

# eXtreme-DataCloud project: Advanced data management services for distributed e-infrastructures

A. Costantini<sup>1</sup>, D. Cesini<sup>1</sup>, D.C. Duma<sup>1</sup>, G. Donvito<sup>2</sup>, M. Viljolen<sup>3</sup>, S. Battaglia<sup>4</sup>, V. Poireau<sup>5</sup>, L. Dell'Agnello<sup>1</sup>, O. Keeble<sup>6</sup>, R. Lemrani<sup>7</sup>, C. Ohmann<sup>4</sup>, J. Marco de Lucas<sup>8</sup>, L. Dutka<sup>9</sup>, P. Fuhrmann<sup>10</sup>, F. Aguilar Gomez<sup>11</sup>

<sup>1</sup> INFN-CNAF, Bologna, Italy

<sup>2</sup> INFN Bari, Bari, Italy

<sup>3</sup> EGI, Netherlands

<sup>4</sup> ECRIN, France

<sup>5</sup> CNRS, France

<sup>6</sup> CERN, Switzerland

<sup>7</sup> IN2P3, France

<sup>8</sup> CSIC, Spain

<sup>9</sup> AGH, Poland

<sup>10</sup> DESY, Germany

<sup>11</sup> Univ. de Cantabria, Spain

E-mail: [alessandro.costantini@cnafe.infn.it](mailto:alessandro.costantini@cnafe.infn.it)

**Abstract.** The development of new data management services able to cope with very large data resources is becoming a key challenge. Such capability, in fact, will allow the future e-infrastructures to address the needs of the next generation extreme scale scientific experiments. To face this challenge, in November 2017 the H2020 eXtreme DataCloud - XDC project has been launched. Lasting for 27 months and combining the expertise of 8 large European research organisations, the project aims at developing scalable technologies for federating storage resources and managing data in highly distributed computing environments. The targeted platforms are the current and next generation e-Infrastructures deployed in Europe, such as the European Open Science Cloud (EOSC), the European Grid Infrastructure (EGI), and the Worldwide LHC Computing Grid (WLCG). The project is use-case driven with a multidisciplinary approach, addressing requirements from research communities belonging to a wide range of scientific domains: High Energy Physics, Astronomy, Photon and Life Science, Medical research. XDC will implement data management scalable services, combining already established data management and orchestration tools, to address the following high level topics policy driven data management based on Quality-of-Service, Data Life-cycle management, smart placement of data with caching mechanisms to reduce access latency, meta-data with no predefined schema handling, execution of pre-processing applications during ingestion, data management and protection of sensitive data in distributed e-infrastructures, intelligent data placement based on access patterns. This contribution will introduce the project, presents the foreseen overall architecture and the developments that are being carried on to implement the requested functionalities.

## 1. Introduction

Lead by INFN-CNAF, the eXtreme DataCloud (XDC) project [1] develops scalable technologies for federating storage resources and managing data in highly distributed computing environments. The provided services are capable of operating at the unprecedented scale required by the most demanding, data intensive, research experiments in Europe and Worldwide. The targeted platforms for the released products are the already existing and the next generation e-Infrastructures deployed in Europe, such as the European Open Science Cloud (EOSC) [2], the European Grid Infrastructure (EGI) [3], the Worldwide LHC Computing Grid (WLCG) [4] and the computing infrastructures that will be funded by the upcoming H2020 EINFRA-12 call. XDC is funded by the H2020 EINFRA-21-2017 Research and Innovation action under the topic Platform-driven e-Infrastructure innovation [5]. It is carried on by a Consortium that brings together technology providers with a proven long-standing experience in software development and large research communities belonging to diverse disciplines: Life Science, Biodiversity, Clinical Research, Astrophysics, High Energy Physics and Photon Science. XDC started on 1st November 2017 and will run for 27 months until January 2020. The EU contribution for the project is 3.07 million euros. XDC is a use case driven development project and the Consortium has been built as a combination of technology providers, Research Communities and Infrastructure providers. New developments will be tested against real-life applications and use cases. Among the high level requirements collected from the Research Communities, the Consortium identified those considered more general (and hence exploitable by other communities), with the greatest impact on the user base and that can be implemented in a timespan compatible with the project duration and funding.

