

# The Arkivum Solution for ARCHIVER

Matthew Addis Arkivum

# Agenda

- Introduction (10 mins)
- Solution Overview (20 mins)
- Solution Training: Part 1 (30 mins)
- Break (5 mins)
- Solution Training: Part 2 (30 mins)
- Solution Training: Part 3 (40 mins)
- Questions and Discussion (30 mins)



# Introduction (10 mins)

lacksquare

# Agenda

- Arkivum Overview
- Product and Services
- Certification, Standards and Good Practice
- Approach to ARCHIVER
- Partnership with Google
- Availability of the solution in EOSC



# About Arkivum

- Founded in 2011 out of the University of Southampton – initial focus on Higher Education Research Data
- VC backed and funded
- Headquartered in Reading, UK
- ~25 full time employees
- ISO 9001 and 27001 certified
- Full SaaS offering launched in 2017
- +50 customers as of 2022 across Higher Education, Heritage, Pharma and Life Sciences



# Arkivum Fully Managed Archiving and Preservation Solution



# Standards, Certification, Quality and Assessment

- Certified to ISO 9001 and ISO 27001
- SoPs for everything we do as part of our QMS
- Validated product releases
- GxP audited by customers, with a 100% success rate
- Mappings available of our product/services to NDSA preservation levels, DPC RAM, Core Trust Seal and FAIR







## Services and Good Practice

Professional and Customer Services

- Requirements analysis, SoPs and workflows
- Onboarding
- Data migrations
- Computer Systems Validation
- Integrations

#### Advocacy, Education, Advice

- Regular blog content
- eBooks
- Monthly webinars
- Research reports
- Conferences and events
- Promotion of LTDP practices into in new domains



**Our Product** 

# Approach to ARCHIVER

- Collaborative and iterative development and testing
- Knowledge sharing across domains (LTDP, RDM, data science)
- Real world tests and benchmarking











- GCP provides a highly scalable infrastructure (compute, storage, networking)
- Connected to GÉANT and NRENs across Europe
- GCP is widely used as a cloud platform for scientific data processing
- Low carbon footprint and good environmental credentials
- Discounts for education/research
- Arkivum is a Google Reseller
- Combined solution in ARCHIVER



 $https://indico.cern.ch/event/1008656/contributions/4232849/attachments/2192041/3705047/physics\_analysis\_google\_meeting-16.02.2021.pdf$ 

# European Open Science Cloud (EOSC) and related projects

- EOSC Marketplace listing
  - Early Adopters Programme
- Ordering possible through EOSC
  - Via the Arkivum website
- Active involvement in past, present and future EOSC events
- Future initiatives and projects
  - EOSC TF DP

	open .oud	Find resource	All resources 🗸	Q My EOSC Marketpl
Access physical & elf Access physical & elf	ifrastructures > Da	ta Storage + Digital Preservation + Arkivum Digital Archiving and Preservation Solution	on	
•	Arkivur	n Digital Archiving and Preservation Solu	ution	
arkivum	Ensuring rese	earch data lives forever		ccess the resource
	ជំដំដំដំ ជំដំដំដំ	(0.0 /5) O reviews Add to comparison Add to favourites		ORDER REQUIRED
	→ Webpage	→ Helpdesk → Helpdesk e-mail	Ask a q	uestion about this resource?
ABOUT DETAILS	REVIEWS (	))		
-iif-			SCIEN	ITIEIC CATEGORISATION
hort, their data management must	align with the FA	Intes of complex data that often needs to be open, shared and accessio IR data management principles.	de lor decades. In	
Current solutions keep this data on o scale.	ageing and hard	to manage in-house systems that are expensive, challenging to maintain	n and are difficult	Sther
kivum provides a digital archiving i anagement use cases of the scier fective and environmentally sustair	and preservation ntific research co nable way.	solution that has been built to meet the varied and challenging long-term mmunity. Our technology is able to preserve petabyte level datasets, in	n data the most cost-	Other
We ensure that data is: - Safeguarded: we archive multiple copies of customers data while maintain data integrity - Findable: easily find the right research data as quickly as possible -Accessible: data is accessible to those who need it, now and in the future - Usable: regardless			GORISATION	
now wild the usia is stored for, th	and contractiloo	nar e mil po rozovoro a lo papilo lo grano lo grano lo lo lotoro i loggi i digital pras	UT YMMUT I	Data Storage

# Solution Overview (20 mins)

Ο

# Agenda

- LTDP of Research Datasets
- Architecture and Arkivum Solution
- Integration into RDM landscape of Repositories and FAIR
- SSO using eduGAIN
- Scalability and Performance
- Deployment options: cloud and on-premise
- Cost-efficiency and Environmental Sustainability



# Long-Term Digital Preservation of Research Data: Functionality



Scientific use cases deployments: https://www.archiver-project.eu/deployment-scenarios

ARCHIVER "current state of the art" report in the context of the EOSC: https://doi.org/10.5281/zenodo.3618215

## Long-Term Digital Preservation of Research Data: Good Practice



### Standards, External Assessment and Certification



Digital Preservation Coalition

# **FAIR Forever?**

Long Term Data Preservation Roles and Responsibilities, Final Report

February 2021 (V.7)

Dr Amy Currie and Dr William Kilbride

ं

EOSC FAIR Forever has received funding from the European Union under the EOSC Secretariat project. EOSC secretariate. The received funding from the European Union's Horizon Programme call H2020-INFRAECSC-05-2018-2019, grant Agreement number 831644. Europeana is an initiative of the European Union, financed by the European Union's Connecting Europe Facility and European Union Member States (https://pro.europeana.eu/ and https://www.europeana.eu)

# Long-Term Digital Preservation (LTDP) Factories



# Arkivum Microservices and Workflows





#### **NDSA Preservation Levels**



# **DPC RAM**

I - Content preservation Processes to preserve the meaning or functionality of the digital content and ensure its continued accessibility and usability over time	3 - Managed	<ul> <li>The organization has implemented a managed process to monitor and plan for accessibility of content over time, for example:</li> <li>Technology watch activities are carried out and 'at risk' content is identified.</li> <li>Technical dependencies are detected and documented.</li> <li>Actions are occasionally carried out to ensure preservation and quality of content such as migration, emulation or modification of creation or capture workflows.</li> <li>Preservation actions occur with an understanding of the properties of the digital object that should be retained to support current and future use cases.</li> <li>All changes to digital content are recorded, including details of when, what, how, why and who.</li> </ul>	• • • • •	The solution supports file format identification using a range of tools (Tika, Siegfried, FIDO). The system records a wide range of file characteristics including file type, date, size, filepath and checksums. Technical metadata extraction and file format characterisation provides detailed technical information on the types of content being stored in the solution. File format verification is supported, e.g. using JHOVE, MediaConch and VeraPDF. File format migrations/normalisations can be done according to configurable rules, e.g. using the Archivematica Format Policy Register. Users can search for content according to file format or other technical characteristics. This is in addition to searching by descriptive metadata. Reports are available on content in the solution, e.g. normalisation reports on file format conversion done during ingest. File format identification, characterisation, validation and normalisation are all recorded in the system audit trail.	The Arkivum solution supports the execution of preservation actions, e.g. file format normalisations, but it is up to the organisation using the solution to assess their risks and decide the preservation plans/policies most suitable for them.
---	----------------	---	-----------	---	---

#### **Repositories and FAIR**



**Fastering Fair Data Pract** 

# 

My dashboard

2019

Export

+-

#### O Preview

You are previewing a new record that has not yet been published.

