

# Dataverse and digital preservation: an overview of different institutions

(Doc. CO23/16) (4 RDM\Repositori de dades\Preservacio\SEIDOR\2306 DataverseDigitalPreservation\_anOverview.docx, 16.06.23)

Report prepared by SEIDOR

for Open Science Area  
Consorci de Serveis Universitaris de Catalunya (CSUC)

June 2023

## Acknowledgements

---

This document has been prepared by SEIDOR jointly with representatives of the Consorci de Serveis Universitaris de Catalunya (CSUC).

We would like to express our sincere appreciation to all experts cited in this document for the time they devoted to informing and sharing their expertise.

Written by Albert Baucells and David Pérez (SEIDOR)

Coordinated by Mireia Alcalá and Lluís Anglada (CSUC - Open Science Area)

Digital version: <http://hdl.handle.net/2072/536532>

*This project is part of the Next Generation Recovery Plan, as a response by the Union to the COVID-19 pandemic under the Operational Programme FEDER of Catalonia 2014-2020 with a grant of €12,73 million.*

Version control

---

<b>Data i versió</b>	<b>Autors</b>
16-05-23 DRAFT	<ul style="list-style-type: none"><li>• Albert Baucells and David Pérez</li></ul>
26-05-23 DRAFT-V1	<ul style="list-style-type: none"><li>• Mireia Alcalá and Lluís Anglada</li></ul>
31-05-23 DRAFT-V3	<ul style="list-style-type: none"><li>• Albert Baucells and David Pérez</li></ul>
13-06-23 DRAFT-V3	<ul style="list-style-type: none"><li>• Mireia Alcalá and Lluís Anglada</li></ul>
14-06-23 V1	<ul style="list-style-type: none"><li>• David Pérez</li></ul>

## Summary

---

<b>Executive summary .....</b>	<b>5</b>
<b>1. Introduction.....</b>	<b>6</b>
<b>2. Objective .....</b>	<b>6</b>
<b>3. Methodology .....</b>	<b>7</b>
<b>4. Analysis of preservation systems.....</b>	<b>7</b>
4.1. dataRepositoriUM	7
4.2. Borealis	8
4.3. DataverseNL	8
4.4. DataverseNO	9
4.5. FAIRData	9
4.6. ODUM	9
<b>5. Conclusions .....</b>	<b>10</b>
<b>6. Bibliography and references .....</b>	<b>12</b>
6.1. CoreTrustSeal certifications documentation	12
6.2. General information	12
6.3. Tools	13
6.4. Others	13
<b>Appendix .....</b>	<b>14</b>
1. Final questions	14
2. Institutions interviews	15
3. Grid	25
4. Recording and transcription of the interviews	26

## Executive summary

---

This document was delivered to CSUC by a SEIDOR IT team with expertise and seniority, with the aim of providing support in the acquisition of useful information for a decision-making process to be carried out in relation to the improvement of a Digital Preservation System. The document presents, in a structured way, the result of a process of data collection from a selection of five internationally recognized institutions. The selection was a choice of six institutions from different countries which are providing digital preservation services mainly to preserve open research data.

Guided by FAIR principles, as well as by main technology trends, the most valuable information required to fulfil the goals of CSUC were sought. A comparison grid was created by accessing, collecting, processing, and documenting open documentation published by target institutions. Open documentation related with the Implementation of the CoreTrustSeal certification and Core Trustworthy Data Repositories Requirements was found to be very valuable and was largely collected. The comparison grid is organized by the functional areas defined by NSDA to determine Levels of Digital Preservation: Storage, Integrity, Control, Metadata and Content. Comparison grid was approved as a helpful tool to CSUC for the decision-making process.

A more detailed analysis of CSUC's overall needs to improve its preservation system dictates that a set of questions be asked in person to team members from each institution whose preservation system is of interest to CSUC. The final part of the project is the collection of detailed information that was not being perceived through the open documentation reviews.

## 1. Introduction

---

CORA.RDR (<https://dataverse.csuc.cat>) is a federated and multidisciplinary repository for the publication of research datasets in FAIR mode (Findable, Accessible, Interoperable and Reusable) following the European Open Science Cloud (EOSC) guidelines.

These guidelines request datasets to be reusable for at least 10 years after publication. Leaving aside other aspects, this requires the digital preservation of datasets. This is based on different processes including having a remote copy system. This problem is common to the different repositories of research data and, specifically, to those who use Dataverse software (the same as CORA.RDR).

To be able to establish a system of remote copies of the datasets of the CORA.RDR and in order to do it in the most efficient way possible, we want to know the characteristics of the systems that perform this function in repositories, not only similar in objectives and function but, in addition, in software, and that is why CSUC want to have knowledge of the digital preservation characteristics of these repositories:

- Borealis (Canada)
- dataRepositoriUM (Portugal)
- DataverseNL (The Netherlands)
- DataverseNO (Norway)
- FAIRData (Finland)
- Odum's Research Data Information Systems (RDIS) (EUA)

From the repositories mentioned above, we want to know the technical operating characteristics, both software and hardware, suitable for replicating the infrastructure in the Catalan case. In an indicative way, we want to know: the technological architecture used (local, in the cloud, etc.), the mechanisms of the ingest process (centralized/decentralized), the relationship between the preservation function and the open repository, etc.

