaded by users directly into the BDSC cluster. By virtue of the BDSC Big Data framework, the AI can effectively learn any reduced data at the beamline that achieves the required quality levels to merit further analysis at the BDSC cluster; it then effectively acts as a cloud system, equipped with AI-assisted fog filtering. It also further learns the parameters selected by the users to analyze the datasets assigned by them to the BDSC for further analysis, providing less experienced users with an automatic unmanned selection system for their datasets, subsequent to AI training, without requiring expert assistance at the beamlines. A fully trained AI at the BDSC can assist users in selecting, filtering, presenting, and submitting to the BDSC analysis cluster only those datasets containing scientifically valuable information, obviating the need for a manual search, which would be impractical in the case of Big Data. To provide a smoother and more reliable pipeline, all the raw data from the beamlines are always automatically stored within the BDSC storage system, where the users assign and flag any data requiring further processing by the BDSC Big Data framework, which is tagged by the BDSC metadata system.

The users can then access all their data, both locally at SSRF, or through an internet-based virtual private network (VPN) with a cluster account, by remotely accessing the BDSC system. Utilising a combination of VPN and active directory (AD) accounts, users can access the BDSC Big Data framework, including SciCat, HPC, and storage systems. Currently, the BDSC uses the SSRF Internet Protocol Security VPN to grant users remote access to the hosts; however, first-time users are required to download, install and configure it. It also needs to be kept open all the time, and reopened for each new access; it only supports a limited number of clients, and it sometimes causes conflicts with other websites using the same VPN. Moreover, it is resource-demanding and may generate instabilities during Big Data transfer. Finally, it complicates the process of remotely invoking the public application programming interface (API). Therefore, to improve the user experience, the BDSC will develop a new platform, which uses a secure sockets layer (SSL) VPN and net balance gates, thereby implementing a new web portal at the SSRF, where users will no longer be requested to install a client, while benefitting from services that provide enhanced performance and compatibility, and further addressing any VPN conflicts. Moreover, taking into account future robotic automation at the beamlines, the BDSC will collect information from the beamline control system at the SSRF, relating to the entire configuration of beamlines for each experiment, and input them within the BDSC metadata catalogue (figure 6). This information will further enrich the BDSC tagging system used to tag all the Big Data produced at the SSRF, which can be used for AI training, enabling it to learn the right beamline configurations, alignments and setup parameters for each experiment from experienced beamline operators, together with robotic automation. In fact, the plan is to interface the BDSC directly with the SSRF beamline

control system, so that the BDSC is not limited to collecting information, but can extendits capabilities to sending input back to the IoT robotics at the beamlines through the beamline control system; this will help to facilitate IoT robotic automation, assisted by AI, at the beamlines. The Big Data processed at the BDSC are then visualized in real-time on the thin client terminals of users, via a user-friendly interface, where users are able to judge the quality of their final results against the scientific publication criteria (figure 6). If the users are not satisfied with the final processed results, they can engage in a new cycle of beamline reconfiguration, changing samples or data reduction, inputting new analysis parameters, or requesting a new BDSC analysis of the same dataset. Moreover, in this case, user decisions regarding the quality of their final results will be captured by our Big Data framework tagging system, and stored within the metadata catalogue, which is then fed to the AI for training by ML. Once fully trained, the AI at the BDSC will be able to fully evaluate, unmanned, the scientific quality of the users' final results, feedback to the beamline control system to initiate a robotic automatic readjustment of the whole experiment at the beamline, and re-run the whole data evaluation cycle, effectively providing the entire SSRF with a full real-time data pipeline. The BDSC Big Data framework can also link user operations and productivity lifecycle at the SSRF, together with their corresponding publications, using DOIs. To offer the most immersive, user-friendly, smooth experience, the BDSC is further equipped with a comprehensive search engine, which allows users to search through all the data collected by the SSRF. This uses an NoSQL cluster engine, which indexes all metadata to accelerate parallel searches. Moreover, the raw data are translated into key-values and linked to designated JavaScript Object Notation (JSON) metadata, which supports customized JSON queries to navigate any key or value. Following its commissioning, the BDSC is now operational, and has already collected raw, reduced, and fully post-processed data, along with their corresponding metadata. However, SciCat only provides a basic and standard title search. Therefore, the BDSC provides system-level hierarchy search capability, covering the whole SSRF; this will be embedded in a future platform, to provide a user-friendly interface, which will be connected to our metadata system, effectively providing users with the ability to virtually search through their logbooks. This provides for a solid basis to build an AI-assisted search engine at the BDSC which encompasses the SSRF, where users can type human-readable sentences, similarly to conventional internet search engines, resulting in AI-assisted suggestions, rather than searching for sample codes and dates. This feature is invaluable to further accelerating the scientific productivity of users, allowing them to quickly navigate to the most relevant data for their final publications, rather than engaging in a tedious manual search through their handwritten logbooks. Based on users' queries, the BDSC search engine can directly return the most relevant final results, which can be visualized, allowing them to immediately and more naturally interpret and adapt them for publication. By virtue of the BDSC Big Data and metadata framework, users can also trace back the original raw data that produced the final results, as well as the original sample codes associated with the labels recorded within their logbooks.

The BDSC infrastructure is developed for large facilities in general, to support the international scientific community; it can therefore be adapted to any synchrotron, and extended to other facilities worldwide, including neutron reactors and XFELs, by means of adaptable user and beamline application layers and plug-in modules [95], thereby providing a robust basis for a more uniform development of superfacility platforms. The BDSC further advocates the adoption of a national unified scientific ID, which can unequivocally identify users nationwide, irrelevant of the facility they are accessing and using. This will promote horizontal data sharing, linking and tagging between different facilities [30, 96] within a single and uniform platform, readily accessible to users. If achieved, this will dramatically improve the user experience at large facilities and AI training for scientific purposes, providing a larger and more diverse dataset basis for training. It will further promote the more widespread adoption of the powerful combined analysis approach [97–101], which allows for the combination of results from different beamlines and facilities, with a view to an enhanced interpretation of scientific results, and a deeper comprehension of natural phenomena, with the added benefit of AI-assisted technologies.

The specifications of the BDSC hybrid HPC cluster are detailed in table 2; it is further equipped with a storage system (7.7 PB hard disk drives, HDDs, and 108 TB solid state drives) designed to store metadata, raw, and processed data from all beamlines at the SSRF for a duration of 6 months, after which the data are transferred to the BDSC tape storage system for long-term archival.

The BDSC utilizes the SSRF network infrastructure, as shown in figure 7. The HPC and network-attached storage clusters are connected using a 100 Gb s$^{-1}$ InfiniBand backbone, which is further connected to a core switch at 100 Gb s$^{-1}$. The beamlines and laboratory switches are connected to the core switch using a 40 Gb s$^{-1}$ network. The beamline network is connected to a 10 Gb s$^{-1}$ staff network, which is secured by a firewall. Finally, the entire SSRF network can be remotely accessed via the China Science and Technology Network, and the China Telecom network.

**Table 2.** Specifications of the BDSC hybrid HPC cluster and storage system at SSRF.

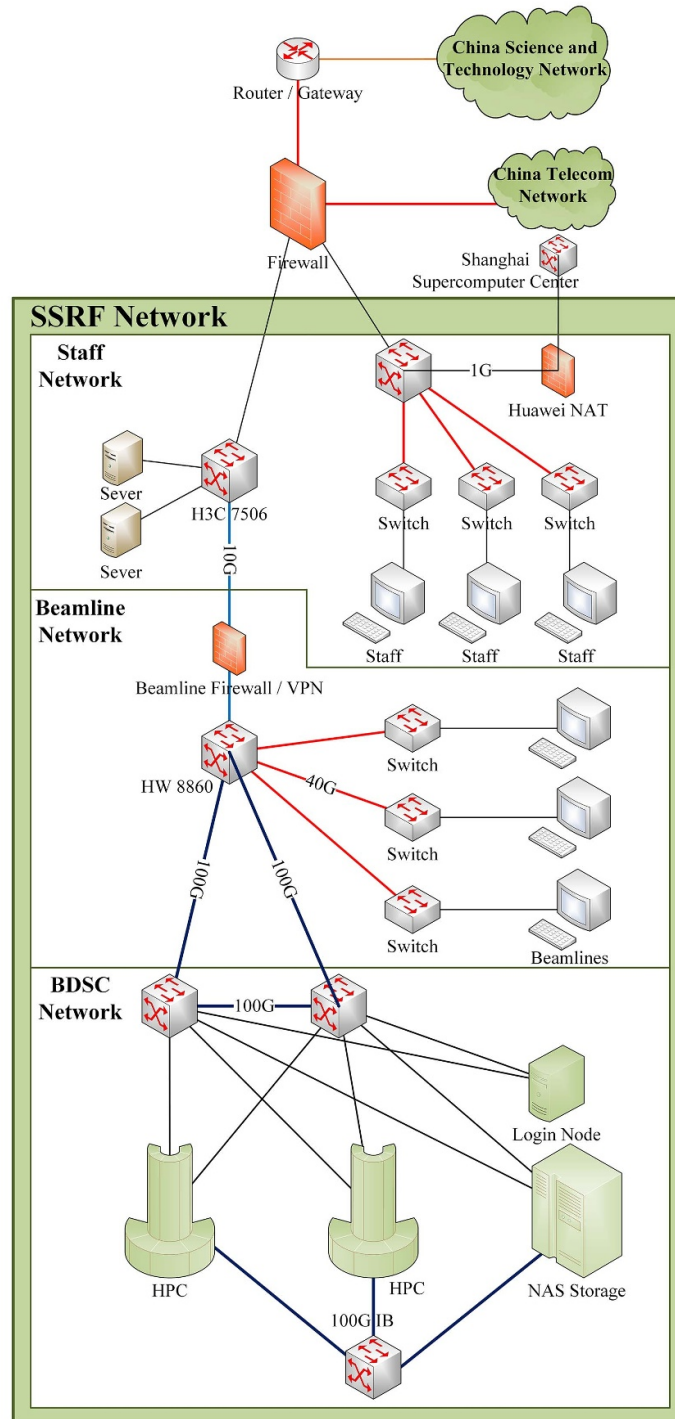| Component | No. | Configuration |
| --- | --- | --- |
| CPU node | 48 | Intel Xeon Gold 6104 (2.3 GHz, 18-Core) * 2 16G DDR4 Memory * 8 |
| GPU node (PCIe) | 4 | Intel Xeon 5118 (2.3 GHz, 12-Core) * 2 NVIDIA Tesla P100 GPU Card *2 16G DDR4 Memory * 8 |
| GPU node (NVLINK) | 1 | Intel Xeon Gold 6132 (2.6 GHz/14-Core) *4 NVIDIA Tesla P100 GPU Card *4 32G DDR4 Memory * 32 |
| Fat node | 1 | Intel Xeon E7-8860v4 (2.2 GHz/18-Core) *16 16G DDR4 Memory * 128 |
| Storage node | 30 | 8 TB SATA Hard Disk * 32 900 GB Solid State Disk *4 |

### 3.5. Beamline integration and controls

To demonstrate the effectiveness of the BDSC framework deployed at SSRF, a real case study has been selected to showcase its present and future potential in terms of increasing users' scientific productivity. The BL17U1 Biological Macromolecular Crystallography Beamline (MX Beamline) at the SSRF has been fully integrated within the BDSC framework, where a real pilot experiment was performed, including users. In fact, the users are the real target of all the BDSC endeavors, which are always focused on user science.