## 2. Project Objectives

The XDC project develops open, production quality, interoperable and manageable software that can be easily plugged into the target European e-Infrastructures and adopts state of the art standards in order to ensure interoperability. The building blocks of the high-level architecture foreseen by the project are organized in a manner to avoid duplication of development effort. All the interfaces and links to implement the XDC architecture are developed exploiting the most advanced techniques for authorization and authentication. Services are scalable to cope with most demanding, extreme scale scientific experiments like those run at the Large Hadron Collider at CERN and the Cherenkov Telescope Array (CTA), both of them represented in the project consortium. The project will enrich already existing data management services by adding missing functionalities as requested by the user communities. The project will continue the effort invested by the now ended INDIGO-DataCloud project [6] in the direction of providing friendly, web-based user interfaces and mobile access to the infrastructure data management services. The project will build on the INDIGO-DataCloud achievements in the field of Quality of Services and data lifecycle management developing smart orchestration tools to realize easily an effective policy driven data management. One of the main objectives of the project is to provide data management solutions for the following use cases:

- Dynamic extension of a computing center to a remote site providing transparent bidirectional access to the data stored in both locations.
- Dynamic inclusion of sites with limited storage capacity in a distributed infrastructure, proving transparent access to the data stored remotely.
- Federation of distributed storage endpoints, i.e. a so-called WLCG Data Lake, enabling fast and transparent access to their data without a-priori copy.

These use cases will be addressed implementing intelligent, automatic and hierarchical caching mechanisms.

### 3. User Communities requirements

In the following paragraphs the user communities represented in XDC are presented together with the list of the main requirements injected into the project development process.

#### 3.1. *LifeWatch*

LifeWatch [7] is the e-Science and Technology European Infrastructure for Biodiversity and Ecosystem Research, that aims to advance science in these disciplines and to address the big environmental challenges as well as to support knowledge-based strategic solutions to environmental preservation. This mission is achieved by providing access to a multitude of datasets, services, tools and computing resources in general, enabling the construction and operation of Virtual Research Environments (VREs). Requirements injected into the project:

- (i) Integration of different tools based on the cloud, that allow the management of the data life cycle, the production of data based on FAIR (+R, Reproducibility) principles,
- (ii) Integration of tools that allow the automation of saving the data produced by the sensors stored in distributed storage solution and the management of the metadata for the integration of the different sources.

#### 3.2. *CTA*

Very high-energy electromagnetic radiation reaches Earth from a large part of the Cosmos, carrying crucial and unique information about the most energetic phenomena in the Universe. CTA [8] (Cherenkov Telescope Array) will answer many of the persisting questions by enabling the detection of more than thousands sources over the whole sky. CTA builds on the proven technique of detecting gamma-ray induced particle cascades in the atmosphere through their Cherenkov radiation, simultaneously imaging each cascade stereoscopically with multiple telescopes, and reconstructing the properties of the primary gamma ray from those images. Requirements injected into the project:

- (i) Intelligent and automated dataset distribution: The CTA computing model should be distributed among four data centers and two geographic array sites. The event reconstruction consists of six levels of data, each with specific and complex policies for data distribution: disk vs tape, nb of versions, nb of replicas per version, etc. Some versions must be automatically erased/removed by the system. Data placement is important to distribute load between the four data centers, not only for storage but also for reconstruction/analysis jobs. The data access rights are complex.
- (ii) Data ingestion preprocessing: The CTA distributed archive is based on the Open Archival Information System (OAIS) ISO standard. Event data are in files (FITS format) containing all metadata. Metadata are extracted from the ingested files, with an automatic filling of the metadata database. Metadata will be used for the further query of archive. The system should be able to manage replicas, tapes, disks, etc, with data from low-level to high-level.