## 2. Objective

---

The report aims to explore how different institutions are utilizing Dataverse and implementing digital preservation strategies. For this reason, is to outline the state of the art of the reference technological models for the digital preservation of digital objects or "Datasets" (set of files and metadata) of the research system.

### 3. Methodology

---

This section provides information on the methodology followed in this report. First, the institution selection was provided by CSUC and included institutions from Europe (dataRepositoriUM from Portugal, DataverseNL from the Netherlands, DataverseNO from Norwegian and DigitalPreservationFI from Finland) and America (Borealis from Canada and Odum's Research Data Information Systems (RDIS) from EUA) institutions.

Based on these systems, information is collected from the different public websites (see 6. Bibliography and references) and a series of questions are drawn up that should be answered. After a discussion by CSUC, some of these questions are discarded and the nine final questions (see appendix 1) asked information for: number of copies, localisation, media type, integrity and checking.

With the information collected, it is checked that there are still gaps and each institution is requested to have a short one-hour interview (see appendix 2) to go deeper with some of the topics. In general terms, these interviews are done with the technical staff of the repository and allow for fill-in the gaps.

Then, a grid (see appendix 3) is developed with different indicators to easily compare the different systems. And this allows us to make the synthesis and final conclusions (see 5. Conclusions).

### 4. Analysis of preservation systems

---

In this section each system is introduced and a short summary of each of them is made.

#### 4.1. dataRepositoriUM

“DataRepositóriUM is an institutional Data Repository to share, publish and manage research data generated and collected by the activity of researchers and in the research units of the University of Minho.”

The UMinho considers that digital preservation is important, but they have not yet been able to develop policies and workflows in the repository. For this reason, no information is collected from this system.

## 4.2. Borealis

“Borealis, the Canadian Dataverse Repository, is a bilingual, multidisciplinary, secure, Canadian research data repository, supported by academic libraries and research institutions across Canada. Borealis supports open discovery, management, sharing, and preservation of Canadian research data.”<sup>1</sup>

The preservation system achieves three copies to five different nodes via OLRC (using a software feature of OpenStack Swift) without controlling it. The production data in OpenStack swift is automatically stored in 5 nodes, one copy in the NFS storage system synced nightly, one copy on-site on tape synced nightly, and one copy of tape off-site at Iron Mountain storage facility in Hamilton (shipped every other day or twice a week). A geographic radius of 450 kilometres away from each other copy is reached. A monthly integrity check of all the files in the Dataverse is performed. The integrity check is based on the MD5 hash formula taken from the Dataverse ingest flow. Integrity-calculated information is stored separately in a database.

## 4.3. DataverseNL

“DataverseNL is a shared service provided by participating institutes and DANS. DataverseNL uses the Dataverse software developed by Harvard University, which is used worldwide. DataverseNL is jointly offered by participating institutes and DANS. Since 2014, DANS has been managing the technical infrastructure; the participating institutes are responsible for managing the deposited data in their dataverses within DataverseNL. Every institute participates in the Dataverse Advisory Board, which determines the policies of the service. The local dataverse managers exchange their experiences regularly, together they develop "good practices" in the field of research data management.”<sup>2</sup>

The preservation system consists of two backup replicas, which are stored in geographically distributed locations within a 20km range. The storage type for the copies is a snapshot. UNF checksum and normalization function and SHA256 integrity function are used. One copy resides in the datacenter of the contracted vendor, located in Almere, another at a local CPD in The Hague. Currently, these copies are run offline (tape) and on disk. Once a month, a bit rot check is done. A great number of random files at a time are checked to ensure data integrity is conserved, this contributes to having an insight into the health of the data. No measurements are made of the total amounts of computing resources spent. The main resource spent is Disk IO. Virus checking is done by using <https://www.clamav.net/>. Integrity information is stored on disk as a file. It is physically stored on the same medium, in a different folder.

---

<sup>1</sup> <https://borealisdata.ca/about/>

<sup>2</sup> <https://dataverse.nl/>



#### 4.4. DataverseNO

“DataverseNO (<https://dataverse.no/>) is a national, generic repository for open research data, owned and operated by UiT The Arctic University of Norway. DataverseNO is aligned with the FAIR Guiding Principles for scientific data management and stewardship. The technical infrastructure of the repository is based on the open-source application Dataverse, which is developed by an international developer and user community led by Harvard University. DataverseNO is CoreTrustSeal certified”<sup>3</sup>.

The preservation system of DataverseNO has three copies of the backup replicas. There are two location copies stored in two separate datacenters within the same building in Oslo, Norway. The two datacenters are separated with a fire resistant wall. In addition, there is one copy stored in a Microsoft Azure Datacenter in The Netherlands. Currently implementing storage of another additional copy on immutable storage.