High energy electrons run in ultrahigh vacuum within the electron storage ring. Each beamline is designed to use the synchrotron radiation from the electronic storage ring, based on different and specific applications and research fields [102]. The main function of a beamline is to select the energy range, and to obtain a high-flux focused beam with an optimal size and small divergence [103]. Next, monochromatic light, white light, or polarized light may be obtained at different beamlines. The user experiments are then performed at the end of the beamlines, where the endstations are located, equipped with detectors, data acquisition systems, and control systems.

At the SSRF MX Beamline, a home-made in-vacuum undulator with 80 periods, and a period length of 2.5 cm, serves as the source. The beamline optical components include a double plane crystal monochromator, and a toroidal mirror that focuses the beam both vertically and horizontally. The x-ray energy range is 5–18 keV, the flux at the sample is $3.8 \times 10^{12}$ phs s$^{-1}$ (at 12.4 keV 240 mA), and the focused beam size is $67 \times 23$ μm$^2$ [36]. A self-integrated diffractometer reduces the sphere of confusion of the rotatory axis to 1 μm. The diffractometer is equipped with an on-axis viewing system capable of providing the ideal resolving power. The area detector has also been upgraded to the newest generation of detectors, EIGER X 16 M, which can collect data at 133 Hz [103, 104]. A Rigaku ACTOR robotic system is also provided, to improve experimental efficiency. The MX Beamline has provided users with a robotic sample mounting system since 2011, to help improve the efficiency of the beamline. This system now allows users to screen approximately 20 crystals per hour. In addition, this sample mounting system is a requirement for remote experiments. In fact, users can utilize all the beamline functions to collect biological macromolecular crystal diffraction data remotely via the internet. In October 2017, the MX Beamline was upgraded with an Eiger X 16 M detector, which can provide a maximum frame rate of 133 Hz. The detector can output files directly in HDF5/Nexus format. This new detector facilitated the implementation of the 'shutterless' data collection method; with this device, hundreds of diffraction datasets can be collected per day. This effectively enhanced the Big Data deluge effect at the beamline. The corresponding manual processing and recording procedures at the beamline are time-consuming and error-prone tasks.

As the most advanced beamline for data management and data processing at the SSRF, the MX Beamline designed an automatic data-processing and experiment information management system, designated Aquarium [104]. Once the dataset collection is completed, Aquarium can submit data processing jobs to a high-performance beamline cluster; it can then automatically process the dataset, from data reduction to model building, if anomalous scattering signals occur, whereby the data processing results can be monitored and investigated through a website module. All experimental information, including sample information, parameters, and data processing results are deposited into a PostgreSQL database, and accessed via the Aquarium website. Since the Lightweight Directory Access Protocol is integrated into Aquarium, all users manage their data in their own directories, and cannot access the data of other users, which ensures data privacy and security. The metadata and the data processing results requiring small storage could be stored within the beamline database and local storage system for a long time; however, the raw data was stored at the MX Beamline for 1 month, owing to the storage limitations of the cluster system locally deployed at the beamline (27 TB flash memory and 48 TB HDDs) and the high-efficiency data collection pipeline (several TB per day).

**Figure 7.** BDSC network infrastructure, interfaced with the Shanghai Supercomputer Center, accessible to external users, along with local users and staff at SSRF, through the SSRF network, the China Science and Technology Network, and the China Telecom network.

The MX Beamline runs on an extant local data pipeline infrastructure, which complicates the deployment of a novel platform, because it cannot be easily modified without disrupting the beamline's intensive operations. Therefore, it is a suitable candidate to demonstrate the flexibility of the BDSC Big Data framework when deployed at a large scientific facility. The BDSC Big Data framework can easily retroactively adapt to existing systems running at the beamlines without disrupting their demanding operations or compromising the centralization scope of the unified BDSC Big Data and metadata framework. Since the SSRF is an in-demand large facility, with both national and international users, its user operations cannot be delayed or disrupted. Therefore, using a Big Data framework that can transparently interface with the data

pipeline infrastructures of different beamlines, while translating them into a unified architecture, based on a universal Big Data and metadata machine-readable format in the background at the BDSC, is essential (figure 8). This approach, therefore, does not require the modification of legacy or current data pipelines linked to the beamlines, which is often either unproductive or infeasible, while fully achieving the BDSC objective of creating a superfacility at the SSRF. The BDSC Big Data Framework has been designed to be unconstrained by detector data formats and beamline data pipelines; as such, it can provide the highest level of architectural adaptability and resilience, while translating all large facility data pipelines, including legacy systems, into a universal machine-readable metadata-supported format, which can be used effectively for tagging all the Big Data, from any source at the SSRF, for AI training, and for IoT robotic operations. The BDSC can therefore be regarded as a real-world translation layer for AI, which can use the BDSC Big Data architecture to scout, capture, and uniformly reformat any data source in the real world into a universal machine-readable language, and then learn from it. The BDSC has augmented the MX Beamline capabilities, both qualitatively and quantitatively, thereby fully integrating the beamline within the BDSC Big Data framework, while accelerating the data evaluation process by a factor of three. Having successfully completed the pilot phase, the BDSC Big Data framework is now fully operational on the MX Beamline, and the beamline users are already fully benefitting from it; the BDSC is now extending its Big Data framework to all the other beamlines at the SSRF.
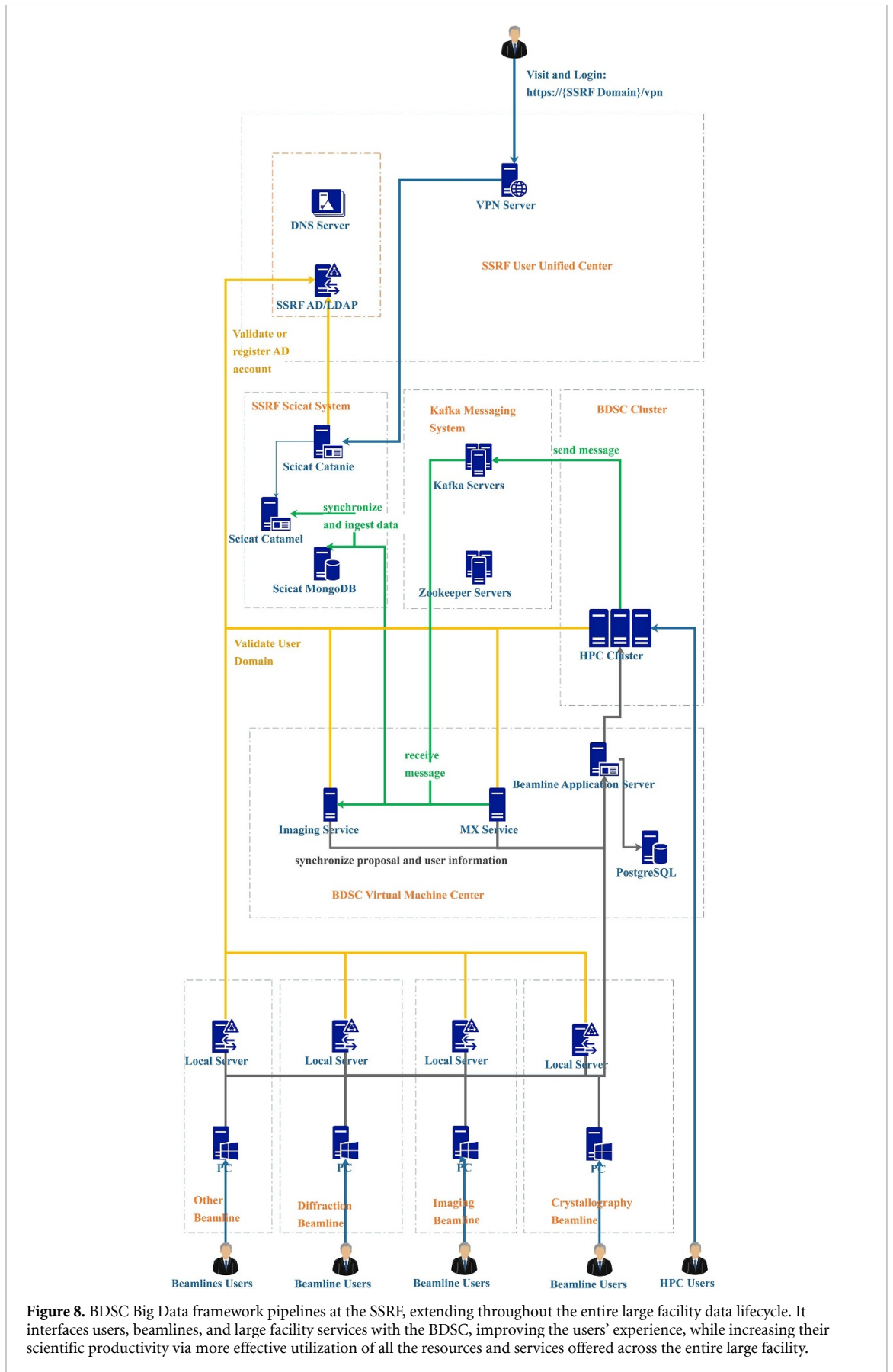
The Central Control Room at the SSRF controls the linear accelerator, booster, storage ring, and insertion devices, and monitors the beamline safety interlock status and some elements of beamline operation status. Each beamline at SSRF has its own local beamline control system, based on the Experimental Physics and Industrial Control System [74, 75, 105, 106], which controls all the equipment, ranging from beamline front-ends (including the insertion devices) to the end-stations, and monitors the operational status of the whole beamline [74, 75]. Currently, the beamline control systems output data as process variables, via the channel access protocol, to operator interface computers or beamline data acquisition pipelines, without remote monitoring or long-term data storage.