arkivum =			< Back to edit			
navigation	Home > C_test_data/CERN/HiggsToBBNtupleProducerTool/ > O_test_data/CERN/H			Versions		
<ul> <li>☐ Dashboard +</li> <li>☆ Admin +</li> </ul>			Published 2019   Version v1  Dataset  Metadata-only  Sample with iet, track and secondary vertex properties	Preview     Only published versions are displaye		
🗅 Files 🗕	Content Type	ID	for Hbb tagging ML studies			
File Explorer Files to Import		ntuple_merged_0.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	HiggsToBBNTuple_HiggsToBB_QCD_RunII_13TeV_MC	Version v1 2		
$\equiv$ Reports +	EVE	ntuple_merged_1.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTooI/Test/ntuple_merged	Duarte, Javier	Details		
🕮 Rules +	-	ntuple_merged_2.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	Citation Style APA -	Resource type		
Manual Ingest     Notifications	-	ntuple_merged_3.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	Duarte, J. (2019). Sample with jet, track and secondary vertex properties for Hbb tagging	Dataset		
	-	ntuple_merged_4.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	ML studies HiggsToBBNTuple_HiggsToBB_QCD_RunII_13TeV_MC [Data set]. CERN Open Data Portal.	Publisher CERN Open Data Portal		
		ntuple_merged_5.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged		Rights		
	RE	ntuple_merged_6.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	Description The dataset consists of particle jets extracted from simulated proton-proton collision events at a center-of-mass energy of	()))) Creative Commons Zero v1. Universal		
	RE	ntuple_merged_7.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	13 TeV generated with Pythia 8. It has been produced for developing machine-learning algorithms to differentiate jets originating from a Higgs boson decaying to a bottom quark-antiquark pair (Hbb) from quark or gluon jets originating from			
	122	<b>ntuple_merged_8.h5</b> /arkivum/higgs1/HiggsToBBNtupleProducerTooI/Test/ntuple_merged	quantum chromodynamic (QCD) multijet production. The reconstructed jets are clustered using the anti-kT algorithm with R=0.8 from particle flow (PF) candidates (AK8 jets).	Export		
	77	files.txt /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/files.txt	The standard L1+L2+L3+residual jet energy corrections are applied to the jets and pileup contamination is mitigated using the charged hadron subtraction (CHS) algorithm. Features of the AK8 jets with transverse momentum pT > 200 GeV and pseudorapidity $ \mathbf{n}  < 2.4$ are provided. Selected features of inclusive (both charged and pseudorapidity $ \mathbf{n}  < 2.4$ are provided.	JSON - Expo		
	FRE	ntuple_merged_10.h5 /arkivum/higgs1/HiggsToBBNtupleProducerToo1/Train/ntuple_merged	0.95 GeV associated to the AK8 jet are provided. Additional features of charged PF candidates (formed primarily by a charged particle track) with $pT > 0.95$ GeV associated to the AK8 jet are also provided. Finally, additional features of			
		ntuple_merged_11.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/ntuple_merged.	4_11.h5 /arkivum/higgs1/Hi ···			
	-	ntuple_merged_12.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/ntuple_merged.	1_12.h5 /arkivum/higgs1/Hi ···			
		ntuple_merged_13.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/ntuple_merged.	I_13.h5 /arkivum/higgs1/Hi ···			
		Entries per page:	: 15 ~ < 1 2 3 () 7 > >>			

# Other Repository Integrations





## CoreTrustSeal and FAIR

						A 1 (meta)dist air efertiveaue by their electricity of the antitric using a stataana commissions protocol. A 11 the protocol sopen, free, and universally implementable A12 the protocol allows for an authentication and authorizatio where necessary. A2 metadata are accessible, even when the data are no longer a K2.0 reservation plan	zea /s context) 1 procedure, /ailable.
Area	Short Requirement Name	FAIR Principle ( <u>Nature</u> )	RDA Indicator	FAIRsFAIR Metric	Evidence Links	II. (meta)data use a formal, accessible, shared, and broadly app knowledge representation. 12. (meta)data use vocabularies that follow FAR principles. 19. (meta)data incluée wulfidin defences to other (meta)data	cable language for
Digital Object Management	14. Data reuse	<ol> <li>(meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.</li> </ol>	RDA-I1-01D Data uses knowledge representation expressed in standardised format (Important) RDA-I1-02D Data uses machine-understandable knowledge representation (Important) RDA-I1-01M Metadata uses knowledge representation expressed in standardized format (Important) RDA-I1-02M Metadata uses machine-understandable knowledge representation (Important)	FsF-I1-01M Metadata is represented using a formal knowledge representation language.	Metadata in the Arkivum solution is accessible in json or CSV format. Metada Metadata can be searched and retrieved using a REST API. Metadata files car Including Json-Id, and these metadata files can be associated to datasets, e.g. Metadata in InvenioRDM is in json format and can be searched and exported vi InvenioRDM landing pages include record citation metadata in machine readab	A (neil) plat include qualities is orbity (metapata). R metaplata in rolly described with a planning of accurate a R1.metaplata) are rolly described with a clear and accessible data with R1.1. (metaplata are rolly described with a clear and accessible data with R1.2. Unorange and a seasociated with detailed provenance. R1.2. Unorange are associated with detailed provenance. R1.2. Unorange and a seasociated with detailed provenance. R1.8. Returns R1.8. Returns	d relevant age license.
		l2. (meta)data use vocabularles that follow FAIR principles	RDA-I2-01M Metadata uses FAIR-compliant vocabularies (Important) RDA-I2-01D Data uses FAIR-compliant vocabularies (Useful)	FsF-I2-01M Metadata uses semantic resources.	Metadata files can be ingested into the Arkivum system in their native format ( Include references to external schemas or vocabs, e.g. schema.org, wikidata. InvenioRDM supports vocabularies. https://inveniosoftware.org/products/rdm/	(ML or Json). These metadata files can oadmap/	
		l3. (meta)data include qualified references to other (meta)data	RDA-I3-01M Metadata includes references to other metadata (Important) RDA-I3-02M Metadata includes references to other data (Useful) RDA-I3-01D Data includes references to other data (useful) RDA-I3-02D Data includes qualified references to other data (Useful) RDA-I3-03M Metadata includes qualified references to other metadata (Important) RDA-I3-04M Metadata include qualified references to other data (Useful)	FsF-I3-01M Metadata includes links between the data and its related entities.	The Arkivum solution supports DataCite relation types to allow entities in the s each other. Each end of the relationship is defined using identifiers. https://support.datacite.org/docs/relationtype_for_citation invenioRDM supports DataCite relation types to describe the relationship(s) of https://inveniordm.docs.cern.ch/reference/metadata/#related-identifiersworks-	ystem (files, datasets etc.) to be related to a record to other resources. 0-n	
	R1. meta( described of accurat attributes		RDA-R1-01M Plurality of accurate and relevant attributes are provided to allow reuse (Essential)	FsF-R1-01MD Metadata specifies the content of the data.	The Arkivum solution supports ingest of metadata files in their original domain provided using indexed and searchable fields (DC, DataCite, custom). Technical from files. Other forms of metadata include file format and checksums. InvenioRDM includes metadata on resource type and subtype using configurabl https://inveniordm.docs.cern.ch/reference/metadata/#resource-type-1 InvenioRDM supports metadata on formats. https://inveniordm.docs.cern.ch/re	specific format. Metadata can also be metadata can be automatically extracted e controlled vocabs. ference/metadata/#formats-0-n	
		R1.3. (meta)data meet domain-relevant community standards	RDA-R1.3-01M Metadata complies with a community standard (Essential) RDA-R1.3-02M Metadata is expressed in compliance with a machine-understandable community standard (Essential) RDA-R1.3-01D Data complies with a community standard (Essential) RDA-R1.3-02D Data is expressed in compliance with a machine-understandable community standard (Important)	FSF-R1.3-01M Metadata follows a standard recommended by the target research community of the data FsF-R1.3-02D Data is available in a file format recommended by the target research community.	Metadata files can be ingested into the Arkivum system in their native format ( metadata standards. Fields in metadata files can be extracted, mapped and in InvenioRDM supports records that include files. One or more of these files cou	KML or json), e.g. in domain specific dexed so metadata is searchable. d be a domain specific metadata format.	