MD5 integrity functions are used.

#### 4.5. FAIRData

“Fairdata services are part of the digital preservation services of the Ministry of Education and Culture, Finland (“Minedu”). The Fairdata services support the research process and management of digital data with some of the components described in the Framework for Open Science and Research. The Fairdata services consist of the following service components: IDA – Research Data Storage, Etsin – Research Data Finder, Qvain – Research Dataset Metadata Tool, auxiliary services, such as Metax, identity management, download component for published data and Digital preservation service for Research Data (including management and packaging).”<sup>4</sup>

The preservation system consists of four total backup replicas: three online copies and one offline. The copies are geospatially distributed into three separate locations (data centers in Finland): the first one is local in disk storage, second and third on tape storage and the fourth and fifth on the dark archive. Related to integrity, the system signed SMIME manifests.

#### 4.6. ODUM

“The Odum Institute Data Archive is a leader in research data stewardship, with over 50 years of experience beginning with the acquisition of the Louis Harris Data Center in 1965. Our longstanding commitment to data access and research transparency has been a driving force behind ongoing efforts to enhance our infrastructure, workflows, and policies to ensure that the

---

<sup>3</sup> <https://info.dataverse.no/>

<sup>4</sup> <https://www.fairdata.fi/en/about-fairdata/fairdata-services/>

data assets in our care remain FAIR—findable, accessible, interoperable, and reusable—now and into the future.

We are home to one of the largest catalogues of social science research data in the U.S. which includes the Harris Polls, North Carolina Vital Statistics, and the most complete collection of 1970s U.S. Census data. In addition, we manage and provide access to the UNC Dataverse, a web-based data repository, that enables scientists, research teams, scholarly journals, and other members of the UNC research community to archive and share their own datasets.”<sup>5</sup>

The preservation system consists of four total amount backup replicas. The copies are housed in geographically distributed storage locations (offsite local server, cloud server in Northern Virginia and Northern California and a copy in the private cloud). MD5 integrity functions are used.

## 5. Conclusions

---

In conclusion, several key findings can be drawn:

### Storage:

- European institutions rely on two to five copies for preservation systems, while Canadian and American institutions use a higher **number of copies**, ranging from five to twelve.
- Using cloud storage allows for a larger number of copies compared to on-premises infrastructure.
- **Key strategies for safeguarding copies** include storing them in different locations, utilizing services from multiple providers, and using different storage media types. Some institutions heavily rely on cloud providers and their certifications, considering the type of media used for storage less important.
- There is a trend to **shift from file/block storage to object storage**, especially in cloud-based digital preservation procedures. The percentage of object storage increases with the adoption of cloud-based preservation. However, this transformation requires software changes and poses some difficulties.

### Integrity:

- Many interviewees believe that integrity checks should be activated in more situations. Dataverse ensure **integrity calculation only at ingestion time**.
- Integrity information, added to content for long-term preservation, is usually **calculated from the entire content** of the ingested file rather than from split parts. Computed integrity information is commonly stored separately from the preserved content. Some

---

<sup>5</sup> <https://odum.unc.edu/archive/#archive1>

systems comply with this requirement, and those relying on cloud providers trust their certifications.

- The **MD5 cryptographic** hash function, previously widely used for integrity validation, is considered too weak. Teams are replacing it with stronger alternatives such as SHA256 checksum verification.
- Teams relying on cloud provider infrastructure are less concerned about knowing when integrity checking is performed compared to those relying on on-premise solutions.

Control:

- **Virus checking** is not considered a top priority by most digital preservation IT teams. Open-source tools for virus checking are challenging to keep up to date. Institutions working with large internet service providers include virus checking in their ingestion processes without significant workload. Virus checking is considered unimportant for executable files rarely included in datasets.

Other:

- The analysis should have included a broader range of Digital Preservation System components, not just limited to Dataverse.

## 6. Bibliography and references

---

### 6.1. CoreTrustSeal certifications documentation

<b>DataverseNL</b>	<a href="https://www.coretrustseal.org/wp-content/uploads/2018/04/Tilburg-University-Dataverse.pdf">https://www.coretrustseal.org/wp-content/uploads/2018/04/Tilburg-University-Dataverse.pdf</a>
<b>DataverseNO</b>	<a href="https://www.coretrustseal.org/wp-content/uploads/2020/03/DataverseNO.pdf">https://www.coretrustseal.org/wp-content/uploads/2020/03/DataverseNO.pdf</a>
<b>ODUM</b>	<a href="https://www.coretrustseal.org/wp-content/uploads/2020/10/Odum-Institute-Data-Archive.pdf">https://www.coretrustseal.org/wp-content/uploads/2020/10/Odum-Institute-Data-Archive.pdf</a>