The BDSC is currently implementing a full integration of all beamline control systems into its centralized BDSC Big Data framework for full IoT robotic automation at all beamlines at the SSRF, assisted by AI, and implemented through ML and neural networks (NNs), thereby relieving the staff engineers and beamline scientists at the SSRF from demanding beamline control system tasks, while allowing them to focus more on user science. With its Big Data framework, the BDSC can provide online real-time remote monitoring of all the beamline systems at the SSRF, thereby leveraging the full potential of the metadata architecture deployed by the BDSC. Furthermore, the operational status of all beamlines, including motor positions, encoders, vacuum status, safety interlock information etc, will be collected, tagged, and uniformly formatted into a machine-readable format and linked to all beamline experimental dataflows and pipelines, centralized at the BDSC, to ensure a universal, unified, and centralized AI-assisted data lifecycle management system, providing feedback and controlling IoT robotic automation at the beamlines. Moreover, the BDSC will manage long-term data storage of all data from all the beamline control systems at the SSRF. Controlling and fully integrating the beamline control systems within the BDSC Big Data framework is the final step for the BDSC to implement a fully IoT robotically-automated and AI-assisted superfacility at the SSRF, based on Big Data analysis, visualization, and tagging from all unrelated data sources and pipelines.
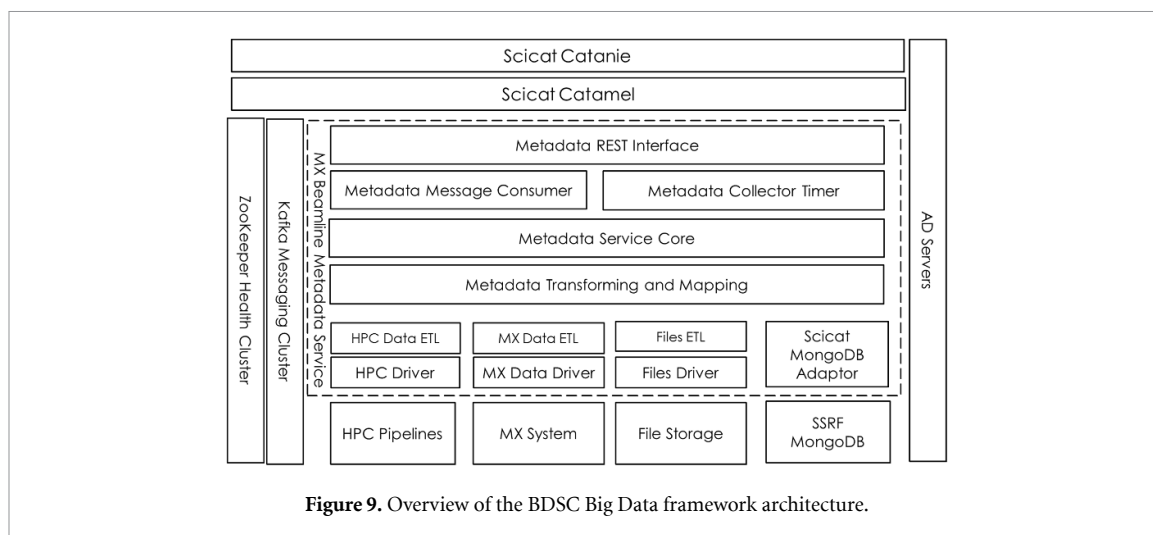
### 3.6. Big Data framework architecture

The BDSC Big Data framework architecture schematic is shown in figure 9. It has been designed and developed to cover all the beamlines within the entire SSRF, including the SSRF Phase-II beamlines under construction, where the MX Beamline has been selected as the pilot beamline to deploy the framework.

The BDSC Big Data framework manages each beamline data and metadata service at the SSRF, including SciCat and MongoDB modules, MX service modules (as well as the service modules for all the other beamlines at the SSRF), MX Beamline infrastructure (as well as the beamline infrastructures for all the other beamlines at the SSRF), HPC pipelines, and file storage, as well as the BDSC unified and universal AD system at the SSRF. The SciCat modules encompass SciCat Catanie and SciCat Catamel, two fundamental components of SciCat, providing a web portal, cataloguing, and data catalogue backend service, respectively. The users visit Catanie; then Catanie invokes Catamel to access the metadata. The BDSC has further implemented AD server components, which use ADs to manage the unified and universal account system for the HPC infrastructure and MX Beamline (as well as for all the other beamlines at the SSRF), software applications, storage, and all function services. The Zookeeper Health Cluster component coordinates the Kafka messaging cluster and other real-time data services. The BDSC Big Data framework architecture implemented these components to manage the overall messaging and health system of the entire infrastructure. The Kafka messaging cluster is used to report the pipeline's status. When a pipeline sends a message, the MX service is triggered to receive this message; it immediately starts the data injection process,

**Figure 8.** BDSC Big Data framework pipelines at the SSRF, extending throughout the entire large facility data lifecycle. It interfaces users, beamlines, and large facility services with the BDSC, improving the users' experience, while increasing their scientific productivity via more effective utilization of all the resources and services offered across the entire large facility.

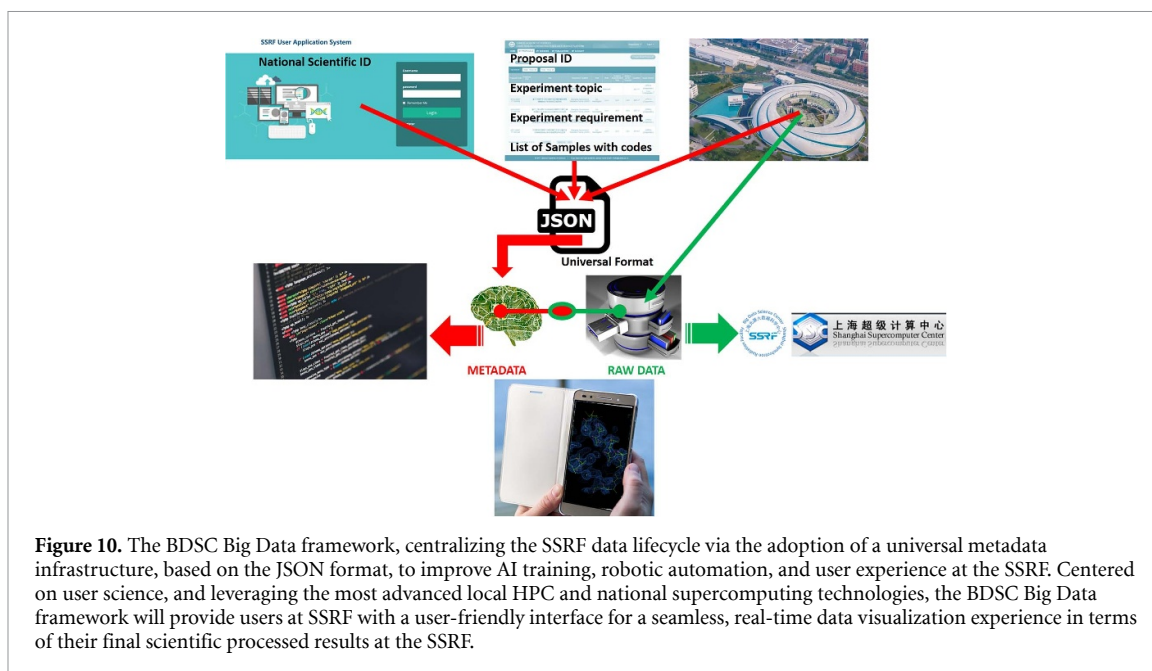**Figure 9.** Overview of the BDSC Big Data framework architecture.

implementing the extract-transform-load (ETL), and finalizing the injection of data into SciCat. The MX Metadata Service is an independent metadata service, scalable, and designed to be deployed for different beamlines. It contains several components, including the Metadata REST interface, Metadata Message Consumer, Metadata Collector Timer, Metadata Service Core, Metadata Transforming and Mapping, ETLs, and Drivers for the Kafka messaging cluster (see figure 9). The service can capture data from HPC pipelines, the MX system and file storage at the SSRF; it can then inject the corresponding metadata into MongoDB, either by using the adaptor, or going through the Catamel component using the mapping definitions. Finally, the data can be displayed using the SciCat Catanie component. In addition, the MX Metadata Service has been designed for maximum flexibility and reliability, because it not only processes real-time metadata from Kafka, but also processes any metadata accidentally lost due to network failures, which is a feature particularly relevant for users during unsupervised experiments or night shifts, when only limited local assistance is available.

The BDSC Big Data framework integrates the SciCat metadata model and MongoDB as components of a wider system, developed by the BDSC to centralize the entire SSRF data lifecycle, including proposals, datasets, files, samples, user IDs, users, and beamline operator parameters, experimental details, beamline and technique details, corresponding publications, etc, into a JSON universal metadata format for AI, ML, NNs, and IoT robotic automation for science (see figure 10). The BDSC Big Data framework provides a seamless data visualization experience of the final scientific processed results directly onto the thin client terminals, via user-friendly graphical user interfaces (GUIs) (figure 10). The BDSC has also implemented a docker deployment environment covering the entirety of the SSRF, integrating the Catamel, Cantanie, and Kafka containers.

However, SciCat has limitations that require improvement. Specifically, SciCat provides the model infrastructure, but does not provide specific metadata ETL or demon services; furthermore, SciCat does not provide scientific metadata structure and specifications; it also does not provide business process management, because data process flow management and data privileges are absent in SciCat. Moreover, SciCat does not provide pipeline services and injection from data sources to HPCs. Finally, SciCat provides only a very basic GUI. The BDSC has therefore developed new functionalities within its Big Data framework to address and overcome these limitations.

To develop the Big Data framework at SSRF, the BDSC used Java, C, and Bash as the main programming languages to support multi-platform and Linux batches, implementing Java Native Access/Java Native Interface and pipe interfaces in Java to communicate with the C and Bash processes. To provide comprehensive programming and configuration models for Java-based website and client-side enterprise applications, the BDSC used Apache Tomcat, Spring Framework, Kafka, Mybatis, and Alibaba DRUID as the supporting components. Furthermore, to manage related component dependencies to develop and build applications to run on servers equipped with OpenJDK 8.0, the BDSC used Maven and Gradle. All of these tools were in turn used to develop the SSRF-MX-Beamline-Service, and the SSRF-Cluster-Job-Interface framework.

As each beamline has specific software applications that run different tasks and generate unique metadata, it is difficult to develop only one service to handle all the different types of data pipelines. Therefore, the BDSC has developed a general framework and components that allow each beamline service component to be extended and used according to each specific and different requirement (figure 11). In the

**Figure 10.** The BDSC Big Data framework, centralizing the SSRF data lifecycle via the adoption of a universal metadata infrastructure, based on the JSON format, to improve AI training, robotic automation, and user experience at the SSRF. Centered on user science, and leveraging the most advanced local HPC and national supercomputing technologies, the BDSC Big Data framework will provide users at SSRF with a user-friendly interface for a seamless, real-time data visualization experience in terms of their final scientific processed results at the SSRF.
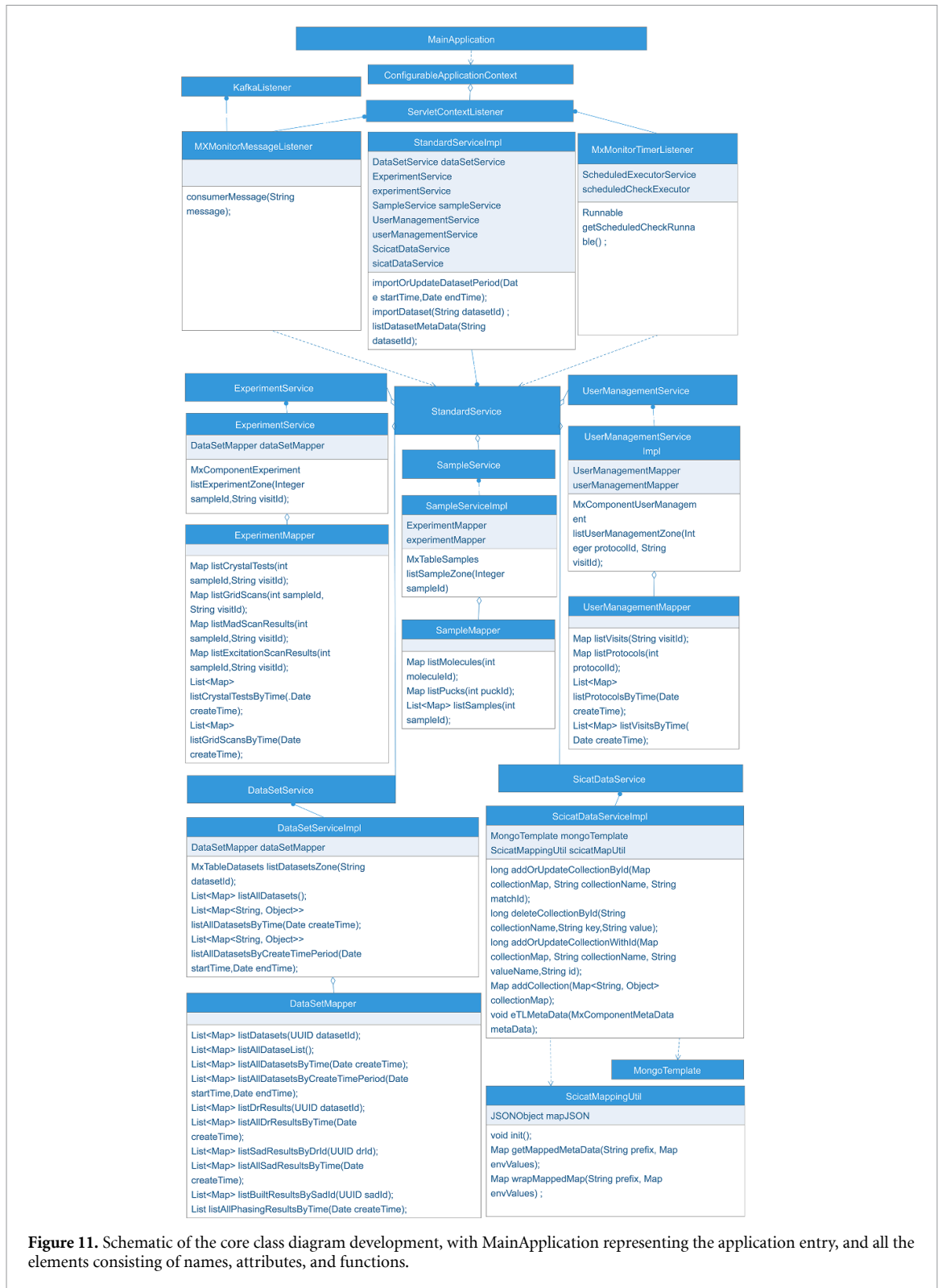
case of the MX beamline application, it provides a PostgreSQL DB server as the metadata source, granting a specific account with only select privileges. The SSRF-MX-Beamline-Service has been developed using data access object (DAO) models to search for data from the PostgreSQL DB server; this includes DatasetMapper, ExperimentMapper, SampleMapper, and UserManagementMapper. These DAO models provide specific interfaces that reflect the MX PostgreSQL DB structure to the Data Entity Java Beans. The implementation of DataSetService, ExperimentService, SampleService, and UserManagementService can then use the corresponding DAO models to publish MX metadata DAO interfaces.