F1. (meta)data are assigned a globally unique and persistent identifier. F2. data are described with rich metadata (defined by R1 below). F3. metadata (dearly and explicit) include the identifier of the data if describes. F4. (meta)data are registered or indexed in a searchable resource. F3.D. bad discovery and identification

A1 (meta)data are retrievable by their identifier using a standardized

E.

# **Cloud-hosted Processing of Archived Research Datasets**

- GCP hosted applications
  - Cloud Functions (serverless code)
  - Cloud Run (containers)
  - GKE (Kubernetes)
  - Compute Engine (VMs)
- Arkivum Webhooks
  - Ingest, preservation, access
- Arkivum REST API
  - including search, get metadata, export files
- Xrootd server
  - Easy integration of scientific apps





≡	Google Cloud	cern-dedicated 👻	Q Search Pro	ducts, resources,	docs (/)		v) ii D	) <b>.</b> ()	: (
٨	Kubernetes Engine	Workloads	C REFRESH	+ DEPLOY	DELETE	6	G OPERATIONS	r 🖻 HELP A	SSISTAN
<b>(</b> ])	Clusters	Cluster	▼ Na	mespace	- RESE	T SAVE			
- 54	Workloads	Wedde de ere deel	Workloads are deployable units of computing that can be created and managed in a cluster.						
A	Services & Ingress	cluster.							
	Applications	OVERVIEW	COST OPTIMIZATI	ON					
⊞	Secrets & ConfigMaps	Ţ Filter (Is s	system object : False (	Filter workloads	3			×ø	III
	Storage	Name 1	`	Status	Туре	Pods	Namespace	Cluster	
:=	Object Browser	reana-ca	che	📀 ОК	Deployment	1/1	default	reana	
	object browser	reana-db	1	🖉 ОК	Deplo				
A	Migrate to containers	reana-m	essage-broker	🛇 ок	State	Cal	na		
0	Backup for GKE	reana-re	source-quota-update	🕗 ок	Cron				
۲	Config Management	reana-ru 47cf-b96	n-batch-c0c4997f-612d 2-eaac1d1c4c24	с- 🕑 ОК	Job				
<b></b>	Protoct	reana-se	rver	📀 ОК	Deplo		Yo	our workf	lows
۲	FIOLECT NEW	reana-tra	efik	📀 ОК	Deplo				
		reana-ui		📀 ОК	Deplo				
		reana-we	orkflow-controller	📀 ОК	Deplo		S	earch	

Your workflows	€ Refreshed at 16:06:06 U	тс
Search	Q	•
Status 🝷	Sort by 👻	
root6-roofit-yadage-kubernetes #1 Finished 11 minutes ago	<b>finished</b> in 1 min 21 sec step 2/2	
<ul> <li>root6-roofit-snakemake- kubernetes #1</li> <li>Finished 14 minutes ago</li> </ul>	<b>finished</b> in 32 seconds step 2/2	
♥ root6-roofit-serial-kubernetes #1 Finished 15 minutes ago	<b>finished</b> in 19 seconds step 2/2	
root6-roofit-cwl-kubernetes #1 Finished 16 minutes ago	<b>finished</b> in 1 min 17 sec step 2/2	



# Portability

- Single software stack that can run on AWS, GCP and OpenStack
- Deployed on GCP and at CERN for ARCHIVER
- Ingest and export of data and metadata in native formats
- Interact with solution using standard protocols (s3 buckets, REST API, eduGAIN, Web UI)

aws

openstack.



Google Cloud

### Open Standards, Open Specifications, Open Source: Deployment



Open Standards, Open Specifications, Open Source: Interaction











# **Bagit & BDbags**







# eduGAIN and AAI

- Arkivum is a member of the UKAMF
- eduGAIN integration uses SimpleSAML php
  - WAYF service for selecting IdP
  - Routing of IdP authentication to systems/tenants
  - Mapping of user's IdP attributes to roles/groups
- Other SSO options possible:
  - LDAP and Active Directory
  - SAML and OpenID
  - Social (Twitter, Facebook, Instagram, Google etc.)



# Scalability and Performance

0

 $\mathbf{O}$ 

0

# Serverless Computing: Scalability, Performance, Efficiency

- Scalable and cost-effective archiving workflows and processing
  - Kubernetes and autoscaling
  - Scale to zero as well as autoscaling for peak loads (pods and nodes)
  - Pre-emptible nodes to reduce costs (up to 70% lower)
  - Terraform and Ansible for provisioning
  - Rancher, Prometheus, Grafana, Kibana for monitoring and analytics
- Microservices
  - Checksum, virus scan, file format identification, caching, replication, unpack ...
  - Stateless and able to run in parallel (can 100TB+ per day)
  - Jobs recorded and tracked











### Example: Ingest and Archiving of Astronomy Datasets



440,000 image files (25TB) ingested as 2780 big data bags within 24hrs



### **Preservation Formats and Access Formats**



Media type	File formats	Preservation format(s)	Access
Audio	AC3, AIFF, MP3, WAV, WMA	WAVE (LPCM)	MP3
Email	PST	MBOX/EML	PDF
Email	MSG	EML	PDF
Office docs and presentations	DOC, WPD, RTF, DOCX, PPTX, PPT	PDF/A	PDF/A
Plain text	TXT	Original format	Original format
Portable Document Format	PDF	PDF/A	Original format
Raster images	BMP, GIF, JPG, JP2*, PCT, PNG*, PSD, TIFF, TGA	TIFF	JPEG
	3FR, ARW, CR2, CRW, DCR, DNG, ERF, KDC,		
Raw camera files/Digital Negative format	MRW, NEF, ORF, PEF, RAF, RAW, X3F	TIFF	JPEG
Spreadsheets	XLS, XLSX	Original format	Original format
Vector images	AI, EPS, SVG	SVG	PDF
	AVI, FLV, MOV, MPEG-1, MPEG-2, MPEG-4, SWF,		
Video	WMV	FFV1/LPCM in MKV	MP4



#### Preferred Formats

Preferred Formats	Format Versions	Format Specifications
ASCII Text	7 bit	ISO/IEC 646:1991 Information technology ISO 7-bit coded character set for information interchange: (http://www.iso.org/iso/catalogue_detail.htm?csnumber=4777 🗷)
Unicode Text	UTF-8	RTF 3629: UTF-8, A Transformation Format of ISO 10646: ( http://tools.ietf.org/html/rfc3629 🗹)
	UTF-16	RFC 2781 UTF-16: An Encoding of ISO 10646: ( http://www.ietf.org/rfc/rfc2781.txt 🕜)
OpenDocument Text Format (ODF)	OpenDocument 1.0	ISO/IEC 26300:2006 Information technology OpenDocument Format for Office Applications (OpenDocument) v1.0: ( http://www.iso.org/iso/iso_catalogue/ catalogue_tc/catalogue_detail.htm?csnumber=43485 🕜)
PDF/A-1	PDF/A-1	ISO 19005-1:2005 Document management Electronic document file format for long-term preservation Part 1: Use of PDF 1.4 (PDF/A-1): (http://www.iso.org /iso/catalogue_detail?csnumber=38920 ⑦)
PDF/A-2	PDF/A-2	ISO 19005-2:2011 Document management Electronic document file format for long-term preservation Part 2: Use of ISO 32000-1 (PDF/A-2): (http://www.iso.org/iso/home/store /catalogue_tc/catalogue_detail.htm?csnumber=50655 ♂)
Acceptable Formats		
Acceptable Formats	Format Versions	Format Specifications

Acceptable Formats	Format Versions	Format Specifications
PDF	PDF 1.7	ISO 32000-1:2008 Document management Portable document format Part 1: PDF 1.7: (http://www.iso.org/iso/catalogue_detail.htm?csnumber=51502 7
	PDF 1.0-1.6	Adobe® Portable Document Format Version 1.6: http://www.adobe.com/devnet /pdf/pdf_reference_archive.html 🗹
Microsoft Word (DOCX) Office Open XML	OOXML Microsoft Word for Windows, version 2007-2010	[MS-OI29500]: Office Implementation Information for ISO/IEC 29500 Standards Support: ( http://msdn.microsoft.com/en-us/library/ee908652%28v=office.12%29 🗹)
Microsoft Word 97 Binary Document Format (DOC)	8.0	[MS-DOC]: Word (.doc) Binary File Format: ( http://msdn.microsoft.com/en-us/library /cc313153%28v=office.12%29.aspx 🖒
#### **Example: Preserving Documents**