#### 3.3. *ECRIN*

ECRIN (European Clinical Research Infrastructure Network) [9] is a not-for-profit intergovernmental organisation that provides support for the development and implementation of multinational clinical research projects in Europe. These are mostly investigator initiated (rather than industry sponsored) clinical trials, run by non-commercial Clinical Trials Units (CTUs) based within universities or hospitals, though ECRIN does also support trials initiated by biotech and medical device small and medium enterprises (SMEs). As a result more and more researchers are making such material (generically, clinical trial data objects) available for sharing with others. The datasets are rarely freely available - instead a variety of access mechanisms

(e.g. individual request and review, membership of pre-authorised groups, or web based self-attestation) are used in combination with different access types (e.g. download versus in-situ perusal). Furthermore the various data objects are stored in a wide variety of different locations making the discoverability of such data an important challenge to address. Requirements injected into the project:

- (i) Availability of metadata managements tools, easy to use, web based that allow researchers to provide the necessary metadata themselves (and therefore most accurately), aggregating that metadata centrally
- (ii) Availability of tools to harvest and map the metadata from existing data repositories, using available APIs and/or data mining techniques
- (iii) Availability of a generic metadata schema that can not only describe the data objects themselves, but also link them to the source trials and provide information on how the data objects can be accessed (ECRIN has recently proposed such a schema, created by extending DataCite).

### 3.4. *WLCG*

The WLCG (Worldwide LHC Computing Grid) [4] is both a worldwide collaboration that counts thousands of researchers and a research distributed e-infrastructure shared by the LHC (Large Hadron Collider) experiments at CERN. WLCG is, in fact, composed by at least four main distinct communities, one per experiment: ALICE, ATLAS, CMS and LHCb. WLCG was historically built for the 4 WLCG experiments (ATLAS; CMS; LHCb, ALICE) but it is now shared, in several sites, with many other High Energy Physics (and astrophysics) Virtual Organizations, at national and international level. This created a very variegated and heterogeneous ecosystem of HEP user communities injecting requirements into the same infrastructure. Thousands of researchers exploit the infrastructure on a daily basis to run their applications workflows. Due to the advent of the High Luminosity (HI-LUMI) LHC, expected to enter into production in 2025, has been estimated that the infrastructure and the provided services will need to handle more than 1.5 ExaByte. This extreme scale poses several challenges, technical and economical, on the data management infrastructure. In fact, the storage resources counts roughly for the 50% of the total infrastructure cost. To support these new computing models, improved infrastructure-level data management services would greatly help in lowering the effort to maintain the experiment frameworks. Requirements injected into the project: The XDC project will provide smart caches mechanisms to be added to the storage systems currently deployed in the WLCG infrastructures. These caching mechanisms will be the basis to realize a distributed storage system deployed across centers connected by fat networks (10Tb in 2025) and operated as a single service. The so called "DataLakes" will simplify the dynamic extension of sites to remote locations and allowing for less datasets distribution. The inclusion into the infrastructure of non standard configuration of sites (i.e. diskless, only tape, HPC farm, etc) will be favored by the availability of such systems. The increased usage of caches will allow less disk usage and a reduced a-priori data transfers among the distributed sites.

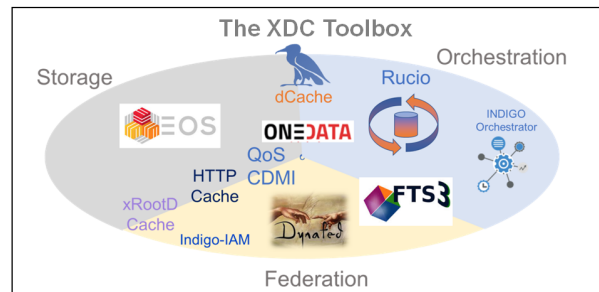
### 3.5. *XFEL*

The European XFEL [10] is a X-RAY Light Source, located in Hamburg, Germany. It is installed mainly in underground tunnels which can be accessed on three different sites. The 3.4 kilometre-long facility will run from the DESY campus in Hamburg to the town of Schenefeld in Schleswig-Holstein. At the research campus in Schenefeld, teams of scientists from all over the world will carry out experiments using the X-ray flashes. Using the X-ray flashes of the European XFEL, scientists will be able to map the atomic details of viruses, decipher the molecular composition of cells, take three-dimensional images of the nanoworld, film chemical reactions, and study