The BDSC provides the file system that deploys the SSRF-MX-Beamline-Service, developed to invoke the file system interface and list the files with the proper rights, so that it can acquire the correct file attribute information in the dataset folders. The SciCat Catamel server provides the REST Interface, and Mongo DB functions as the metadata warehouse. The BDSC has developed the ScicatMappingUtil as a component of the SSRF-MX-Beamline-Service and the corresponding metadata definition files. The ScicatMappingUtil component can read and update the current metadata hash-map to match and generate new metadata in terms of metadata definition files. With the SSRF-MX-Beamline-Service, the BDSC also developed the ScicatDataService component, which provides ETL tools. It is designed to transfer multi-level MX metadata objects to a one-level string of key-value of the hash-map object; it then invokes ScicatMappingUtil to match and generate Dataset, Proposal, and OrigDatablock metadata in SciCat. Finally, it adds or updates Proposal, Dataset, and OriginDataBlock to Scicat MongoDB. Within the SSRF-MX-Beamline-Service, the BDSC has developed a standard service component that provides functions to list multi-level data objects, using MX metadata DAO interfaces to import metadata, together with Dataset ID and time period information, to SciCat. Within the SSRF-MX-Beamline-Service, the BDSC has also developed two listener components, the MXMonitorMessageListener, and the MxMonitorTimerListener. The MxMonitorMessageListener has been developed to read messages from the Kafka message cluster, and invoke standard services to import metadata with Dataset ID to SciCat. The MxMonitorTimerListener has been developed to import metadata, with time period information defined within the configuration file, to SciCat. The SSRF-Cluster-Job-Interface has been developed by the BDSC to handle the integration with the MX Slurm job file. When the MX beamline application completes the processing job on the BDSC HPC cluster, the SSRF-Cluster-Job-Interface is invoked, thereby sending a Dataset-ID message to the Kakfa message cluster, ensuring that the SSRF-MX-Beamline-Service will receive this Dataset-ID message in real-time.

The BDSC further introduced massive parallelization at SSRF, implementing data pipeline parallelizing task submissions from the MX Beamline system to the BDSC HPC clusters (figure 12). The HPC clusters can thus schedule and run multi-jobs with different pipelines in parallel, where each pipeline generates its metadata in parallel, while sending messages to the messaging cluster, also in parallel. Each message can then be received and processed by the MX Service, which can be deployed on multiple servers. Furthermore, metadata can be injected into MongoDB clusters via multithread connections.
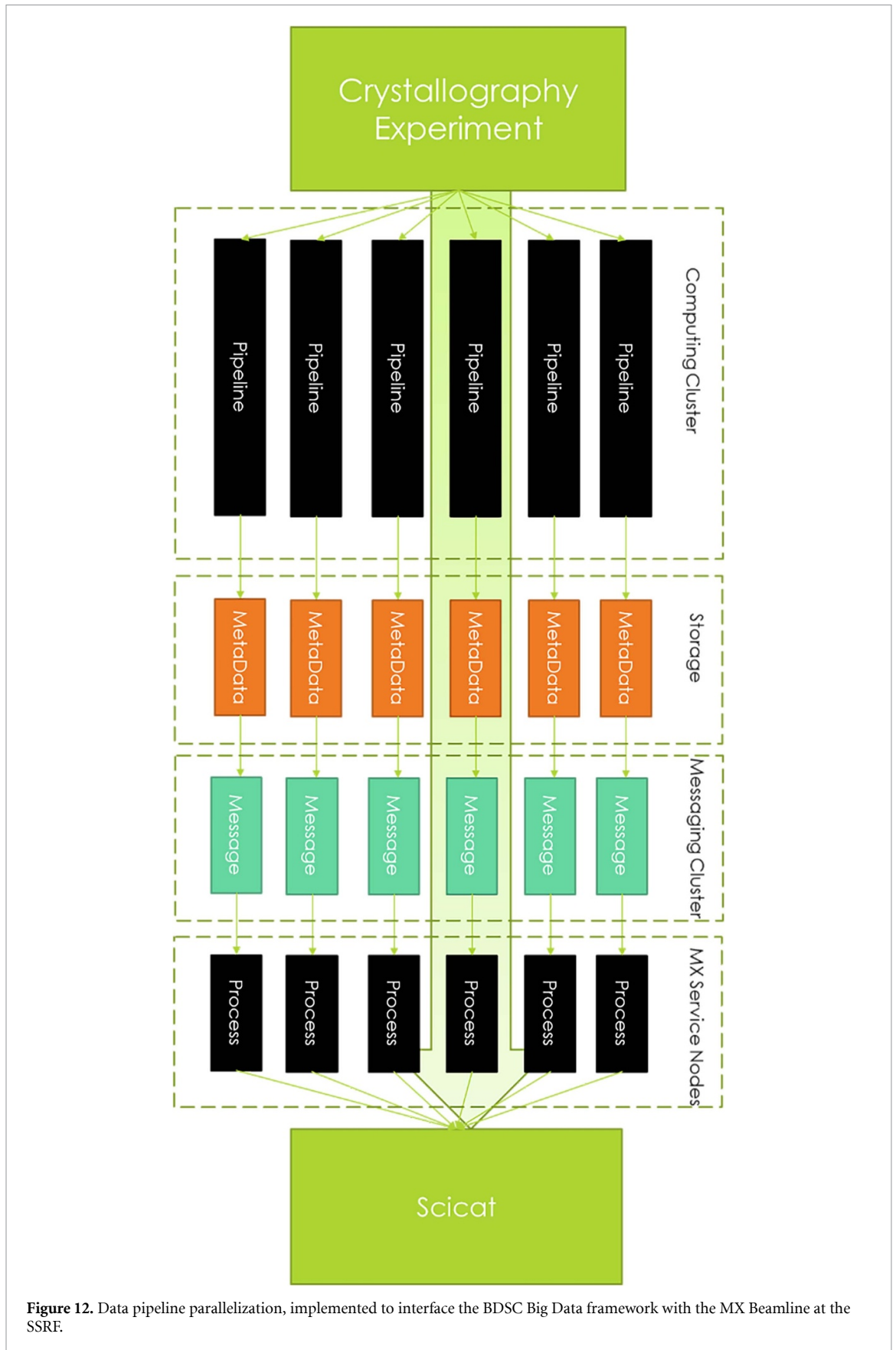
As previously mentioned, the MX Beamline has an existing automatic data-processing and experiment information management system, designated Aquarium, where a relational database, PostgreSQL, is used,

**MainApplication**

**ConfigurableApplicationContext**

**KafkaListener**

**ServletContextListener**

**MXMonitorMessageListener**

consumerMessage(String message);

**StandardServiceImpl**

DataSetService dataSetService
ExperimentService experimentService
SampleService sampleService
UserManagementService userManagementService
ScicatDataService sicatDataService

importOrUpdateDatasetPeriod(Date startTime,Date endTime);
importDataset(String datasetId) ;
listDatasetMetaData(String datasetId);

**MxMonitorTimerListener**

ScheduledExecutorService scheduledCheckExecutor

Runnable getScheduledCheckRunnable() ;

**ExperimentService**

**ExperimentService**

DataSetMapper dataSetMapper

MxComponentExperiment
listExperimentZone(Integer sampleId,String visitId);

**ExperimentMapper**

Map listCrystalTests(int sampleId,String visitId);
Map listGridScans(int sampleId, String visitId);
Map listMadScanResults(int sampleId,String visitId);
Map listExcitationScanResults(int sampleId,String visitId);
List<Map> listCrystalTestsByTime(.Date createTime);
List<Map> listGridScansByTime(Date createTime);

**StandardService**

**SampleService**

**SampleServiceImpl**

ExperimentMapper experimentMapper

MxTableSamples
listSampleZone(Integer sampleId)

**SampleMapper**

Map listMolecules(int moleculeId);
Map listPucks(int puckId);
List<Map> listSamples(int sampleId);

**UserManagementService**

**UserManagementService Impl**

UserManagementMapper userManagementMapper

MxComponentUserManagement
listUserManagementZone(Integer protocolId, String visitId);

**UserManagementMapper**

Map listVisits(String visitId);
Map listProtocols(int protocolId);
List<Map> listProtocolsByTime(Date createTime);
List<Map> listVisitsByTime( Date createTime);

**DataSetService**

**DataSetServiceImpl**

DataSetMapper dataSetMapper

MxTableDatasets listDatasetsZone(String datasetId);
List<Map> listAllDatasets();
List<Map<String, Object>> listAllDatasetsByTime(Date createTime);
List<Map<String, Object>> listAllDatasetsByCreateTimePeriod(Date startTime,Date endTime);

**DataSetMapper**

List<Map> listDatasets(UUID datasetId);
List<Map> listAllDatasetList();
List<Map> listAllDatasetsByTime(Date createTime);
List<Map> listAllDatasetsByCreateTimePeriod(Date startTime,Date endTime);
List<Map> listDrResults(UUID datasetId);
List<Map> listAllDrResultsByTime(Date createTime);
List<Map> listSadResultsByDrId(UUID drId);
List<Map> listAllSadResultsByTime(Date createTime);
List<Map> listBuiltResultsBySadId(UUID sadId);
List listAllPhasingResultsByTime(Date createTime);