Done

- 349k GOV DOCS files in batches of 1000 files, including Archivematica to generate AIPs and DIPs
- 142k normalised files (Office file formats -> PDF/A)

# Cloud Deployment, Costs, Environmental Sustainability

#### Leveraging the Scalability and Efficiency of Cloud Infrastructure



#### ENERGY

#### Recalibrating global data center energy-use estimates

Growth in energy use has slowed owing to efficiency gains that smart policies can help maintain in the near term

#### By Eric Masanet<sup>1,2</sup>, Arman Shehabi<sup>3</sup>, Nuoa Lei<sup>1</sup>, Sarah Smith<sup>3</sup>, Jonathan Koomey<sup>4</sup>

ata centers represent the information backbone of an increasingly digitalized world. Demand for their services has been rising rapidly (1), and data-intensive technologies such as artificial intelligence, smart and connected energy systems, distributed manufacturing systems, and autonomous vehicles promise to increase demand further (2). Given that data centers are energyintensive enterprises, estimated to account for around 1% of worldwide electricity use. these trends have clear implications for global energy demand and must be analyzed rigorously. Several oft-cited yet simplistic analyses claim that the energy used by the world's data centers has doubled over the past decade and that their energy

demand for data center services rises rapidly, so too must their global energy use. But such extrapolations based on recent service demand growth indicators overlook strong countervailing energy efficiency trends that have occurred in parallel (see the first figure). Here, we integrate new data from different sources that have emerged recently and suggest more modest growth in global data center energy use (see the second figure). This provides policy-makers and energy analysts a recalibrated understanding of global data center energy use, its drivers, and near-term efficiency potential. Assessing implications of growing demand for data centers requires robust understanding of the scale and drivers of global data center energy use that has eluded many policy-makers and energy analysts. The reason for this blind spot is a historical lack of "bottom-up" information

As demand for data centers rises, energy efficiency improvements to the IT devices and cooling systems they house can keep energy use in check.

Bottom-up analyses tend to best reflect this broad range of factors, generating the most credible historical and near-term energyuse estimates (7). Despite several recent national studies (8), the latest fully replicable bottom-up estimates of global data center energy use appeared nearly a decade ago. These estimates suggested that the worldwide energy use of data centers had grown from 153 terawatt-hours (TWh) in 2005 to between 203 and 273 TWh by 2010, totaling 11 to 1.5% of clobal electricity use (9).

Since 2010, however, the data center landscape has changed dramatically (see the first figure). By 2018, global data center workloads and compute instances had increased more than sixfold, whereas data center internet protocol (IP) traffic had increased by more than 10-fold (1). Data center storage capacity has also grown rapidly, increasing by an estimated factor of 25 over the same time period (1, 8). There has been a tendency among analysts to use such service demand trends to simply extrapolate earlier bottom-up energy values, leading to unreliable predictions of current and future global data center energy use (3-5). They might, for example, scale up previous bottom-up values (e.g., total data center energy use in 2010) on the basis of the growth rate of a service demand indicator (e.g., growth in global IP traffic from 2010 to 2020) to arrive at an estimate of future energy use (e.g., total data center energy use in 2020).

But since 2010, electricity use per computation of a typical volume server-the workhorse of the data center-has dropped by a factor of four, largely owing to processorefficiency improvements and reductions in idle power (10). At the same time, the watts per terabyte of installed storage has dropped by an estimated factor of nine owing to storage-drive density and efficiency gains (8). Furthermore, growth in the number of servers has slowed considerably owing to a fivefold increase in the average number of compute instances hosted per server (owing to virtualization), alongside steady reductions in data center power usage effectiveness (PUE, the total amount











#### Green Data Centers, Renewable Energy







#### Green Data Centers, Renewable Energy





https://www.google.com/about/datacenters/gallery/#hamina-exterior-landscape

https://cloud.withgoogle.com/region-picker/

#### Costs, Resource Consumption, Carbon Footprint



#### **Benchmarking and Metrics**

- Execute real world scenarios
- Record parameters
  - execution time
  - data volumes, number of files
  - type of activity (ingest, export, preservation)
- Extract costs and resource consumption from cloud provider
- Extract carbon footprint from cloud provider
- Add short-term and long-term storage
  - Upload/export buckets, caching, archive buckets
- Calculate metrics

Ingest	\$2.2 per TB
ingest	0.1 kgCO <sub>2</sub> eq
Long torm storage	\$30 per TB-year
Long-term storage	0.7 kgCO <sub>2</sub> eq per TB-year

# Summary

•

· O

•

0

0

• )

#### The Arkivum Solution for ARCHIVER

- Highly scalable LTDP capable of ingesting and preserving 100TB+ per day
- Co-locate scientific applications with archived data
- Integration with InvenioRDM for creating/publishing landing pages
- Serverless computing: only consume what's needed and when it's needed
- Cost-efficient and minimized carbon-footprint
- Deployment using GCP, AWS and on-premise
- Provided as a fully managed service / SaaS solution
- Supports LTDP requirements and models (DPC RAM and NDSA levels)
- Supports TDR and FAIR requirements and models (CoreTrustSeal and FAIR)









# Solution: Part 1

Ο

0

0

0

lacksquare

0

## Agenda

- Authentication and Authorisation (user interface and API)
- Overview of the Arkivum Web UI
- Configuration and Self-Service
- Getting data into the Arkivum solution



- Local Accounts
- eduGAIN

	BOSON		
Username or email simon.bostock@arkivum.com	Log In	AIN	
Password  Log In	Access t	Arkivum Service     Choose Your Institution     Recent institutions	
		UK federation test IdP test.ukfederation.org.uk	>
oson.archiver.arkivum.net/auth/realms/boson/broker/saml/login?cli		CERN Service Provider Proxy cern.ch Universidad de Sevilla us.es	>

#### Authorisation – User Roles

System Activities	Read Only	User	Superuser	Admin
Dashboard view of datapools				
Search screen and any datasets associated to those assigned datapools		$\checkmark$		
Download files and metadata				
Upload files into datapools (#1)		$\checkmark$		
Add and remove existing retentions to objects and collections				
Request preservation				
View retention rules (#1)				
Initiate bulk exports				
Import and export metadata				
Create retention rules				

#### Authorisation – User Roles

Reporting	Read Only	User	Superuser	Admin
Ingest Report				
Preservation Report		$\checkmark$		$\checkmark$
Normalisation Report		$\checkmark$		
Hold Report		$\checkmark$		$\checkmark$
Unhold Report				$\checkmark$
Processing Report				$\checkmark$
Export Report				
Audit Trail				$\checkmark$
Deletion Report				

Approval Workflows	Read Only	User	Superuser	Admin
Approve/reject ingests				
Approve/reject deletion requests				
Approve/reject Bulk Exports				

## https://vimeo.com/723625238/b5a3ee264a

## AAI Summary

- Arkivum is a member of the UKAMF
- Arkivum is a registered SP
  - REFEDS R&S, DP CoCo, Sirtfi
- Implementation using simpleSAML and keycloak
- WAYF allows users to select their home IdP
- Routing of authenticated users to specific tenants
  - e.g. user from CERN IdP -> cern.archiver.arkivum.net
- Mapping of user attributes to roles
  - e.g. User X with attribute Y that has value Z -> ROLE\_ADMIN
  - eduPersonEntitlement (part of REFEDS R&S bundle)



#### **Dashboard Overview**



#### **Dashboard Overview**

https://vimeo.com/722917823/17c05d473a

#### Arkivum Dashboard

- System usage: data volumes, throughput, processes
- Administration: datapools, metadata configuration, buckets, retentions and holds
- Ingest and Export: via buckets, direct upload/download
- Search and Navigation: browse and find datasets, files, records
- Reports: ingest, preservation, exports, deletion, retentions, holds, audit log
- Notifications: approvals, alerts
- Web based, cross-platform
- Everything in the UI is also available through a REST API