### 6.2. General information

<b>Borealis</b>	<a href="https://spotdocs.scholarsportal.info/plugins/servlet/mobile?contentId=228196817#content/view/228196817">https://spotdocs.scholarsportal.info/plugins/servlet/mobile?contentId=228196817#content/view/228196817</a>
	<a href="https://spotdocs.scholarsportal.info/display/DAT/B.+Successful+CTS+Certification+Applications">https://spotdocs.scholarsportal.info/display/DAT/B.+Successful+CTS+Certification+Applications</a>
<b>DataverseNL</b>	<a href="https://dans.knaw.nl/en/preservationplan/1">https://dans.knaw.nl/en/preservationplan/1</a>
<b>DataverseNO</b>	<a href="https://site.uit.no/dataverseno/about/policy-framework/preservation-policy/">https://site.uit.no/dataverseno/about/policy-framework/preservation-policy/</a>
	<a href="https://site.uit.no/dataverseno/about/policy-framework/preservation-policy/preservation-plan/">https://site.uit.no/dataverseno/about/policy-framework/preservation-policy/preservation-plan/</a>
	<a href="https://site.uit.no/dataverseno/">https://site.uit.no/dataverseno/</a>
<b>FAIRDATA - CSC IT Center for Science (Finland)</b>	Digital Preservation with added resilience in Finland
	<a href="https://www.fairdata.fi/en/">https://www.fairdata.fi/en/</a>
	<a href="https://digitalpreservation.fi/mets-validator">https://digitalpreservation.fi/mets-validator</a>
<b>ODUM</b>	<a href="https://www.coretrustseal.org/wp-content/uploads/2020/10/Odum-Institute-Data-Archive.pdf">https://www.coretrustseal.org/wp-content/uploads/2020/10/Odum-Institute-Data-Archive.pdf</a>

<b>FAIRDATA - CSC IT Center for Science (Finland)</b>	Digital Preservation with added resilience in Finland
	<a href="https://www.fairdata.fi/en/">https://www.fairdata.fi/en/</a>
	<a href="https://digitalpreservation.fi/mets-validator">https://digitalpreservation.fi/mets-validator</a>
<b>UMinho</b>	<a href="https://www.um.edu.mt/library/oar/bitstream/123456789/92709/1/University%20of%20Minho%20Data%20Repository.pdf">https://www.um.edu.mt/library/oar/bitstream/123456789/92709/1/University%20of%20Minho%20Data%20Repository.pdf</a>

### 6.3. Tools

<b>Libnova</b>	<a href="https://www.libnova.com/">https://www.libnova.com/</a>
<b>TwoRavens Dataverse Project</b>	<a href="https://guides.dataverse.org/en/5.9/installation/r-rapache-tworavens.html">https://guides.dataverse.org/en/5.9/installation/r-rapache-tworavens.html</a>
<b>Digital Preservation Coalition</b>	<a href="https://www.dpconline.org/handbook/technical-solutions-and-tools/tools">https://www.dpconline.org/handbook/technical-solutions-and-tools/tools</a>
<b>Apps, Dataverse Project</b>	<a href="https://guides.dataverse.org/en/latest/api/apps.html">https://guides.dataverse.org/en/latest/api/apps.html</a>
<b>Archivematica</b>	<a href="https://www.archivematica.org/es/">https://www.archivematica.org/es/</a>

### 6.4. Others

<b>Dataverse Project</b>	<a href="https://guides.dataverse.org/en/latest/">https://guides.dataverse.org/en/latest/</a>
<b>Wikipedia</b>	<a href="https://es.wikipedia.org/wiki/Open_Archival_Information_System">https://es.wikipedia.org/wiki/Open_Archival_Information_System</a>
<b>NDSA</b>	<a href="https://ndsa.org/publications/levels-of-digital-preservation/">https://ndsa.org/publications/levels-of-digital-preservation/</a>
<b>CoreTrustSeal</b>	<a href="https://www.coretrustseal.org/">https://www.coretrustseal.org/</a>

## Appendix

---

### 1. Final questions

1. How many complete copies the preservation system achieves?
2. In how many different geographic locations are the copies stored?
  - a. Interest in knowing whether they are local copies, private cloud, or public cloud b. Also, interest in knowing if copies access is online or offline.
3. How many of the copies have different disaster threat?
4. Is there, at least, one copy on a different storage media type?
  - a. Interest in knowing the type of media: optical media, hard disk, tape.
5. How often is object integrity checked? Is integrity verification performed for the full content? Does backup system allow for checking integrity of a part of the content?
6. What integrity calculation cryptographic functions are used?
7. What amount of computing resources do integrity checking consume?
8. What amount of computing resources do virus checking consume?
9. Is integrity information managed as special metadata? Is integrity information stored in a different destination from standard metadata?