**SicatDataService**

**ScicatDataServiceImpl**

MongoTemplate mongoTemplate
ScicatMappingUtil scicatMapUtil

long addOrUpdateCollectionById(Map collectionMap, String collectionName, String matchId);
long deleteCollectionById(String collectionName,String key,String value);
long addOrUpdateCollectionWithId(Map collectionMap, String collectionName, String valueName,String id);
Map addCollection(Map<String, Object> collectionMap);
void eTLMetaData(MxComponentMetaData metaData);

**MongoTemplate**

**ScicatMappingUtil**

JSONObject mapJSON

void init();
Map getMappedMetaData(String prefix, Map envValues);
Map wrapMappedMap(String prefix, Map envValues) ;

**Figure 11.** Schematic of the core class diagram development, with MainApplication representing the application entry, and all the elements consisting of names, attributes, and functions.
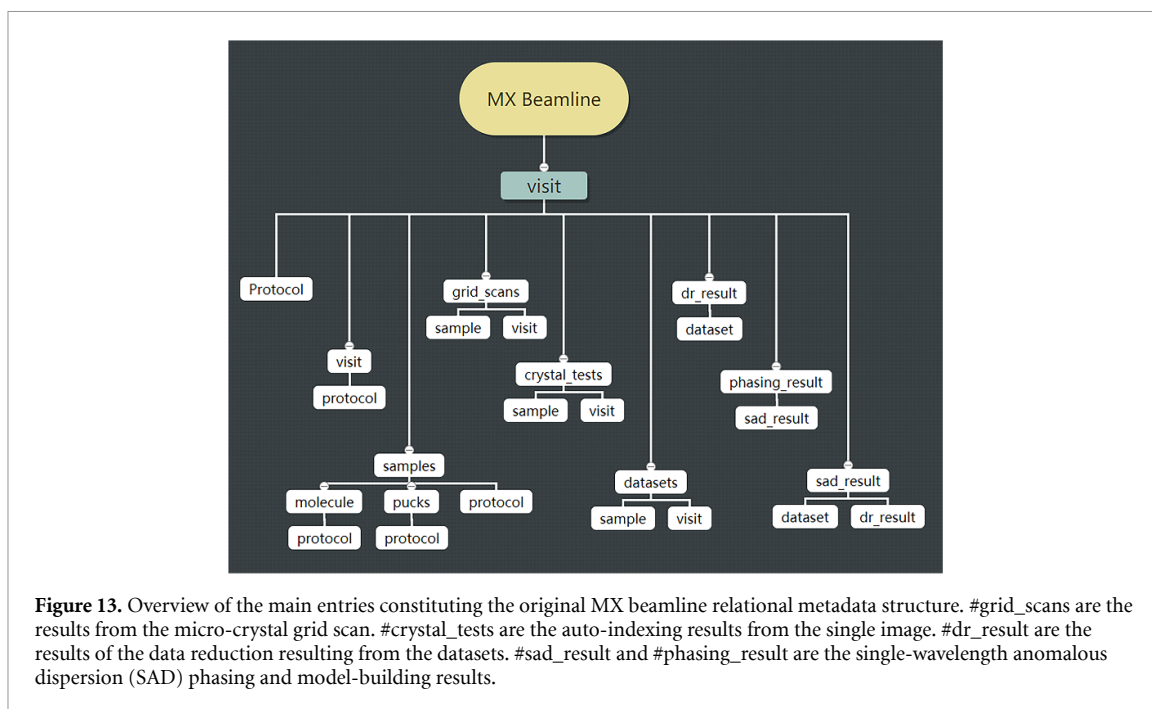
with an entry rate of approximately 400 metadata per dataset into the database. The original metadata structure of the MX Beamline is shown in figure 13. The metadata are hierarchically organized, based on the #visit (where the hashtag symbol # is used to highlight all labels used for the metadata), which labels each independent experiment session, from which all the metadata can be catalogued into: #protocol, #visit, #samples, #grid_scans, #crystal_tests, #datasets, #dr_result, #sad_result and #phasing_result.

To deploy the BDSC Big Data framework without affecting the normal operations of the beamline, the BDSC developed the MX Metadata Service, using ETL to capture all the MX Beamline metadata and inject them into SciCat and MongoDB. The metadata structure from the MX Beamline was then reorganized and re-categorized, thereby creating a unified, universal, and uniform framework for the SciCat metadata

**Figure 12.** Data pipeline parallelization, implemented to interface the BDSC Big Data framework with the MX Beamline at the SSRF.

ingestion system. Figure 14 presents an overview of this framework, where the #scientific_metadata contains the most relevant and important experimental information. As shown in figure 14, #user_management includes #visit, #protocol, and other metadata originating from the Proposal Management System and Unified Account Authentication Service; #sample includes detailed information regarding samples,

**Figure 13.** Overview of the main entries constituting the original MX beamline relational metadata structure. #grid_scans are the results from the micro-crystal grid scan. #crystal_tests are the auto-indexing results from the single image. #dr_result are the results of the data reduction resulting from the datasets. #sad_result and #phasing_result are the single-wavelength anomalous dispersion (SAD) phasing and model-building results.

originating from the Proposal Management System and manual input from users; #experiment includes the information linked to beamline operational status, methods, detectors, experimental parameters, which originate from the beamline control system, and all the GUIs interacting with the users at the SSRF; #dataset includes the information from all the data pipelines and the tuning parameters at SSRF used for data processing, analysis, modeling, simulation, and reconstruction, as well as distributed storage mapping and the tagging system for AI-assisted technologies implementing ML and IoT robotic automation, linking them with the raw data and final processed results for data visualization, all converging into the BDSC (figure 10).

Moreover, the BDSC has developed a customized SQL to NoSQL mapping module, as shown in figure 15. The metadata mapping service supports mapping extensions from all the MX Beamline processed metadata, including Proposal Mapping, Dataset Mapping, and OrigDataBlock Mapping.

The MX Beamline Metadata Service has been developed by the BDSC to have a loose coupling with MX Beamline software applications. Thus, using the API, it can access and read the SQL database also simply with the 'select' permission. The BDSC Big Data Framework exhibits a great degree of flexibility, since, in the case of MX Beamline applications, the parameters and database have been changed or upgraded; the metadata service can then automatically generate a new format that can be injected into SciCat. This only requires the creation of a link with the primary key and the hierarchical JSON format. Moreover, it is equipped with a static metadata mapping configuration file, which can be modified at any time to remap the SciCat ingestion system, should SciCat itself be upgraded.

## 4. AI and supercomputing

### 4.1. The role of AI/ML

As stated above, the deep learning revolution [6] has transformed various application domains, including those dealing with image datasets. Although the techniques within the domain of ML include various techniques, such as deep learning, it is not an overstatement to say that deep learning is the primary reason for recent advances in AI. ML techniques, in general, can be categorized into (a) supervised, (b) unsupervised, (c) self-supervised, and (d) reinforcement learning. Each of these techniques have their pros and cons; for more detailed descriptions of these models, please refer to [107].

Two of these methods are worth discussing here. Firstly, supervised learning, which includes most of the deep learning techniques; this is a method of teaching an ML system (or model) using examples. For example, a specific feature in an x-ray image can be identified if the system can be provided with several examples of such features. The examples with marked features are referred to as the training dataset, and the marked features are referred to as labels. The success of deep learning relies on two aspects: the amount of training data, and labels (or the quality of labels). Although obtaining labels for generic images, such as animals or day-to-day objects, is possible by means of citizen science or similar efforts, securing labels for

**Figure 14.** Detailed overview of the new MX beamline non-relational metadata structure, following its redesign within the BDSC Big Data framework at the SSRF, using the SciCat model system.

scientific datasets is a challenging task. It not only requires expert knowledge, but also significant time to label them to ensure that they can be utilized by various ML frameworks.

Although unsupervised and self-supervised techniques can offer advantages over supervised methods, their efficiency, particularly with regard to output, is arguably inferior to that of supervised methods, at least with respect to the current set of problems. Admittedly, these conclusions are often made on a case-by-case basis, particularly in the absence of any generic benchmarking system. Although frameworks like MLPerf provide a generic picture, they do not include the scientific domain; therefore, it is difficult to draw any firm conclusion.

Both BDSC and SSRF recognize this issue, and therefore emphasize the collection of curated data for all experiments. Accordingly, the data repository is carefully monitored to ensure that the outputs from the superfacility are fully qualified, high-quality labelled datasets. Admittedly, the exact volume of the labelled datasets per beamline is yet to improve, but we are confident that soon it will become populated enough to

**Figure 15.** Relational to non-relational MX Beamline metadata mapping structure within the BDSC Big Data framework.

use for ML. Although we are actively investigating other techniques, such as self-supervised techniques, we would like to ensure that we have a fail-safe option.

In addition to applying ML on scientific datasets, the application of ML for actually performing experiments using robots is another focus of the BDSC. Automation using robots relies on reinforcement learning and multi-agent policy-based systems, rather than supervised techniques. The resulting systems are complex, and require long training periods before they can be deployed. We refer the readers to [108] for a more detailed description of reinforcement learning.

ML systems differ in terms of the way in which they are implemented and consume data. Although one cannot enforce a common format for data consumption, it is possible to ensure that datasets are stored in a readily consumable format. For instance, the majority of beamlines produce their final results in proprietary formats, which often involves a significant amount of time for conversion, and the import and export of data formats. Both BDSC and SSRF recognize the challenges in this area; as such, the data repository is free from any proprietary format. It contains very generic formats, such as images and tables. These formats are readily ingested into any AI/ML framework, and easily readable. Such common formats also ensure that the techniques are easily transferable, e.g., via transferred learning approaches.

Finally, as discussed above, metadata capturing is a significant activity within the Superfacility. Examples of the AI/ML cases presented in section 3 highlight how data collected from the superfacility can be helpful, from dataset analyses to controlling the beamlines. Therefore, the BDSC is now creating an architecture for the deployment of ML at the metadata level; hence, metadata is regarded as having similar significance to scientific datasets.

To estimate the AI/ML resources required for the SSRF, the BDSC ran statistical calculations based on the HPC usage from January 2020 to November 2020. The current users of the BDSC are BL08U (soft x-ray spectromicroscopy and soft x-ray interference lithography), BL09U (high-resolution and wide energy range photoemission spectroscopy—Dreamline), BL10U (time-resolved ultra-small-angle x-ray scattering and biosafety P2 protein crystallography), BL13W (x-ray imaging and biomedical applications), BL17U (macromolecular crystallography), and BLCTL (tomography) beamlines. Users from different beamlines are also submitting their tasks, including graphic data processing, material spectral analysis, molecular dynamics simulations, and *ab initio* quantum simulations.

BL17U is the main implemented user of the HPC at the BDSC. The data are processed automatically through six pipelines on the HPC, indicating that at least three computation nodes are usually used to parallel process the tasks submitted from BL17U.

The most frequent submitted data processing work involves the Aquarium pipelines from BL17U, with X-ray Detector Software (XDS) programs. XDS programs are used for processing single-crystal monochromatic diffraction data, recorded by the rotation method for the 'reduction' of two-dimensional data images ('frames') obtained from crystals irradiated with monochromatic x-rays. XDS programs can process data images from CCD-, imaging-plate-, multiwire-, and pixel-detectors in a variety of formats, as well as from multi-segment detectors assembled from several rectangular components in arbitrary arrangements. XDS automatically splits its tasks for concurrent execution via several remote nodes in the networked file system (NFS) environment at the BDSC HPC cluster; each computation node comprises a shared memory multiprocessor system.