## Configuration

- Datapools
- Metadata namespaces
- Buckets

#### arkivum = navigation ŝ for Dashboard Admin Datapool Name Datapool Path Datapools /arkivum Arkivum Namespace /arkivum-preserved Arkivum Preserved Buckets Saved Queries /deep\_archive Deep Archive **Relation Kinds** Default ற Files Frequent Access /frequent\_access Reports D Rules /ilcdocs ILCDOCS Manual Ingest /invenio-rdm InvenioRDM /magic1 MAGIC1

Preservation Ena...

false

true

false

false

false

false

false

false

#### Datapools

arkivum	=						<b>A</b> *	Matthew Add	dis ~
navigation ด Dashboard	+								+
Admin		Datapool Name	Datapool Path	Preservation Ena	Location Set	Metadata Namespaces	Webhook URL	Encrypted	
Datapools     Namespace		Arkivum	/arkivum	false	quickaccess	technical,http://www.arkivum.com/xs		true	
Buckets		Arkivum Preserved	/arkivum-preserved	true	quickaccess	technical,http://www.arkivum.com/xs		true	
Saved Queries		Deep Archive	/deep_archive	false	deeparchive	technical,http://www.arkivum.com/xs		false	
Files	+	Default		false	quickaccess	technical,http://www.arkivum.com/xs		true	
Reports	+	Frequent Access	/frequent_access	false	quickaccess	technical,http://www.arkivum.com/xs		false	
	+	ILCDOCS	/ilcdocs	false	quickaccess	technical,http://www.arkivum.com/xs		false	
<ul> <li>Manual Ingest</li> <li>Notifications</li> </ul>		InvenioRDM	/invenio-rdm	false	quickaccess	dataciteTypes,dataciteDescriptions,te	http://invenio-rdm-1.arc	true	•••
		MAGIC1	/magic1	false	quickaccess	technical,magic_telescopes_stereo		true	

Datapool	Safeguarding	Preservation	InvenioRDM	Escrow	Read Only Access	User Access	Admin Access
Datapool 01					User 1	User Group 1	Admin Group
Datapool 02				$\checkmark$	User 2	User Group 2	Admin Group
Datapool 03						User Group 3	Admin Group

## Datapools

- Where data will be stored for the long-term
  - Frequent access or deep archive buckets
- What metadata fields and rules should be applied
  - Dublin Core, DataCite, domain specific
- What processes to run when data is ingested
  - Safeguarding and preservation
- Automatic cleanup
  - E.g. remove datasets after successful ingest
- File encryption
  - In addition to cloud storage, user provided keys
- User access
  - Users can be given roles that allow access to specific datasets
- Webhooks
  - Called at the end of the ingest and/or preservation process



## Datapools

Add new datapool	×
Name * Deep Archive	Path * /deep_archive
Encrypted Preserved	
Location set	Quota
Deep Archive v	1 ТВ
Namespaces	
Available namespaces	Selected namespaces
URL	URL
http://www.arkivum.com/xsd/atom	http://www.arkivum.com/xsd/dublincore
http://www.arkivum.com/xsd/isadg	
http://www.arkivum.com/xsd/esignature	
• mandatory fields	Discard Apply

#### Metadata Namespaces

arkivum	≡						
navigation 🖻 Dashboard	+	¢٩					
Admin	-	Namespace	Prefix	Aggregation Types	Grouped	Editable	Mandatory
Datapools		Dublin Core	dc	All types	false	true	false
Buckets		Usability	atom	All types	false	true	false
Saved Queries		ISAD(G)	isadg	0	false	true	false
Files	+	ESignature	eSignatures	F	true	false	false
Reports	+	Identifier	identifiers	All types	true	true	false
🕮 Rules	+	Technical	technical	F	false	false	false
Manual Ingest A Notifications		DataCite Creators	dataciteCreators	C, F, O	true	true	false
		DataCite Identifiers	dataciteldentifiers	C , F , O	true	true	false
		DataCite Types	dataciteTypes	C , O , F	false	true	false
		DataCite Descriptions	dataciteDescriptions	C , O , F	true	true	false
		DataCite	datacite	C , F , O	false	true	false
		MAGIC Telescopes Stereo	magic_telescopes	F , C , O	false	true	false

#### Metadata Namespaces

- Fields
  - Title, Subject, Date etc.
- Field Types
  - String, Date, Integer
- Repeatable or Single
  - Titles, Identifiers etc.
- Grouped
  - First Name, Last Name, ORCID, Affiliation
- Editable
  - Write once or updateable by users

			×
Uri •			
magic_telescopes_stereo_	metadata		
Prefix *			
magic_telescopes			
Label *			
MAGIC Telescopes Stereo			
Mandatory	Grouped	C Edito	ıble
Aggregation types			
Aggregation types			
All types			
🜔 0 🌔 F 🌔 C			
identifier - String Editable			
fname - String <sup>Editable</sup>			
version - Integer Editable			
run - Integer <sup>Editable</sup>			
subrun - Integer <sup>Editable</sup>			
telescope - String Editable			
year - Integer Editable			
month - Integer <sub>Editable</sub>			

#### **Buckets**

arkivum =				
navigation	*			
Admin	Location	Location Type	Cloud Provider	Bucket Name
Datapools Namespace	cern-ingest	Ingest	GCP	arkivum-cern-ingest
• Buckets	cern-export	Export	GCP	arkivum-cern-export
Saved Queries	XRootD Server	Export	S3 Compatible Storage	cern-export
Relation Kinds	Location 1	Content Archive	GCP	arkivum-cern-loca-1
≡ Reports +	Location 2	Content Archive	GCP	arkivum-cern-loca-2

#### **Buckets**

- Ingest, Storage, Export
- AWS, GCP, on-prem
  - GCP: standard, nearline, coldline, archive
  - AWS: standard, glacier, deep glacier
- Provided by Arkivum or external
- Automatic cleanup of files
- Grouped together for long-term storage
  - E.g. 1 copy GCP frequent access, 1 copy GCP deep archive, 1 copy Azure escrow

Change clean up settings X	Change clean up settings X
Ingest auto clean up •	Export retain duration (ISO 8601 format - '0' to disable) * P3D
* mandatory fields Cancel Save	* mandatory fields Cancel Save



## Ingest of Datasets

- Ingest via buckets
- Ingest via REST API
- Ingest via Web UI
- Optional use of bagit
- Optional use of archive containers (tar, zip, 7z)
- Optional ingest of metadata









**Bagit & BDbags** 

le Ingestion								
Data Pool Arkivum			~					
Small Uploads	(j)							
Drag & Drop		Upload Queue						
Drag & Drop files here		Queue length: 2						
	•	Name	Size	Progress	Status	Actions		
		SCAN-9709037.tif	1.039 MB			1. Upload	Ø Cancel	t Remove
		SCAN-9709037.pdf	3.786 MB			1 Upload	© Cancel	🛍 Remove
		Queue progress						
		* Please do not refresh your browser while using drag and drop to upload content						
						1 Upload All	⊘ Cancel A	II 🖞 Remove

## Ingest Workflow using Buckets



#### Dataset Ingest

https://vimeo.com/722918646/46a9059c74

#### **Optional Ingest Approval Workflow**

- Users can upload content, but not ingest it
- Ingest Approvers get notified of new datasets
- Approvers can review datasets and approve/reject
- Datasets ingested
- Decisions and actions in audit trail
- Notifications when ingest complete



# Solution: Part 2

Ο

0

O

0

lacksquare

0

## Agenda

- Search and Navigation
- Export of data and metadata
- Providing and updating metadata