## 2. Institutions interviews

<b>Name</b>	<b>Borealis</b>
<b>Country</b>	<b>Canada</b>
<b>Software</b>	<b>Dataverse</b>
<b>Analysis:</b>	
<b>1. How many complete copies the preservation system achieves?</b>	
<p><i>The system is achieving at least three copies. And it can be up to five</i></p> <p><i>Anytime anything is deposited into Dataverse, it is automatically copied into backup and also into the preservation system. Then that copy in preservation is replicated at least three times across the Ontario Library research cloud, which is cloud-based. The infrastructure has nodes at five different universities in Ontario,</i></p> <p><i>Daily backup of all files to tape using IBM Tivoli Storage Manager (TSM)</i></p> <p><i>For active files:</i></p> <p><i>Seven versions of a file are available for restore for 30 days</i></p> <p><i>If a file has not been modified for over 30 days, the most recent version of the file is retained permanently in backup</i></p> <p><i>The six previous versions of a file are discarded after 30 days</i></p> <p><i>For deleted files:</i></p> <p><i>The most recent version of a deleted file is available for restore for 60 days</i></p> <p><i>Two copies of the tape backup are retained onsite and one copy is retained offsite</i></p>	
<b>2. In how many different geographic locations are the copies stored?</b>	
<p><b>a) Interest in knowing whether they are local copies, private cloud or public cloud</b></p> <p><b>b) Also, interest in knowing if copies access is online or offline</b></p>	
<p><i>Copies are in different nodes. One of them is the University of Toronto. The other ones are, in Queen's University, which is based in Kingston Ontario, Guelph University, which is in Guelph Ontario, the University of Ottawa which is based in Ottawa. The other one is in McMaster based in Hamilton Ontario. Copies are replicated. Across those nodes.</i></p> <p><i>A geographic radius of like 450 kilometers away from each other</i></p>	
<b>3. How many of the copies have different disaster threat?</b>	
<p><i>This can be deduced from the two topics above.</i></p>	
<b>4. Is there, at least, one copy on a different storage media type?</b>	
<p><b>a) Interest in knowing the type of media: optical media, hard disk, tape.</b></p>	
<p><i>Backup uses tape</i></p> <p><i>One in a private cloud it is not known</i></p>	
<b>5. How often is object integrity checked? Is integrity verification performed for the full content? Does the backup system allow for checking integrity of a part of the content?</b>	
<p><i>A monthly check of all the files in the dataverse is performed</i></p> <p><i>They check against all the files in the cloud storage</i></p> <p><i>They have been running the file integrity for two years and they have found 124 corrupted files, and they found that most of them were uploaded with zero</i></p> <p><i>Integrity checks are stored separately from the content</i></p>	
<b>6. What integrity calculation cryptographic functions are used?</b>	
<p><i>Dataverse uses MD5 and their storage uses MD5</i></p>	
<b>7. What amount of computing resources do integrity checking consume?</b>	
<p><i>There aren't measures. It takes 24 hours to run, six terabytes of data</i></p>	
<b>8. What amount of computing resources do virus checking consume?</b>	
<p><i>When there is an upload, there's virus checking happening, there is no checking on the storage side</i></p>	
<b>9. Is integrity information managed as special metadata?</b>	

<b>Is integrity information stored in a different destination from standard metadata?</b>
<i>It is stored separately, the record of each fixity check (both positive and negative) is stored in an internal MySQL database</i>

<b>Name</b>	DataverseNL
<b>Country</b>	Holanda
<b>Software</b>	Dataverse
<b>Website</b>	https://dataverse.nl
<b>Related institutions</b>	4TU.ResearchData NIOO-KNAW Trimbos Instituut Utrecht University Vrije Universiteit Amsterdam Hanzehogeschool Groningen University of Applied Sciences Leiden University Avans University of Applied Sciences Tilburg University Protestantse Theologische Universiteit Maastricht University University of Groningen University Medical Center Utrecht Hogeschool Rotterdam Fontys Hogeschool Rijksdienst voor het Cultureel Erfgoed

**Analysis:**

<b>1. How many complete copies the preservation system achieves?</b>
Has two total amount of backups replicas
<b>Quantitative information of interest:</b>
<ul style="list-style-type: none"> <li>Amount of backup replicas</li> </ul>
Two
<b>Additional or reinforcement information of interest:</b>
<ul style="list-style-type: none"> <li>Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that ended up with the current amount of backup replicas?</li> </ul>
There is not additional information

<b>2. In how many different geographic locations are the copies stored?</b>
a) Interest in knowing whether they are local copies, private cloud or public cloud
b) Also, interest in knowing if copies access is online or offline
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>Copies are online/offline</li> <li>Storage location strategy and physical separation of copies (local copies, private cloud, or public cloud)</li> </ul>
<b>Geographically distributed storage locations</b>
20 km
<b>Additional or reinforcement information of interest:</b>
<ul style="list-style-type: none"> <li>Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that led to the current geographic location of the copies and the local vs cloud strategy?</li> </ul>
There is not additional information