processes such as those occurring deep inside planets. The expectation is to make derived data (images) available for further processing at the Kurchatov Institute NRC as fast as possible after the raw data has been calibrated at DESY. As the European XFEL started in Sep 2017, this process is not yet in place. So the XDC work will be an integral part of setting up the data analysis chain for the collaboration with the Moscow Kurchatov Institute NRC and the XFEL central facilities

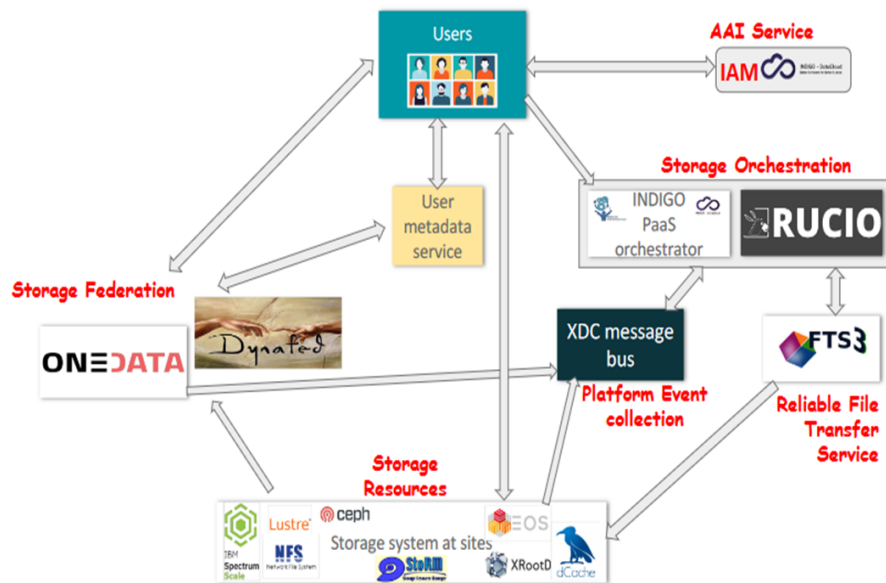
#### 4. XDC Overall architecture

The XDC project aims at providing advanced data management capabilities that require the execution of several tasks and the interaction among several components and services. Those capabilities should include but are not limited to QoS management, preprocessing at ingestion and automated data transfers. Therefore global orchestration layer is needed to take care of the execution of those complex workflows. Figure 1 highlights the main components and their role among the three different levels: Storage, Federation, and Orchestration.



**Figure 1.** XDC main components and related roles.

Figure 2 highlights the high level architecture of the XDC project by describing the components and the related connections.



**Figure 2.** High level architecture of the XDC project.

#### *4.1. XDC Orchestration system*

In the XDC project the global orchestration layer is needed to take care of the execution of those complex workflows. The orchestration covers two essential aspects:

- The overall control, steering and bookkeeping including the connection to compute resources
- The orchestration of the data management activities like data transfers, and data federation.