## Search and Navigation

navigation	Home > P_C_ILCDOCrecords/ilcdoc-8710-1651491995/_0 > AIP_ILCDOCrecords/ilcdoc-8710-1651491995/_reduced > AIP_C_ILCDOCrecords/ilcdoc-8710-1651491995/_0									
🛱 Dashboard +										
尊 Admin +	<b>₿</b>	+ Q Search								
රා Files –	> C_ILCDOCrecords/ilcdoc-8737-165:	Content Type	ID Di							
File Explorer	P_C_ILCDOCrecords/ilcdoc-8710-16514		ControlsCommodityCommutingUCNoto_Eg28g4c2_72go_4ff4_877g_oo4008							
Files to Import	AIP_ILCDOCrecords/ilcdoc-8710-16	PCF	arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivum-pr							
≡ Reports +	AIP_C_ILCDOCrecords/ilcdoc-8	1000	ControlsCommodityComputingILCNote.doc_1 /arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivum-pr/ /ark							
II Rules +	AIP_C_ILCDOCrecords/ilcdoc-8	and.	metadata-ilcdoc-8710.xml /arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivum-pr/ /ark							
Annual Ingest Annual Ingest	<ul> <li>C_ILCDOCrecords/ilcdoc-8710-1651</li> <li>P. C. ILCDOCrecords/ilcdoc-9492-1652</li> </ul>	magitcreate.log /arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivu								
-	> P_C_ILCDOCrecords/ilcdoc-9492-/_0	TXT	bibdoc.txt /arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivum-pr							
	<ul> <li>P_C_ILCDOCrecords/ilcdoc-8710-1652</li> <li>P_C_ILCDOCrecords/ilcdoc-8710-16514</li> </ul>	-	dc.xml /arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivum-pr							
	P_C_ILCDOCrecords/ilcdoc-8600-1652	727	sip.json /arkivum-preserved/b6ea0e05-cc48-440f-b891-896f80eebeb9/data/objects/arkivum-pr/ /ark							
	<ul> <li>P_C_ILCDOCrecords/ilcdoc-8710-16514</li> <li>P_C_ILCDOCrecords/ilcdoc-8720-1652</li> </ul>									
	> P_C_ILCDOCrecords/ilcdoc-8737/_0 *									

#### Search and Navigation

https://vimeo.com/723627686/2153b90960

#### Hierarchical Data Structures

- Datasets can be structured using PCDM
- Structure defined in metadata
- Tree viewer shows hierarchies and links
- Breadcrumbs and other navigation aids




### Searching

- Search box and search expressions (Booleans, fuzzy matches, wildcards etc.)
- Query builder
- Filter by datapools, fields, type of entity (files, collections etc.)
- Searches only return results that user is entitled to see
- Save and re-run queries
- Easy export results of a search
- REST API or UI
- Backed by ElasticSearch allowing use of DSL or Query Strings

n	nvigation ිn Dashboard +	٩٩			
	🌣 Admin –	Query ID	Query Name	Query String	Datapools
	Datapools Namespace	62b06ca235aefe005428	Crab Nebula 1 Jan 2019	(metadata.magic_telescopes.year:"2019")AND(	Default, MAGIC1,
	Buckets	62b07389334e7d04ad7b	Perseus 31 Dec 2019	(Perseus)AND (metadata.magic_telescopes.ye	Default, MAGIC1,
	Saved Queries				

#### Export

- Files and/or Metadata
- With/without bagit
- Choice of buckets
- Webhooks and Notifications
- Request/Approve workflow
- Export reports and audit trail

Bulk Export Request	×
Export Path (optional)	
dataset	
Location: CERN Export Bucket  Export Metadata Type: Content and Metadata  Export Format: BagIt  Include URI in metadata export Include technical metadata in metadata export	

Dashboard +	÷				
향 Admin +	Export ID	ID	Export Status	Export Requested	Last Modified Time
□ Files + ■ Reports _	62b0741a334e7d04ad7b57c2		Success	2022-06-20 13:20:26	2022-06-20 13:24:54
Audit Trail	62b05871334e7d04ad7b4d00	C_Calibratedvl	Success	2022-06-20 11:22:25	2022-06-20 11:24:12
Hold Report	62b0575c334e7d04ad7b44fe	C_Calibratedv1	Success	2022-06-20 11:17:48	2022-06-20 11:22:51
Unhold Report Processing Report	62ac8a02334e7d04ad7b2a5b	C_Calibratedv1	Success	2022-06-17 14:04:50	2022-06-17 14:09:09
Ingest Report	62ac7ad2334e7d04ad7b1503	C_Calibratedvl	Success	2022-06-17 13:00:02	2022-06-17 13:04:29
Preservation Report	62a98e70d8e4601ec217265d	2022-04-08_20-29_Calibra	Success	2022-06-15 07:46:56	2022-06-15 07:49:17
Deletion Report	62a98a62a3d6fc37d9ff9b9f	C_Calibratedvl	Success	2022-06-15 07:29:38	2022-06-15 07:36:17

#### Export

https://vimeo.com/723623926/b82a560f3d

#### Export Approval Workflow

- Optional workflow
- Allows control over costs and access
- Can be combined with integrations, e.g. xrootd and InvenioRDM

Dear Matthew Addis,

An export request 62b08724434bc871120b47e8 has been submitted by Bob Researcher and requires approval to proceed.

Please click on this link to review the request and approve or reject it, or log on to the Arkivum system to review your notifications.

This request must be actioned by 2022-07-20 14:41:40 or it will be automatically rejected and returned to Bob Researcher .

Number of files in request: 1

Export path: 6948

Export bucket: cern-export

Exported files P\_C\_ILCDOCrecords/ilcdoc-6948-1652866857/

Regards

Arkivum Support

You have been sent this mail because you have an account in Arkivum.com. Please contact us to change your email preferences.

#### Providing and Updating Metadata

- Metadata can be provided as XML, json or CSV
- Metadata fields are indexed and searchable
- Metadata can be extracted from domain specific metadata files
- Metadata fields can be mapped, e.g. to DublinCore or DataCite
- Rules can be applied to enforce metadata constraints
- Metadata can be provided when data is ingested, or added afterwards
- Metadata can be used to define the structure of datasets
- Metadata can be used to link datasets or files to each other

Data Files
Metadata Files (native format)
ark-manifest.json (metadata mapping)

## Dataset

#### Metadata Example





JSON Raw Data Headers	
Save Copy Collapse All Expand All	7 Filter JSON
internal:	"SCAN-9709037"
<b>v</b> 1:	
report_number:	"UCRL-8417"
creation_date:	"1990-01-27T00:00:00"
▼ imprint:	
date:	"1958"
▼ oai:	
▼ indicator:	
0:	"cerncds:FULLTEXT"
1:	"cerncds:SCAN:FULLTEXT"
2:	"cerncds:SCAN"
value:	"oai:cds.cern.ch:104881"
owner:	"PUBLIC"
▼ subject:	
source:	"SzGeCERN"
term:	"Mathematical Physics and Mathematics"
<pre>~ cataloguer_info:</pre>	
library:	"CER01"
hour:	"1801"
modification_date:	"20091110"
creation_date:	"19900127"
agency_code:	"SZGECERN"
humber_ol_citations:	0
<ul> <li>title:</li> </ul>	"Notos on statistics for physicists"
entre:	Notes on statistics for physicists
• persistent_identifiers_keys.	"recid"
1.	"system number"
2.	"system_control_number"
3:	"oai"
v system number:	
value:	"000008580CER"
location:	null
✓ files:	
<b>v</b> 0:	
comment:	null
status:	
▼ magic:	
0:	"PDF document, version 1.3"
1:	"application/pdf; charset=binary"
2:	"PDF document, version 1.3"
3:	"application/pdf; charset=binary"
4:	"application/pdf"
description:	"Access to fulltext document"
✓ url:	"http://cds.cern.ch/record/104881/files/SCAN-9709037.pdf"
eformat:	".pdf"