<b>3. How many of the copies have different disaster threat?</b>
--



<p>There is not information</p> <p><b>Quantitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Main threat for each copy</i></li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about specific threats of the copies? This will allow CSUC whether they have similar threats.</i></li> </ul>
<p><b>4. Is there, at least, one copy on a different storage media type?</b></p> <p>a) <b>Interest in knowing the type of media: optical media, hard disk, tape.</b></p> <p>Copies storage type is snapshot.</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>List of backup types</i></li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current list of used backup types?</i></li> </ul> <p>Snapshot</p>
<p><b>5. How often is object integrity checked? Is integrity verification performed for the full content? Does the backup system allow for checking integrity of a part of the content?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Full contents / partial content (integrity information)</i></li> <li>• <i>Integrity check frequency</i></li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy/ technique to manage full/ partial integrity information?</i></li> </ul>
<p><b>6. What integrity calculation cryptographic functions are used?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Cryptographic function (integrity calculation)</i></li> <li>- Universal Numerical Fingerprint (UNF)</li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current cryptographic function to calculate and verify the integrity information?</i></li> <li>• <i>What software tools are involved in the generation/ verification of the integrity information?</i></li> </ul>
<p><b>7. What amount of computing resources do integrity checking consume?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Consumption of computing resources for integrity verification</i></li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to integrity verification?</i></li> </ul>
<p><b>8. What amount of computing resources do virus checking consume?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Consumption of computing resources for virus checking</i></li> </ul>

<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to the virus inspection?</i></li> </ul>
<p><b>9. Is integrity information managed as special metadata?</b>  <b>Is integrity information stored in a different destination from standard metadata?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Integrity information (metadata) storage strategy</i></li> </ul>
<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy about the way integrity information is treated in relation to metadata?</i></li> </ul>

<b>Name</b>	<b>DataverseNO</b>
<b>Country</b>	<b>Norway</b>
<b>Software</b>	<b>Dataverse</b>
<b>Website</b>	<b>https://dataverse.no/</b>
<b>Related institutions</b>	<b>Inland Norway University of Applied Sciences</b> <b>MF Norwegian School of Theology, Religion and Society</b> <b>Nord University</b> <b>Norwegian University of Life Sciences (NMBU)</b> <b>NTNU – Norwegian University of Science and Technology</b> <b>UiT The Arctic University of Norway</b> <b>Nofima</b> <b>University of Agder</b> <b>University of Bergen</b> <b>University of Oslo</b> <b>University of Stavanger</b> <b>VID Specialized University</b> <b>Western Norway University of Applied Sciences</b> <b>Østfold University College</b>

**Analysis:**

<p><b>1. How many complete copies the preservation system achieves?</b></p> <p>Has two total amount of backups replicas</p> <p><b>Quantitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• <i>Amount of backup replicas</i></li> </ul> <p>Three (for data files; two for database/metadata records)</p> <p><i>Additional or reinforcement information of interest:</i></p> <ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that ended up with the current amount of backup replicas?</i></li> </ul> <p>There is not additional information</p>
<p><b>2. In how many different geographic locations are the copies stored?</b></p> <p>a) Interest in knowing whether they are local copies, private cloud or public cloud</p> <p>b) Also, interest in knowing if copies access is online or offline</p>

<p>There are two location copies stored in two separate datacenters within the same building in Oslo, Norway. The two datacenters are separated with a fire-resistant wall. In addition, there is one copy stored in a Microsoft Azure Datacenter in The Netherlands. Currently implementing storage of another additional copy on local, immutable storage.</p> <p>a) Two copies are in private cloud at partner institution in Oslo, Norway. One copy is in private Virtual Data Center in public cloud (Microsoft Azure) in The Netherlands.</p> <p>b) The copies access is online.</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Copies are online/ offline</li> <li>• Storage location strategy and physical separation of copies (local copies, private cloud, or public cloud)</li> </ul> <p>Three copies are stored online. Currently implementing storage of fourth copy on immutable storage.</p> <p><b>Additional or reinforcement information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that led to the current geographic location of the copies and the local vs cloud strategy?</li> </ul> <p>The choice of storage mode and location is based on the UiT cloud deployment policy.</p>
<p><b>3. How many of the copies have different disaster threat?</b></p> <p>There is not information</p> <p><b>Quantitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Main threat for each copy</li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about specific threats of the copies? This will allow CSUC whether they have similar threats.</li> </ul>
<p><b>4. Is there, at least, one copy on a different storage media type?</b></p> <p>a) Interest in knowing the type of media: optical media, hard disk, tape.</p> <p>All copies are on hard disk. Currently implementing storage of another additional copy on local, immutable storage.</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• List of backup types</li> </ul> <p>Back up by the storage at partner institution in Oslo, Norway, and in Azure blob storage in virtual datacenter in The Netherlands.</p> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current list of used backup types?</li> </ul> <p>Snapshot</p>
<p><b>5. How often is object integrity checked? Is integrity verification performed for the full content? Does backup system allow for checking integrity of a part of the content?</b></p> <p>At file ingest. After that relying on file integrity support from cloud provider.</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Full contents / partial content (integrity information)</li> <li>• Integrity check frequency</li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy/ technique to manage full/partial integrity information?</li> </ul>
<p><b>6. What integrity calculation cryptographic functions are used?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Cryptographic function (integrity calculation)</li> </ul> <p>- Universal Numerical Fingerprint (UNF)</p> <p>- MD5</p>