Consequently we have decided to split the responsibilities between two different components: the INDIGO Orchestrator [11] and Rucio [12]. The INDIGO PaaS Orchestrator, the system wide orchestration engine, is a component of the PaaS layer that allows to instantiate resources on Cloud Management Frameworks (like OpenStack and OpenNebula) and Mesos clusters. It takes the deployment requests, expressed through templates written in TOSCA YAML Profile 1.0 [13], and deploys them on the best cloud site available. The Rucio project, the data management orchestration subsystem, is the new version of ATLAS Distributed Data Management (DDM) system services for allowing the ATLAS collaboration to manage the large volumes of data, both taken by the detector as well as generated or derived, in the ATLAS distributed computing system. Rucio is used to manage accounts, files, datasets and distributed storage systems. Those two components, the PaaS Orchestrator and Rucio, provide different capabilities and can complement each other to offer a full set of features to meet the XDC requirements. Rucio implements the data management functionalities missing in the INDIGO Orchestrator: the Orchestrator will make use of those capabilities to orchestrate the data movement based on policies. Rucio will be integrated in the INDIGO Orchestrator as a plugin to be used to steer the data movement. The Orchestrator will be the main entry point for the user requests and, in particular, it will also interact directly with the storage backends in order to get information about the data availability in order to trigger the right processing flow

#### *4.2. XDC Quality-of-Service implementation*

The idea to provide scientific communities or individuals with the ability to specify a particular quality of service when storing data, e.g. the maximum access latency or minimum retention policy, was introduced within the INDIGO-DataCloud project. In XDC, the QoS concept is envisioned to consistently compliment all data related activities. In other words, whenever storage space is requested, either manually by a user or programmatically by a framework, the quality of that space can be negotiated between the requesting entity and the storage provider

#### *4.3. Caching within XDC*

In this section we consider how the XDC architecture treats the storage and access of data, building a hierarchy of components whose goal is to maximise the accessibility of data to clients while minimising global infrastructure costs. The architecture considers a set of multi-site storage systems, potentially accessed through caches, both of which are aggregated globally through a federation. While large, multi-site storage systems may hold the majority of a community custodial data, XDC does not foresee that they will necessarily host all the compute capacity. In particular, CPU-only resource centres or cloud procurements must be supported. Provision must therefore be made to ensure wide area data access is as efficient as possible. Such resource centres may access custodial data through a standalone proxy cache. The simplest cache envisaged by XDC, a kind of "minimum viable product", would have the following characteristics: Read-only operation, Fetch data on miss with service credentials, Data can be chunked or full-file, Manage cache residency, evicting data when necessary, HTTP frontend with group-level authorisation. Beyond this scenario, certain extensions will be investigated: Write support, Synchronisation and enforcement of ACLs, Namespace synchronisation (allowing discovery of non-resident data), Call-out to FTS for data transport, QoS support, Integration

with notifications, such as primary storage notifying a file deletion. To such purpose, various technologies are available to the project to serve as the basis of an implementation:

- The system runs native dCache [14] or EOS, but operates in a "caching mode" staging data in when a cache miss occurs.
- A service such as Dynafed [15] will be augmented to initiate data movement. While it would only hold only metadata, it would use a local storage system for holding the data
- A standalone HTTP cache could be built from existing web technology, such as nginx [16], modified for horizontal scalability and relevant AAI support.

#### *4.4. XDC data management and new developments*

Data management functionality for end users will be also available via the Onezone component of the Onedata platform [17]. Onezone will provide single-sign on authentication and authorization for users, which will be able to create access tokens to perform data access activities via the Web browser, REST API or using Onedata POSIX virtual filesystem. Onezone will enable federating multiple storage sites by deploying Oneprovider services on top of the actual storage resources provisioned by the sites. For the purpose of job scheduling and orchestration, Onedata will communicate with Indigo Orchestrator component by means of a message bus, allowing the orchestrator to subscribe for events related to data transfers and data access. This will allow Orchestrator to react to changes in the overall system state (e.g. a new file in a specific directory or space, data distribution changes initiated by manual transfers, cache invalidation or on-the-fly block transfers). In order to ensure federation level Quality of Service Onedata will expose an interface providing information on data access latency and block based file location. Onedata will be also responsible for definition of federation level authentication and authorization aspects of data access, based on OpenID Connect [18]. Upon login, each user will be able to generate access tokens, which can be used for accessing and managing data using REST API or mounting the virtual POSIX filesystems on computational nodes. Furthermore, Onedata will provide a federated level encryption key management service, allowing users to securely upload symmetric encryption keys (e.g. AES-256). On the data access layer, Onedata will provide a WebDAV [19] storage interface, to enable integration of other HTTP transfer based components such as FTS [20] or EOS [21] to make the data managed by these components to be accessible in a unified manner via the POSIX virtual filesystem provided by Onedata. Furthermore, Onezone, the entry point to the data management aspects of the platform, will allow for semi-automated creation of data discovery portals, based on metadata stored in the federated Oneprovider instances and on a centralized ElasticSearch engine indexing the metadata. This solution will allow the communities to create custom indexes on the data and metadata, provide customizable styles and icons for their users and defining custom authorization rights based on user classes (public access, access on login, group access, etc.).