The most frequently submitted simulation software is Large-scale Atomic/Molecular Massively Parallel Simulator (LAMMPS). LAMMPS is free and open-source software, distributed under the terms of the GNU General Public License, and developed by Sandia National Laboratories (a National Nuclear Security Administration research and development laboratory in the United States). LAMMPS comprises a classical molecular dynamics code that focuses on material modeling, and uses the message passing interface for parallel communication. LAMMPS can potentially be used for solid-state materials (metals, semiconductors), soft matter (biomolecules, polymers), and coarse-grained or mesoscopic systems; in addition, it supports accelerated performance on the graphics processing unit (GPU) nodes.

A total of 36 187 tasks were submitted in the BDSC HPC cluster, as of November 2020. On average, 3500 tasks are submitted per month. The compatibility of the HPC software is improved according to user requirements from different beamlines. For example, we have solved some Linux library dependency bugs, and improved the submission procedure. Consequently, beamline users are more likely to use the HPC at the BDSC, rather than their local clusters, at the beamlines. On average, approximately 15%–20% of the computational resources of the BDSC are consumed by all the submitted tasks. Here, the CPU available peaks are estimated to be at 80% capacity of all CPUs in the BDSC HPC cluster. This is due to the consideration that if all the CPUs are at capacity, every submitted task needs to wait in a queue, which could typically take 1–3 h. Therefore, the computation peaks are set to be a bit lower than the ideal computation capacity of the HPC at the BDSC. Most submitted tasks (about 80%) require at least one computation node. Other submitted tasks (about 20%) require approximately three computation nodes. The largest require 12 nodes, i.e., approximately one third of the available CPU resources of the BDSC HPC cluster (maximum 46 nodes). This implies that additional CPU resources are required, due to the sudden emergence of high-CPU-consuming tasks.

According to our estimation, two beam lines consume approximately 15% of the computational resources of the HPC cluster at BDSC. The current computational resources are estimated to be available for 13 beam lines. If AI/ML computation tasks are also included to optimize work from each beam line, and their computation usage is similar to the current workloads, the computational resources must therefore be doubled. This is due to the consideration that for each submitted data processing task, another AI/ML accelerating task must be submitted. For example, an AL/ML accelerating algorithm would quickly predict the score (quantity) of the processed graphic data, which would provide additional information to the users. Users can make a decision to continue their work (if the predicted score is good enough) or stop and restart a new experiment by selecting another sample (if the predicted score is totally unacceptable). Consequently, according to our estimation, the current HPC resources are available to support AI/ML applications for 6–7 beamlines.

Currently, two AI/ML-oriented projects, aiming to accelerate and optimize crystallography and diffraction image processing for the BL17U1 and BL13W1 beamlines, are being conducted at the BDSC. Based on the collected metadata and raw images, both unsupervised and supervised learning approaches are implemented for different imaging processing purposes, as follows.

(a) The first includes cluster analysis for the automatic classification of crystallography images [109]. Various XDS images of different macromolecules, such as proteins, polymers, nucleic acids, etc, are collected from previous user experiments at the beamline, and stored in the database. Images in the database are automatically clustered into different hierarchical groups, based on their pattern similarity. By establishing a proper hierarchical tree, the XDS images from a new experiment would automatically be classified into an existing group. Users would be able to compare their processing parameters (such as the exposure time, distance, energy, etc) with the existing metadata in this specific group. The image quality scores of the new XDS images would also be compared with the average score and average user preference (using

**Figure 16.** The BDSC supercomputing architecture. The BDSC can elastically scale-out the scientific data workloads of users at the SSRF to the national supercomputer centers. The BDSC is directly interfaced with the SSC, and it is designed to modularly scale-out its users' workloads to any other supercomputer nationwide.

the like/dislike system) in this specific group. Users can then decide whether to keep their image data for further processing, or restart a new experiment with the recommended parameters.

(b) The other involves real-time image quality prediction [110]. The image quality scores fluctuate, owing to the real-time equipment and environmental conditions at the beamlines. Based on the time series prediction algorithms of artificial NN models, the image quality variation pattern in each experiment would be predicted in real-time, according to the current obtained and processed imaging qualities. For example, the quality of the first 50% of the images are calculated, followed by the automatic prediction of the remaining 50%. Given that users prefer images with high quality scores, they would be recommended to process the image data with predicted high quality in the remaining 50% of the images in the data processing pipelines. The high-quality standard is selected to be above the average image quality score for a specific image group, according to its image pattern similarity in the automatic cluster analysis. This time series prediction strategy could significantly accelerate the entire image processing by reducing the low-quality image processing time. The daily quality pattern of every beamline would be stored in the database, and automatically analyzed. The predicted quality variation patterns of the beamlines would be compared with their real running status. It would also be helpful to measure the instrument performance, and to detect possible instrument failure in a full maintenance cycle.

### 4.2. Supercomputing

The BDSC has been designed to interface with national supercomputers. It can access national supercomputer centers through their public interfaces, allowing the entire SSRF to transparently submit tasks, jobs, and transfer data to these national supercomputers. To seamlessly adapt all the SSRF data pipelines with a single interface to different supercomputer facilities, the BDSC is designing a unified, universal, and uniform middleware platform, capable of elastically scaling-out the workload from the local BDSC HPC clusters to the national supercomputer centers, in a manner entirely transparent to the users (figure 16).

The BDSC has designed the SSRF data pipelines to run on HPC clusters, agnostic with regard to their locations, interfacing delocalized and distributed HPC clusters with the software applications within beamlines, and services via a unified and centralized middleware platform. The SSRF has been connected to

the Shanghai Supercomputer Center (SSC) through the VPN and 1 Gb fiber-optic network, with the MX Beamline software applications and metadata services deployed on virtual machines at the BDSC, capable of accessing the supercomputer through LANs, internet VPNs, and fiber-optic networks. The MX Beamline software applications integrate the pipeline scheduler; they can query and submit processing tasks to idle supercomputers, using FTPS/SFTP protocols to retrieve data. The pipeline then finalizes its job, and sends a message to the messaging cluster. The MX Metadata Service monitors and receives the messages automatically, followed by managing the corresponding metadata that are properly injected into SciCat.

With the successful and complete integration of the MX Beamline within the BDSC Big Data framework, the BDSC is now extending its Big Data framework to cover all the other beamlines at the SSRF.

## 5. Future developments

The BDSC aims to further expand and upgrade its current infrastructure with proficient and comprehensive solutions, rich in features, with several prototyping solutions already in development.

After the successful integration of the MX Beamline at the SSRF into the BDSC Big Data framework, the BDSC is now extending its framework to include the entire SSRF. It is extending its operations to the x-ray imaging and Laue micro diffraction beamlines, to centralize all the Big Data and metadata from all the unrelated sources at the SSRF through a unified and uniform data pipeline, translating them into a uniform and universal machine-readable format, which can be directly used for AI training through ML.

Moreover, the BDSC is now focusing on dramatically increasing user interactions with the Big Data and metadata framework, to provide the AI for scientific research at the SSRF with more reliable supervised training. Therefore, the BDSC is developing a new universal user interface that can track and tag SSRF user interactions through the entire data lifecycle of a superfacility, while providing a smoother and more immersive user experience, and enhancing user-friendly features.

The BDSC is further engaged in improving the link between the overall user data lifecycle at the SSRF, and the final published results. Therefore, the BDSC is upgrading and enhancing its DOI framework, providing both manual and automatic avenues for the collection of information on user publications, and connecting them with all the corresponding data produced by user experiments at the SSRF. To achieve this objective, the BDSC is interfacing its Big Data framework directly with data citation solutions such as DataCite [111, 112], and adopting findable, accessible, interoperable, and reusable [113] data principles as an integral part of its infrastructure.

The BDSC is also investing substantial resources to contribute to the development of AI for science at large facilities, which will be based on the BDSC Big Data framework, which was engineered to increase the training of NNs; this training will benefit from extensive user interactions and supervision, with the aim of perfecting AI functionalities based on real user cases, along with the scientific research requirements of users, at large facilities. The AI at SSRF will therefore be fed with all the metadata from the entire large facility lifecycle, from proposal to publication, which, as described in previous sections, will be collected by our tagging system from all the unrelated sources at the SSRF; these will be used to train the AI through ML. To enhance the effectiveness and reliability of the AI's supervised training, the BDSC has launched the Research Partnership Program initiative at the SSRF, with the aim of involving expert academicians and users in the training of the AI for science at SSRF. This will enable an extensive knowledge transfer from more experienced users to less experienced users at the SSRF by means of AI-assisted technologies, thus increasing the overall scientific productivity of all users at the SSRF, irrespective of their experience level with large facility operations. The Research Partnership Program will also be extended to support industrial applications at SSRF, thereby also strengthening the bridge between discoveries and inventions at large facilities. Moreover, through this initiative, the BDSC will organize workshops, conferences, and university courses aimed at educating future scientists, providing them with a wider and deeper knowledge of Big Data for science at large facilities.

Furthermore, the BDSC is focusing its resources on large-scale CPU to GPU migrations at the SSRF, while utilizing its GPU cluster to enhance the benefits provided by the adoption of GPU at the SSRF. This migration is particularly useful with regard to those data-intensive beamlines still using outdated applications, with algorithms optimized only for CPU. Transitioning to GPU-optimized applications will be of great significance for scientific research; for instance, the GPU-based ultra-fast online processing framework [114] at the Karlsruhe Institute of Technology can achieve high-quality tomographic reconstructions at a speed 37 times higher than traditional CPU-based solutions. Another example is provided by scalable heterogeneous adaptive real-time ptychography, a high-performance GPU-based software environment at the advanced light source, which expedites time-consuming ptychography reconstructions, and provides quasi-real-time feedback on the results [23, 115]. The BDSC is already deploying GPU migration for the x-ray imaging beamlines at the SSRF. The new GPU-based tomographic reconstruction software at the SSRF is, in fact,

designed and developed for HPC environments with client/server architectures, which can be fully integrated into automatic data pipelines within the BDSC Big Data framework.

One of the most csignificant prospects for the BDSC comprises the robotic automation of the entire SSRF, with a particular focus on beamlines. Accordingly, the BDSC is devising a large-scale strategy to implement robotic automation through IoT at the SSRF, further equipping the beamlines with sensors and automatic tagging capabilities, thereby deploying full, facility-wide automation. This will also include the full integration of the SSRF beamline control systems within the BDSC infrastructure, along with a real-time unmanned automatic AI-assisted beamline feedback system, trained by users in terms of data interpretation and experimental configurations, and the beamline operatives in terms of beamline configurations, alignments, and setups. This is fundamental to automatically set up the best possible experiments and measurements for users at the SSRF in real-time, thus allowing them to return to their home institutions having collected complete, directly usable data, invaluable to their research and publications. Moreover, fully automatic experiments will further contribute to real-time experiments, data pipelines, and data interpretation at the SSRF, overcoming the need for pausing user measurements for long periods.