<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current cryptographic function to calculate and verify the integrity information?</li> <li>• What software tools are involved in the generation/ verification of the integrity information?</li> </ul>
--

<p><b>7. What amount of computing resources do integrity checking consume?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Consumption of computing resources for integrity verification</li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to integrity verification?</li> </ul>
--

<p><b>8. What amount of computing resources do virus checking consume?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Consumption of computing resources for virus checking</li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to the virus inspection?</li> </ul>
--

<p><b>9. Is integrity information managed as special metadata?</b></p> <p><b>Is integrity information stored in a different destination from standard metadata?</b></p> <p>There is not information</p> <p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Integrity information (metadata) storage strategy</li> </ul> <p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy about the way integrity information is treated in relation to metadata?</li> </ul>
---

<b>Name</b>	<b>FAIRDATA CSC IT Center for Science</b>
<b>Country</b>	<b>Finlandia</b>
<b>Software</b>	<b>Fairdata</b>
<b>Website</b>	<b><a href="https://www.fairdata.fi/en/dps-for-research-data/">https://www.fairdata.fi/en/dps-for-research-data/</a></b>
<b>Related institutions</b>	<b>CSC – IT Center for Science Ltd.</b>

**Analysis:**

<p><b>1. How many complete copies the preservation system achieves?</b></p> <p>Has four total amount of backups replicas</p> <p><b>Quantitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Amount of backup replicas</li> </ul> <p>Four</p> <p><b>Additional or reinforcement information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that ended up with the current amount of backup replicas?</li> </ul>
--

There is not additional information
<p><b>2. In how many different geographic locations are the copies stored?</b>  <b>a) Interest in knowing whether they are local copies, private cloud or public cloud</b>  <b>b) Also, interest in knowing if copies access is online or offline</b></p>
<p><i>Quantitative &amp; qualitative information of interest:</i></p> <ul style="list-style-type: none"> <li>• Copies are online/offline</li> <li>• Storage location strategy and physical separation of copies (local copies, private cloud, or public cloud)</li> <li>• 3 online copies &amp; one offline copy</li> <li>• EMP-protected deep underground shelter, with airgap</li> <li>• geospatially distributed into 3 separate locations (data centers in Finland)</li> <li>• 1st (local disk storage) + 2nd (Tape storage) + 3rd (Tape storage) + 4th/5th (Dark archive)</li> </ul>
<p><i>Additional or reinforcement information of interest:</i></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that led to the current geographic location of the copies and the local vs cloud strategy?</li> </ul>
There is not additional information
<p><b>3. How many of the copies have different disaster threat?</b></p>
There is not information
<p><b>Quantitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Main threat for each copy</li> </ul>
<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about specific threats of the copies? This will allow CSUC whether they have similar threats.</li> </ul>
<p><b>4. Is there, at least, one copy on a different storage media type?</b>  <b>a) Interest in knowing the type of media: optical media, hard disk, tape.</b></p>
<p><i>Quantitative &amp; qualitative information of interest:</i></p> <ul style="list-style-type: none"> <li>• List of backup types</li> </ul>
Snapshot (virtual machine)
<p><i>Reinforcing or additional information of interest:</i></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current list of used backup types?</li> </ul>
<p><b>5. How often is object integrity checked? Is integrity verification performed for the full content?  Does the backup system allow for checking integrity of a part of the content?</b></p>
There is not information
<p><i>Quantitative &amp; qualitative information of interest:</i></p> <ul style="list-style-type: none"> <li>• Full contents / partial content (integrity information)</li> <li>• Integrity check frequency</li> </ul>
<p><i>Reinforcing or additional information of interest:</i></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy/ technique to manage full/partial integrity information?</li> </ul>
<p><b>6. What integrity calculation cryptographic functions are used?</b></p>
<p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Cryptographic function (integrity calculation)</li> </ul>
Signed SMIME manifests
<p><i>Reinforcing or additional information of interest:</i></p>

<ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current cryptographic function to calculate and verify the integrity information?</i></li> <li>• <i>What software tools are involved in the generation/ verification of the integrity information?</i></li> </ul>
---

<b>7. What amount of computing resources do integrity checking consume?</b>
There is not information
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Consumption of computing resources for integrity verification</i></li> </ul>
<b>Reinforcing or additional information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to integrity verification?</i></li> </ul>

<b>8. What amount of computing resources do virus checking consume?</b>
There is not information
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Consumption of computing resources for virus checking</i></li> </ul>
<b>Reinforcing or additional information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to the virus inspection?</i></li> </ul>