#### End Result

arkivum	=		🔎 Matthew Addis 🗸
navigation	Home > C_test_data/CERN/104881/ > 104881 > 104881-or	ginal-data	
🖬 Dashboard ۞ Admin	+ + B	+ C Search	
🗅 Files	- ~ C_test_data/CERN/104881/ ··· Con	tent Type ID	Display Name
<ul> <li>File Explorer</li> <li>Files to Import</li> </ul>	<ul> <li>104881 ***</li> <li>104881-metadata ***</li> </ul>	SCAN-9709037.tif /deep_archive/104881/original_scans	/scan-9709037.tif /deep_archive/104 ····
■ Reports	+ 104881-access-data •••• •••		Related Entities View / Edit metadata
C Rules	+ O_test_data/CERN/104881/ •••		Edit Display Name
♀ Manual Ingest ↓ Notifications	<ul> <li>C_ilcdoc-6949/ •••</li> <li>C_Calibratedv1CrabNebula</li> <li>C_test_data/CERN/12351/ •••</li> </ul>	Dublin Core	Edit 🖉 Bulk export
	Related entities for 104881-access-data	X Title Notes on statistic	s for physicists Add To New Object
	isDerivedFrom	Subject	
	104881-access-data 104881-access-data ↓ 104881-original-data 104881-original-data	•••• Mathematical Phy Description	rsics and Mathematics

#### Metadata Updates

- Add metadata to existing datasets using a metadata file, e.g. ark-manifest.json
- Export metadata file, update file, and re-ingest updated metadata file
- GET/POST/PUT metadata using REST API
- Edit metadata in the UI
- Extract metadata fields from native metadata files already in the archive
- Trigger external scripts/applications that extract/generate metadata using webhooks, Cloud Events, Arkivum REST API



#### Round-trip, dataset migration, exit strategy

- Export datasets into a bucket
- Bagit for dataset integrity checks
- Original data files, metadata files
- Normalised preservation versions of files (if created)
- ark-manifest.json description



#### Round trip dataset

https://vimeo.com/723624529/0c4b8419e1

#### **Retention Schedules and Deletion**

- Dataset deletion follows a request/approve workflow
- Deletion needs one approvals
- Retention schedules can create deletion events

Add new rule	×
ID *	Duration *
EPSRC 10 year	Years 🗘 Days 🗘 Hours 🗘 Minute 🗘
Туре	Action *
Collection ~	Delete ~
Default rule	
Reason	
Delete if dataset not access for 10 years	
	Air



# Solution: Part 3

 $\bigcirc$ 

0

O

0

lacksquare

0

#### Agenda

- File Format Normalisation
- Long-term storage: replication, tiers, fixity checks
- Access to archived data using xrootd
- Integration with InvenioRDM
- Co-locating applications with the Arkivum solution for pre and post processing of archived data
- Summary





#### **Preservation Formats and Access Formats**



Media type	File formats	Preservation format(s)	Access
Audio	AC3, AIFF, MP3, WAV, WMA	WAVE (LPCM)	MP3
Email	PST	MBOX/EML	PDF
Email	MSG	EML	PDF
Office docs and presentations	DOC, WPD, RTF, DOCX, PPTX, PPT	PDF/A	PDF/A
Plain text	TXT	Original format	Original format
Portable Document Format	PDF	PDF/A	Original format
Raster images	BMP, GIF, JPG, JP2*, PCT, PNG*, PSD, TIFF, TGA	TIFF	JPEG
	3FR, ARW, CR2, CRW, DCR, DNG, ERF, KDC,		
Raw camera files/Digital Negative format	MRW, NEF, ORF, PEF, RAF, RAW, X3F	TIFF	JPEG
Spreadsheets	XLS, XLSX	Original format	Original format
Vector images	AI, EPS, SVG	SVG	PDF
	AVI, FLV, MOV, MPEG-1, MPEG-2, MPEG-4, SWF,		
Video	WMV	FFV1/LPCM in MKV	MP4



#### Preferred Formats

Preferred Formats	Format Versions	Format Specifications
ASCII Text	7 bit	ISO/IEC 646:1991 Information technology ISO 7-bit coded character set for information interchange: (http://www.iso.org/iso/catalogue_detail.htm?csnumber=4777 🗷)
Union do Taut	UTF-8	RTF 3629: UTF-8, A Transformation Format of ISO 10646: ( http://tools.ietf.org/html/rfc3629 🗹)
Unicode Text	UTF-16	RFC 2781 UTF-16: An Encoding of ISO 10646: ( http://www.ietf.org/rfc/rfc2781.txt 🕜)
OpenDocument Text Format (ODF)	OpenDocument 1.0	ISO/IEC 26300:2006 Information technology OpenDocument Format for Office Applications (OpenDocument) v1.0: ( http://www.iso.org/iso/iso_catalogue/ catalogue_tc/catalogue_detail.htm?csnumber=43485 (27)
PDF/A-1	PDF/A-1	ISO 19005-1:2005 Document management Electronic document file format for long-term preservation Part 1: Use of PDF 1.4 (PDF/A-1): (http://www.iso.org /iso/catalogue_detail?csnumber=38920 ⑦)
PDF/A-2	PDF/A-2	ISO 19005-2:2011 Document management Electronic document file format for long-term preservation Part 2: Use of ISO 32000-1 (PDF/A-2): (http://www.iso.org/iso/home/store /catalogue_tc/catalogue_detail.htm?csnumber=50655 ♂)
Acceptable Formats		
Acceptable Formats	Format Versions	Format Specifications

Acceptable Formats	Format Versions	Format Specifications
PDF	PDF 1.7	ISO 32000-1:2008 Document management Portable document format Part 1: PDF 1.7: (http://www.iso.org/iso/catalogue_detail.htm?csnumber=51502 7
	PDF 1.0-1.6	Adobe® Portable Document Format Version 1.6: http://www.adobe.com/devnet /pdf/pdf_reference_archive.html 🗹
Microsoft Word (DOCX) Office Open XML	OOXML Microsoft Word for Windows, version 2007-2010	[MS-Ol29500]: Office Implementation Information for ISO/IEC 29500 Standards Support: ( http://msdn.microsoft.com/en-us/library/ee908652%28v=office.12%29 🗹)
Microsoft Word 97 Binary Document Format (DOC)	8.0	[MS-DOC]: Word (.doc) Binary File Format: ( http://msdn.microsoft.com/en-us/library /cc313153%28v=office.12%29.aspx 🖒