As previously mentioned, the BDSC is also advocating for the adoption of a national unified scientific ID. This will simplify the user experience and provide the constituent base for a national large scientific facility grid based on a unified access, linking, and sharing system. The BDSC Big Data framework can be extended to other large facilities; therefore, users can manage all their data from different facilities using a common and unified interface. This will support wider knowledge sharing between users at large facilities nationwide. Accordingly, any new scientific theory, algorithm, application, or framework can be easily and seamlessly pushed as a grid update for the Big Data framework nationwide. This, in turn, will automatically update the local frameworks of large facilities nationwide, while making scientific knowledge upgrades immediately available to all users, thereby augmenting horizontal user knowledge transfer with up-to-date science.

In addition, the BDSC intends to actively engage in promoting various initiatives aimed at adapting and migrating the BDSC Big Data Framework to other facilities worldwide.

Finally, the BDSC will support, with its Big Data framework, all the research works under the umbrella of the Zhangjiang Laboratory, offering its scientific Big Data capabilities to strengthen multidisciplinary research within the Zhangjiang Laboratory infrastructure, further consolidating and promoting co-operation and sharing, connecting different types of expertise within the wide network of research groups and experts at the Zhangjiang Laboratory.

## 6. Conclusions

The BDSC aims to establish itself as a solid international reference within the field of Big Data science at large facilities. Therefore, it has implemented a Big Data framework to uniformly format all the unrelated data pipelines at large scientific facilities into a machine-readable format for AI training through ML, thereby providing feedback to the IoT at the beamlines to achieve fully robotically automated beamlines at the SSRF.

With the clear target of deploying a superfacility at SSRF, the BDSC has implemented an advanced Big Data and metadata framework for scientific research, which constitutes the backbone of the superfacility itself, to augment the SSRF with AI using ML and robotic automation via IoT.

Another objective of the BDSC includes providing a state-of-the-art user experience at the SSRF by means of real-time Big Data analysis, interpretation, and visualization. The BDSC has therefore implemented a metadata tagging system at the SSRF for AI, capable of supporting robotic automation to reconfigure user experiments in real time at the beamlines. This will enable users to focus their scientific investigations on those results that produce scientifically meaningful and valuable data for their final publications, while dramatically increasing overall scientific productivity at the SSRF, together with its data infrastructure performance.

Accordingly, the BDSC acts as a scientific algorithm cloud, strongly science-centric, to solve modern scientific problems, while interfacing with national Supercomputer facilities, so as to maximize the scientific achievements of all users at the SSRF.

Moreover, it aims to become a major contributor to the international scientific community, by creating a fully automated, AI-assisted Big Data solution, tailored to the needs of large scientific facilities, advancing their international state-of-the-art capabilities with a scientific Big Data framework, capable of being augmented with AI and robotic automation.

The BDSC has therefore enhanced the deployment of AI for science at the SSRF, by virtue of the implementation of the BDSC Big Data and metadata framework, thereby facilitating AI training through ML, while prototyping based directly on real user experiments in order to guarantee that it is user-science-centric. The BDSC Big Data framework has introduced data centralization at the SSRF via universal and uniform data pipeline formatting, which can translate any unrelated data source at the SSRF

into machine-readable data; these data can be directly accessed for AI training. To further support SSRF users and their scientific achievements, the BDSC has made available its Big Data framework to host user software, optimizing them for HPC architecture. Therefore, users can use their own software onsite at the SSRF accessing this framework to finalize their synchrotron studies and package them into a scientific publication. The BDSC has also devised a future upgrade for the hosting system of its Big Data framework, whereby users will be offered the possibility to adapt and migrate their own software to the HPC architecture. Therefore, the BDSC will optimize SSRF user analysis software to utilize all the advanced, massive capabilities offered by present and future BDSC HPC clusters. The process of user software optimization will also include the development of a layer that can directly and efficiently interface the user's software with the BDSC Big Data and metadata framework, ensuring that the outcomes of the optimized user software are machine-readable for the direct and immediate training of the AI at the SSRF.

The ambitious architecture of the BDSC will create a solid foundation from which to build a larger Big Data framework for science, extending throughout the entire Zhangjiang Laboratory, which will promote the easy and rapid integration of advanced scientific software applications, while delivering a facility-wide multimodal approach for the benefit of users. Moreover, the BDSC multimodal strategy will be crucial for integrating those beamlines that do not produce Big Data within the Big Data scientific framework, but would benefit from a more complete data interpretation, as offered by the combined analysis approach.

## Data availability statement

The data that support the findings of this study are openly available at the following DOI: https://doi.org/10.6084/m9.figshare.14312819.v1

## Acknowledgments

## ORCID iDs

Chunpeng Wang ⓘ https://orcid.org/0000-0002-7026-7726
Feng Yu ⓘ https://orcid.org/0000-0002-9502-3277
Yiyang Liu ⓘ https://orcid.org/0000-0002-0897-1149
Xiaoyun Li ⓘ https://orcid.org/0000-0003-1976-6180
Jige Chen ⓘ https://orcid.org/0000-0001-6954-6413
Jeyan Thiyagalingam ⓘ https://orcid.org/0000-0002-2167-1343
Alessandro Sepe ⓘ https://orcid.org/0000-0002-2320-9398

## References

[1] Bell G, Hey T and Szalay A 2009 *Science* **323** 1297
[2] Foster I, Ananthakrishnan R, Blaiszik B, Chard K, Osborn R, Tuecke S, Wilde M and Wozniak J 2015 *Big Data High Perform. Comput.* **26** 117–32
[3] Assunção M D, Calheiros R N, Bianchi S, Netto M A S and Buyya R 2015 *J. Parallel Distrib. Comput.* **79–80** 3–15
[4] Kumar A, Boehm M and Yang J 2017 *Proc. 2017 ACM Int. Conf. on Management of Data* (*Chicago, Illinois, USA*) pp 1717–22
[5] Toby B H, Gürsoy D, De Carlo F, Schwarz N, Sharma H and Jacobsen C J 2015 *Synchrotron Radiat. News* **28** 15–21
[6] Sejnowski T J 2018 *The Deep Learning Revolution* (Cambridge, MA: MIT Press)
[7] Maddison D R, Swofford D L and Maddison W P 1997 *Syst. Biol.* **46** 590–621
[8] Klosowski P, Koennecke M, Tischler J Z and Osborn R 1997 *Physica* B **241–243** 151–3
[9] Könnecke M *et al* 2015 *J. Appl. Crystallogr.* **48** 301–5
[10] Chen C L P and Zhang C-Y 2014 *Inf. Sci.* **275** 314–47
[11] Ushizima D M *et al* 2016 *JOM* **68** 2963–72
[12] Hey T, Tansley S and Tolle K M 2009 *The Fourth Paradigm: Data-intensive Scientific Discovery* (Redmond, WA: Microsoft Research)
[13] Wang C, Steiner U and Sepe A 2018 *Small* **14** 1802291
[14] Hexemer A, Parkinson D and Tull C 2015 *Synchrotron Radiat. News* **28** 2–3

[15] Shane C R, Declerck T, Draney B, Lee J, Paul D and Skinner D CUG 2017 (available at: https://cug.org/proceedings/cug2017_proceedings/includes/files/pap165s2-file1.pdf)

[16] Troutman K Superfacility framework advances photosynthesis research (available at: www.nersc.gov/news-publications/nersc-news/science-news/2019/superfacility-framework-advances-photosynthesis-research/)

[17] Bard D The superfacility concept (available at: https://anchor.fm/nersc-news/episodes/The-Superfacility-Concept-Debbie-Bard-Interview-e5a5th)

[18] Black D Superfacility—how new workflows in the DOE office of science are changing storage requirements (available at: https://insidehpc.com/2016/05/superfacility/)

[19] Kincade K ESnet paves way for HPC 'superfacility' real-time beamline experiments (available at: www.es.net/news-and-publications/esnet-news/2015/esnet-paves-way-for-hpc-superfacility-real-time-beamline-experiments/)

[20] Snavely C The NERSC superfacility project: a technical overview (available at: www.nersc.gov/assets/GPUs-for-Science-Day/14-cory-snavely.pdf)

[21] Bard D Supercomputing and the scientist: how HPC and large-scale data analytics are transforming experimental science (available at: https://insidehpc.com/2019/09/asupercomputing-and-the-scientist-how-hpc-and-analytics-are-transforming-experimental-science/)

[22] Bard D and Snavely C Superfacility and gateways for experimental and observational data (available at: www.nersc.gov/assets/Uploads/1200-Superfacility-presentation.pdf)

[23] Donatelli J *et al* 2015 *Synchrotron Radiat. News* **28** 4–9

[24] Parkinson D Y *et al* 2016 *AIP Conf. Proc.* **1741** 050001

[25] Bethel E W 2017 *2017 IEEE 13th Int. Conf. on e-Science (E-science) (Auckland, New Zealand)* pp 462–4

[26] NERSC superfacility (available at: www.nersc.gov/research-and-development/superfacility/)

[27] Thomas M, Kleese-van Dam K, Marshall M J, Kuprat A, Carson J, Lansing C, Guillen Z, Miller E, Lanekoff I and Laskin J 2015 *Synchrotron Radiat. News* **28** 10–4

[28] Zwart P *et al* 2015 *Synchrotron Radiat. News* **28** 22–7

[29] Johnson I 2015 *Synchrotron Radiat. News* **28** 28–9

[30] Bicarregui J, Matthews B and Schluenzen F 2015 *Synchrotron Radiat. News* **28** 30–5

[31] Boehnlein A, Matthews B, Proffen T and Schluenzen F 2015 *Synchrotron Radiat. News* **28** 43–7

[32] Gehrke R, Kopmann A, Wintersberger E and Beckmann F 2015 *Synchrotron Radiat. News* **28** 36–42

[33] Jiang M, Yang X, Xu H, Zhao Z and Ding H 2009 *Chin. Sci. Bull.* **54** 4171

[34] Yin L, Tai R, Wang D and Zhao Z 2016 *J. Vac. Soc. Japan* **59** 198–204

[35] Tian F *et al* 2015 *Nucl. Sci. Tech.* **26** 030101

[36] Qi-Sheng W, Feng Y, Sheng H, Bo S, Kun-Hao Z, Ke L, Zhi-Jun W, Chun-yan X, Si-Sheng W and Li-Feng Y 2015 *Nucl. Sci. Tech.* **26** 12–7

[37] Xie H *et al* 2015 *Nucl. Sci. Tech.* **26** 020102

[38] Yang S, Wang L, Zhao J, Xue C, Liu H, Xu Z, Wu Y and Tai R 2015 *Nucl. Sci. Tech.* **26** 010101

[39] Yang T *et al* 2015 *Nucl. Sci. Tech.* **26** 020101

[40] Yu H *et al* 2015 *Nucl. Sci. Tech.* **26** 050102

[41] Zhang L *et al* 2015 *Nucl. Sci. Tech.* **26** 060101

[42] Zhang L *et al* 2015 *Nucl. Sci. Tech.* **26** 040101

[43] Xue C *et al* 2010 *Rev. Sci. Instrum.* **81** 103502

[44] Li N, Li X, Wang Y, Liu G, Zhou P, Wu H, Hong C, Bian F and Zhang R 2016 *J. Appl. Crystallogr.* **49** 1428–32

[45] Liu G, Li Y, Wu H, Wu X, Xu X, Wang W, Zhang R and Li N 2018 *J. Appl. Crystallogr.* **51** 1633–40

[46] Zhou X-J, Zhu H-C, Zhong J-J, Peng W-W, Ji T, Lin Y-C, Tang Y-Z and Chen M 2019 *Nucl. Sci. Tech.* **30** 182