## **5. Conclusions**

In the present contribution the XDC objectives, starting from the technology gaps that currently prevent effective exploitation of distributed computing and storage resources by many scientific communities, have been discussed and presented. Those objectives are the real driver of the project and derive directly from use cases, and the related needs, presented by the scientific communities involved in the project itself, covering areas such as Physics, Astrophysics, Bioinformatics, and others. Starting from the above assumptions, the overall structure of the project have been presented by emphasizing its components, typically based upon or extend established open source solutions, and the relations among them

Moreover, XDC project can complement and integrate with other running projects and communities and with existing multi-national, multi-community infrastructures. As an example,

XDC is collaborating with the Designing and Enabling E-Infrastructures for intensive Processing in Hybrid Data Clouds (DEEP-Hybrid-DataCloud) [22] project aimed at promoting the integration of specialized, and expensive, hardware under a Hybrid Cloud platform, and targeting the evolution of the corresponding Cloud services supporting these intensive computing techniques to production level.

As an added value both projects (XDC and DEEP-HDC) have the common objective to open new possibilities to scientific research communities in Europe by supporting the evolution of e-Infrastructure services for exascale computing. Those services are expected to become a reliable part of the final solutions for the research communities available in the European Open Science Cloud Service Catalogue.

## 6. References

- [1] Web site: [www.extreme-datacloud.eu](http://www.extreme-datacloud.eu)
- [2] Web site: <https://ec.europa.eu/research/openscience/index.cfm?pg=open-science-cloud>
- [3] Web site: <https://www.egi.eu/>
- [4] Web site: [wlcg.web.cern.ch/](http://wlcg.web.cern.ch/)
- [5] Web site: <http://ec.europa.eu/research/participants/portal/desktop/en/opportunities/h2020/topics/einfra-21-2017.html>
- [6] Web site: <https://www.indigo-datacloud.eu>
- [7] Web site: <https://www.lifewatch.eu>
- [8] Web site: <https://www.cta-observatory.org>
- [9] Web site: <https://www.ecriin.org>
- [10] Web site: <https://www.xfel.eu>
- [11] Web site: [www.indigo-datacloud.eu/paas-orchestrator](http://www.indigo-datacloud.eu/paas-orchestrator)
- [12] Web site: <https://rucio.cern.ch/>
- [13] TOSCA Simple Profile in YAML Version 1.0. Edited by Derek Palma, Matt Rutkowski, and Thomas Spatzier. 27 August 2015. OASIS Committee Specification Draft 04 / Public Review Draft 01
- [14] Web site: [www.dcache.org](http://www.dcache.org)
- [15] Web site: [lcgdm.web.cern.ch/dynafed-dynamic-federation-project](http://lcgdm.web.cern.ch/dynafed-dynamic-federation-project)
- [16] Web site: <https://www.nginx.com/>
- [17] Web site: [ondata.org](http://ondata.org)
- [18] Web site: <https://openid.net/connect/>
- [19] Web site: [www.webdav.org/](http://www.webdav.org/)
- [20] Web site: [information-technology.web.cern.ch/services/file-transfer](http://information-technology.web.cern.ch/services/file-transfer)
- [21] Web site: [eos.web.cern.ch](http://eos.web.cern.ch)
- [22] Web site: <https://deep-hybrid-datacloud.eu/>