<b>9. Is integrity information managed as special metadata? Is integrity information stored in a different destination from standard metadata?</b>
There is not information
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Integrity information (metadata) storage strategy</i></li> </ul>
<b>Reinforcing or additional information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy about the way integrity information is treated in relation to metadata?</i></li> </ul>

<b>Name</b>	<b>ODUM</b>
<b>Country</b>	<b>Estats Units</b>
<b>Software</b>	<b>Dataverse</b>
<b>Website</b>	<b>https://odum.unc.edu</b>
<b>Related institutions</b>	<b>Odum institute for research in social science</b>

**Analysis:**

<b>1. How many complete copies the preservation system achieves?</b>
Has four total amount of backups replicas
<b>Quantitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Amount of backup replicas</i></li> </ul>
Four
<b>Additional or reinforcement information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that ended up with the current amount of backup replicas?</i></li> </ul>
There is not additional information

<p><b>2. In how many different geographic locations are the copies stored?</b>  <b>a) Interest in knowing whether they are local copies, private cloud or public cloud</b>  <b>b) Also, interest in knowing if copies access is online or offline</b></p>
<ul style="list-style-type: none"> <li>- Copies housed in geographically distributed storage locations             <ul style="list-style-type: none"> <li>- Offsite local server</li> <li>- Amazon S3 storage in Northern Virginia</li> <li>- Amazon S3 storage in Northern California</li> <li>- A copy in the Unitrends private cloud</li> </ul> </li> </ul>
<p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Copies are online/ offline</li> <li>• Storage location strategy and physical separation of copies (local copies, private cloud, or public cloud)</li> </ul>
<p><b>Additional or reinforcement information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC regarding this point of interest of CSUC, as well as documentation about the decision-making process that led to the current geographic location of the copies and the local vs cloud strategy?</li> </ul> <p>There is not additional information</p>
<p><b>3. How many of the copies have different disaster threat?</b></p> <p>There is not information</p>
<p><b>Quantitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Main threat for each copy</li> </ul>
<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about specific threats of the copies? This will allow CSUC whether they have similar threats.</li> </ul>
<p><b>4. Is there, at least, one copy on a different storage media type?</b>  <b>a) Interest in knowing the type of media: optical media, hard disk, tape.</b></p>
<p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• List of backup types</li> </ul> <p>Snapshot (virtual machine)</p>
<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current list of used backup types?</li> </ul>
<p><b>5. How often is object integrity checked? Is integrity verification performed for the full content?</b>  <b>Does the backup system allow for checking integrity of a part of the content?</b></p> <p>There is not information</p>
<p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Full contents / partial content (integrity information)</li> <li>• Integrity check frequency</li> </ul>
<p><b>Reinforcing or additional information of interest:</b></p> <ul style="list-style-type: none"> <li>• Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy/ technique to manage full/partial integrity information?</li> </ul>
<p><b>6. What integrity calculation cryptographic functions are used?</b></p>
<p><b>Quantitative &amp; qualitative information of interest:</b></p> <ul style="list-style-type: none"> <li>• Cryptographic function (integrity calculation)</li> </ul> <p>MD5 (data and metadata)</p>
<p><b>Reinforcing or additional information of interest:</b></p>

<ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current cryptographic function to calculate and verify the integrity information?</i></li> <li>• <i>What software tools are involved in the generation/ verification of the integrity information?</i></li> </ul>
---

<b>7. What amount of computing resources do integrity checking consume?</b>
There is not information
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Consumption of computing resources for integrity verification</i></li> </ul>
<b>Reinforcing or additional information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to integrity verification?</i></li> </ul>

<b>8. What amount of computing resources do virus checking consume?</b>
There is not information
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Consumption of computing resources for virus checking</i></li> </ul>
<b>Reinforcing or additional information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation, such as monitoring data or traces, that can be shared with the CSUC about the consumption of IT resources dedicated to the virus inspection?</i></li> </ul>

<b>9. Is integrity information managed as special metadata?</b>
<b>Is integrity information stored in a different destination from standard metadata?</b>
There is not information
<b>Quantitative &amp; qualitative information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Integrity information (metadata) storage strategy</i></li> </ul>
<b>Reinforcing or additional information of interest:</b>
<ul style="list-style-type: none"> <li>• <i>Is there documentation that can be shared with CSUC about the data preservation decision-making process that led to the current strategy about the way integrity information is treated in relation to metadata?</i></li> </ul>



### 3. Grid

[https://docs.google.com/spreadsheets/d/e/2PACX-1vSiArVX1M67IDwIyhPuccn\\_-RgTFUeA8mF2pGzNu\\_xEm0HjNsMskFoUPO10cOcTjA/pubhtml#](https://docs.google.com/spreadsheets/d/e/2PACX-1vSiArVX1M67IDwIyhPuccn_-RgTFUeA8mF2pGzNu_xEm0HjNsMskFoUPO10cOcTjA/pubhtml#)

#### **4. Recording and transcription of the interviews**

The recordings and transcriptions were sent to CSUC; however, they are not included in this report.