[47] Tai R SSRF Phase-II Beamline Project: status and progress (available at: http://sri2018.nsrrc.org.tw/site/userdata/1157/paper/C1.2-0756.pdf)

[48] Sun B, Wang Y, Liu K, Wang Q and He J 2019 *AIP Conf. Proc.* **2054** 060028

[49] Qisheng W, Chunyan X, Kunhao Z, Liu K and Jianhua H 2019 *AIP Conf. Proc.* **2054** 060033

[50] Chen Z-H, Sun F-F, Zou Y, Song F, Zhang S, Jiang Z, Wang Y and Tai R-Z 2018 *Nucl. Sci. Tech.* **29** 26

[51] Li A, Jiang H, Wang H, Zhang Z, He Y, Zhao G and Shu D 2017 *Proc.SPIE* (https://doi.org/10.1117/12.2273518)

[52] Deng B, Ren Y, Wang Y, Du G, Xie H and Xiao T 2013 Full-field x-ray nano-imaging at SSRF *SPIE* (https://doi.org/10.1117/12.2035589)

[53] Zhongmin X, Limin J, Xiangjun W, Yajun T and Wei L 2017 *Proc.SPIE* **10389**

[54] Li Z, Fan Y, Xue L, Zhang Z and Wang J 2019 *AIP Conf. Proc.* **2054** 060040

[55] Zhang Q-L, Tian S-Q, Jiang B-C, Xu J-P and Zhao Z-T 2016 *Chin. Phys.* C **40** 037001

[56] Tian S Q, Zhang M Z, Zhang Q L, Jiang B C and Zhao Z T 2015 *6th Int. Particle Accelerator Conf., IPAC2015 (Richmond, VA, USA)* (https://doi.org/10.18429/JACoW-IPAC2015-MOPJE009)

[57] Zhao Z T, Yin L X, Leng Y B, Jiang B C, Tian S Q and Zhang M Z 2015 *6th Int. Particle Accelerator Conf., IPAC2015 (Richmond, VA, USA)* (https://doi.org/10.18429/JACoW-IPAC2015-TUPJE023)

[58] Xi S, Borgna L S and Du Y 2015 *J. Synchrotron Radiat.* **22** 661–5

[59] Chen G, Chu S, Sun T, Sun X, Zheng L, An P, Zhu J, Wu S, Du Y and Zhang J 2017 *J. Synchrotron Radiat.* **24** 1000–5

[60] Crankshaw D, Bailis P, Gonzalez J E, Li H, Zhang Z, Franklin M J, Ghodsi A and Jordan M I 2014 The missing piece in complex analytics: low latency, scalable model management and serving with velox (arXiv:1409.3809) (Accessed 1 September 2014)

[61] Miao H, Li A, Davis L S and Deshpande A 2017 ModelHub: towards unified data and lifecycle management for deep learning *2017 IEEE 33rd Int. Conf. on Data Engineering (ICDE)* pp 571–82 (arXiv:1611.06224)

[62] Deslippe J, Essiari A, Patton S J, Samak T, Tull C E, Hexemer A, Kumar D, Parkinson D and Stewart P 2014 *Proc. 9th Workshop on Workflows in Support of Large-Scale Science (New Orleans, LA, USA)* pp 31–40

[63] Venkatakrishnan S, Mohan K A, Beattie K, Correa J, Dart E, Deslippe J R, Hexemer A, Krishnan H, MacDowell A A and Marchesini S 2016 *Electron. Imaging* **2016** 1–7

[64] Bicer T, Gursoy D, Kettimuthu R, Foster I T, Ren B, Andrede V D and Carlo F D 2017 *2017 IEEE 13th Int. Conf. on e-Science (e-Science) (Auckland, New Zealand)* pp 59–68

[65] Blaiszik B, Chard K, Chard R, Foster I and Ward L 2019 *AIP Conf. Proc.* **2054** 020003

[66] Blair J, *et al* 2014 *SPIE optical engineering + applications* p 9

[67] von Laszeski G *et al* 2000 *Real-time Analysis, Visualization, and Steering of Microtomography Experiments at Photon Sources* (IL, US: Argonne National Lab)

[68] Chard K, Tuecke S and Foster I 2014 *IEEE Cloud Comput.* **1** 46–55

[69] Zhao Z T, Yin L X, Leng Y B, Zhang W Z, Jiang B C and Tian S Q 2013 *IPAC2013* (*Shanghai, China*) (available at: https://accelconf.web.cern.ch/IPAC2013/papers/mopea045.pdf)

[70] He J and Zhao Z 2014 *Natl Sci. Rev.* **1** 171–2

[71] Zhang L, Zhao M, Zeng G and Zhang L 2019 *Manage. Rev.* **31** 279–88

[72] Wang B, Guan Z, Yao S, Qin H, Nguyen M H, Yager K and Yu D 2016 *2016 New York Scientific Data Summit (NYSDS) (New York, NY, USA)* pp 1–5

[73] Liu Z, Bicer T, Kettimuthu R and Foster I 2019 Deep learning accelerated light source experiments (arXiv:1910.04081) (Accessed 1 October 2019)

[74] Zheng L F, Liu P, Zhang Z H, Hu C, Mi Q R, Wu Y F, Gong P R, Zhu Z X and Li Z 2010 *AIP Conf. Proc.* **1234** 805–8

[75] Zhao Y, Zhou Y, Hu C, Zhang X and Zhang Z 2018 *2018 5th Int. Conf. on Systems and Informatics (ICSAI) (Nanjing, China)* pp 134–8

[76] Toby B H *et al* 2009 *J. Appl. Crystallogr.* **42** 990–3

[77] Wang Q, Sun B, Zhou H, Wang Z, Yu F and He J 2019 *Nucl. Instrum. Methods Phys. Res.* A **914** 42–5

[78] Chard R, Chard K, Alt J, Parkinson D Y, Tuecke S and Foster I 2017 *2017 IEEE 37th Int. Conf. On Distributed Computing Systems Workshops (ICDCSW) (Atlanta, GA, USA)* pp 389–94

[79] Paul A K, Tuecke S, Chard R, Butt A R, Chard K and Foster I T 2017 *2nd Joint Int. Workshop on Parallel Data Storage & Data Intensive Scalable Computing Systems* (*Denver, CO, USA*) pp 49–54

[80] Ramakrishnan L and Canon R 2013 *4th SC Workshop on Petascale (Big) Data Analytics: Challenges and Opportunities* (available at: http://sc13.supercomputing.org/sites/default/files/WorkshopsArchive/pdfs/wp163s1.pdf)

[81] Wozniak J M, Chard K, Blaiszik B, Osborn R, Wilde M and Foster I 2015 *2015 IEEE/ACM 2nd Int. Symp. on Big Data Computing (BDC) (Limassol, Cyprus)* pp 51–60

[82] Flannery D *et al* 2009 *2009 5h IEEE Int. Conf. on e-Science (Oxford, UK)* pp 201–7

[83] DLS ICAT (available at: https://icat.diamond.ac.uk/)

[84] ISIS ICAT (available at: https://data.isis.stfc.ac.uk/)

[85] ILL ICAT (available at: https://explorer.ill.eu/)

[86] Tang M, Zhang J, Li Y, Du R, Yan L, Wang Z, Qi F and Tian H Data management and user data portal at CSNS (available at: https://indico.cern.ch/event/773049/contributions/3474469/)

[87] SciCat project—data acquisition, management and publication (available at: https://scicatproject.github.io/)

[88] Krahl R Using ICAT for research data management at HZB (available at: https://th.fhi-berlin.mpg.de/meetings/meta2019/uploads/Meeting/Krahl.pdf)

[89] HZB ICAT (available at: https://icatproject.org/collaboration/facilities/hzb/)

[90] Fisher S M, Phipps K and Rolfe D J 2013 *Proc. 5th Int. Workshop on Science Gateways for Life Sciences, Ser. IWSG 2013* (available at: http://ceur-ws.org/Vol-993/paper6.pdf)

[91] PaNdata ICAT (available at: http://pan-data.eu/ICAT)

[92] Shoaib S and Brian M 2004 *CCLRC Scientific Metadata Model: Version 2* (available at: https://epubs.stfc.ac.uk/manifestation/485/csmdm.version-2.pdf)

[93] Yang E, Matthews B and Wilson M 2013 *Future Gener. Comput. Syst.* **29** 612–23

[94] Matthews B, Sufi S, Flannery D, Lerusse L, Griffin T, Gleaves M and Kleese K 2010 *Int. J. Digit. Curation* **5** 106–18

[95] Pandolfi R J *et al* 2018 *J. Synchrotron Radiat.* **25** 1261–70

[96] UmbrellaID (available at: https://umbrellaid.org/)

[97] Gallagher-Jones M *et al* 2014 *Nat. Commun.* **5** 3798

[98] Nam K-W, Bak S-M, Hu E, Yu X, Zhou Y, Wang X, Wu L, Zhu Y, Chung K-Y and Yang X-Q 2013 *Adv. Funct. Mater.* **23** 1047–63

[99] Grzechnik A, Meven M, Paulmann C and Friese K 2020 *J. Appl. Crystallogr.* **53** 9–14

[100] Grolimund D *et al* 2011 *J. Anal. At. Spectrom.* **26** 1012–23

[101] Fahrnbauer F, Rosenthal T, Schmutzler T, Wagner G, Vaughan G B M, Wright J P and Oeckler O 2015 *Angew. Chem. Int. Ed.* **54** 10020–3

[102] Attwood D 2007 *Soft X-Rays and Extreme Ultraviolet Radiation: Principles and Applications* (Cambridge: Cambridge University Press)

[103] Wang Q-S *et al* 2018 *Nucl. Sci. Tech.* **29** 68

[104] Yu F *et al* 2019 *J. Appl. Crystallogr.* **52** 472–7

[105] Dalesio L R, Hill J O, Kraimer M, Lewis S, Murray D, Hunt S, Watson W, Clausen M and Dalesio J 1994 *Nucl. Instrum. Methods Phys. Res.* A **352** 179–84

[106] The Experimental Physics and Industrial Control System (available at: https://epics-controls.org/)

[107] Russell S and Norvig P 2010 *Artificial Intelligence: A Modern Approach* (Upper Saddle River, NJ: Prentice Hall)

[108] Sutton R S and Barto A G 2018 *Reinforcement Learning: An Introduction* (Cambridge MA: MIT Press)

[109] Estivill-Castro V 2002 *SIGKDD Explor. Newsl.* **4** 65–75

[110] Lin J, Keogh E, Lonardi S and Chiu B 2003 *Proc. 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery* (*San Diego, California*) pp 2–11

[111] Brase J 2009 *2009 4th Int. Conf. on Cooperation and Promotion of Information Resources in Science and Technology* pp 257–61

[112] DataCite (available at: https://datacite.org/)

[113] Wilkinson M D *et al* 2016 *Sci. Data* **3** 160018

[114] Vogelgesang M, Chilingaryan S, Santos T D and Kopmann A 2012 *2012 IEEE 14th Int. Conf. on High Performance Computing and Communication & 2012 IEEE 9th Int. Conf. on Embedded Software and Systems Liverpool, UK* pp 824–9

[115] Marchesini S, Krishnan H, Daurer B J, Shapiro D A, Perciano T, Sethian J A and Maia F R N C 2016 *J. Appl. Crystallogr.* **49** 1245–52