#### File Format Normalisation

navigation	Home > P_C_pll.zip > AlP_pll.zip_reduced > AlP_C_pll.	zip	
🛱 Dashboard +			
尊 Admin +	<b>路</b> …		+ Q Search
ා Files –		TXT	/arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
File Explorer	Name ↓↑ Modified Date ↓↑	NLS.	planetary_data.xls /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
Files to Import	<ul> <li>P_C_pll.zip ····</li> <li>AIP_pll.zip_complete ····</li> </ul>	222	11_Jupiter_FC-97ee4a17-0334-48a8-8805-49ar if6eda05.pdf /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7 45ba59/data/objectr/arkivum-pr
E Reports +	AIP_pll.zip_reduced ····	PDF	11_Jupi er_FC.pdf /arkivum_reserved/67643272_ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
다 Rules + 우 Manual Ingest	AIP_C_pll.zip ••••	PDF	62211main_Jupiter_Lithograph-e33284b8-0e05-4f15-b9e5-e06e90205558 /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
	C_Calibratedv1LSI+612019_08_2	905	62211main_Jupiter_Lithograph.pdf /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d997.551.59/data/objects/arkivum-pr
	C_ilcdoc-6949/ ***		Jupiter-6a29682a-12a3-4170-8f21-38c4a207db9f.pdf /arkivum-preserved/67643273-ddcb-49b1-b3d6-5/9a7d15ba59/data/objects/arkivum-pr
	<ul> <li>C_CalibratedvlCrabNebula2019</li> <li>C_test_data/CERN/12351/ ···</li> </ul>	Por	Jupiter-99f2df86-e349-4d99-91f8-6936d330563b.pdf /arkivum-preserved/67643273-ddcb-49b1-b3d6_5d9a7d15ba59/datr.robjects/arkivum-pr
	C_test_data/CERN/HiggsToBBNtupleProducerTool/		Jupiter.docx /arkivum-presenra/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
	<ul> <li>P_C_ILCDOCrecords/ilcdoc-9492/_0 ···</li> <li>P_C_ILCDOCrecords/ilcdoc-8737-1652431608/_0 ··</li> </ul>	FILE	Jupiter.msg /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
	P_C_ILCDOCrecords/ilcdoc-8710-1651491995/_0 ···	PPT	Jupiter.ppt /arkivum-preserve/67643273-ddcb-49b1-b3d6-5d9a7ct1Eba59/data/objects/arkivum-pr
	<ul> <li>P_C_ILCDOCrecords/ilcdoc-9492-/_0 ***</li> </ul>	FILE	Jupiter-7e61ab37-8371-4469-bd1f-329b1fd57595.wav /arkivum-preserved/67643273-ddcb-49b1-b3d6_549a7d15ba59/dcta/objects/arkivum-pr
	<ul> <li>P_C_ILCDOCrecords/ilcdoc-8710-1652431598/_0 ···</li> <li>P_C_ILCDOCrecords/ilcdoc-8710-1651491991/ 0 ···</li> </ul>	MEL	Jupiter.mp3 /arkivum-preserve/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
	P_C_ILCDOCrecords/ilcdoc-8600-1652431552/_0 ··		743610main_pia17007-full_full.jpg /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/objects/arkivum-pr
	> P_C_ILCDOCrecords/ilcdoc-8710-1651491996/_0 ···		PIA19048_ip.jpg /arkivum-preserved/67643273-ddcb-49b1-b3d6-5d9a7d15ba59/data/obiects/arkivum-pr





#### **Email Normalization Workflow**



#### Long-term Archival Storage in Buckets



## Storage Options and Fixity Checks

- Multiple copies of each file
- Multiple cloud providers
- Combine different tiers of storage
- Multiple checksums for each file
- Data integrity checks

√ ↓↑ Event

Fixity

Fixity

Fixity

Fixity

Fixity

Fixity

Fixity

**Event Type Code** 

OBJECT\_FIXITY\_CHECK\_SUCCESS

OBJECT\_FIXITY\_CHECK\_SUCCESS

OBJECT\_FIXITY\_CHECK\_SUCCESS

OBJECT\_FIXITY\_CHECK\_SUCCESS

OBJECT\_FIXITY\_CHECK\_SUCCESS

OBJECT\_FIXITY\_CHECK\_SUCCESS

OBJECT\_FIXITY\_CHECK\_SUCCESS

• When data first stored

ID

62b1b47106d6c474d83b968d

62b1b47106d6c474d83b9689

62b1b47106d6c474d83b969b

62b1b47106d6c474d83b9689

62b1b47106d6c474d83b969b

62b1b47306d6c474d83b96ae

62b1b47106d6c474d83b9695

Action

Fixity performed

Time Stamp

2022-06-21 12:09:31

2022-06-21 12:09:31

2022-06-21 12:09:31

2022-06-21 12:09:31

2022-06-21 12:09:31

2022-06-21 12:09:31

2022-06-21 12:09:31

• Periodic checks

Туре:		F	
Errors Present:		No	
MD5: af66a3bf39d1415	8d21c9ce9ab9280a6		
SHA-256: c82eeeb55c0d63	34700e7913bb02e6efeb313f96	Sef60f9f335f1afc3eed213df4	
SHA-512: d53efdfe8771dab	9b4d2c028ee596fa6f61cc4aa	a297606aacd2c8d3fe3554d262d19d1e993ddd8bb1dfaa3e5fb40afbd4528	a42852c49d05d8be96fb95a0a9
Adler-32: b8a6c672			
Metadata processing	g Caching Indexing I	Integrity check Encryption Replication Fixity	
Location: Location	on 2		
Status			
Success			
Success Last Checked Time 2022-06-20 11:29:19			
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined			
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined			
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Location	on 1		
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Location Status Success	on 1		
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Locati Status Success Last Checked Time 2022-06-20 11:29:19	on 1		
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Locatio Status Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined	on 1 Fixity Notifications		
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Locati Status Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined	on 1 Fixity Notifications Datapool	Last Fixity Time Check	Status
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Locatio Status Success Last Checked Time 2022-06-20 11:29:16 Number of Retries undefined	on 1 Fixity Notifications Datapool Arkivum	Last Fixity Time Check 2022-06-20 09:16:55	Status Clear
Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined Location: Locatio Status Success Last Checked Time 2022-06-20 11:29:19 Number of Retries undefined	on 1 Fixity Notifications Datapool Arkivum Arkivum Preserved	Last Fixity Time Check 2022-06-20 09:16:55 2022-06-21 12:09:31	Status Clear Clear

#### Datapools and Storage

https://vimeo.com/722930836/9659bb1784

# Publication, Discovery, Analysis and Reuse of Research Data

 $\cap$ 

 $\square$ 

0

# 

My dashboard

2019

Export

+-

#### O Preview

You are previewing a new record that has not yet been published.

arkivum =			< Back to edit			
navigation	Home > C_test	_data/CERN/HiggsToBBNtupleProducerTool/ > O_test_data/CERN/H		Versions		
<ul> <li>☐ Dashboard +</li> <li>☆ Admin +</li> </ul>	tana / 2202000,000 0000000000000000000000000		Published 2019   Version v1	Preview     Only published versions are displayed.		
🗅 Files 🗕	Content Type	ID	for Hbb tagging ML studies			
File Explorer Files to Import		ntuple_merged_0.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	HiggsToBBNTuple_HiggsToBB_QCD_RunII_13TeV_MC	Version v1 2		
$\equiv$ Reports +	EVE	ntuple_merged_1.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTooI/Test/ntuple_merged	Duarte, Javier	Details		
🕮 Rules +	-	ntuple_merged_2.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	Citation Style APA -	Resource type		
Manual Ingest     Notifications	-	ntuple_merged_3.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	Duarte, J. (2019). Sample with jet, track and secondary vertex properties for Hbb tagging	Dataset		
	-	ntuple_merged_4.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	ML studies HiggsToBBNTuple_HiggsToBB_QCD_RunII_13TeV_MC [Data set]. CERN Open Data Portal.	Publisher CERN Open Data Portal		
		ntuple_merged_5.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged		Rights		
	RE	ntuple_merged_6.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	Description The dataset consists of particle jets extracted from simulated proton-proton collision events at a center-of-mass energy of	Oreative Commons Zero v1.0		
	RE	ntuple_merged_7.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Test/ntuple_merged	13 TeV generated with Pythia 8. It has been produced for developing machine-learning algorithms to differentiate jets originating from a Higgs boson decaying to a bottom quark-antiquark pair (Hbb) from quark or gluon jets originating from	Universal		
	122	<b>ntuple_merged_8.h5</b> /arkivum/higgs1/HiggsToBBNtupleProducerTooI/Test/ntuple_merged	quantum chromodynamic (QCD) multijet production. The reconstructed jets are clustered using the anti-kT algorithm with R=0.8 from particle flow (PF) candidates (AK8 jets).	Export		
	77	files.txt /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/files.txt	The standard L1+L2+L3+residual jet energy corrections are applied to the jets and pileup contamination is mitigated using the charged hadron subtraction (CHS) algorithm. Features of the AK8 jets with transverse momentum $pT > 200$ GeV and pseudoranidity $ p  < 2.4$ are provided. Selected features of inclusive (both charged and pseudoranidity $ p  < 2.4$ are provided.	JSON - Expo		
	FRE	ntuple_merged_10.h5 /arkivum/higgs1/HiggsToBBNtupleProducerToo1/Train/ntuple_merged	0.95 GeV associated to the AK8 jet are provided. Additional features of charged PF candidates (formed primarily by a charged particle track) with $pT > 0.95$ GeV associated to the AK8 jet are also provided. Finally, additional features of			
	-	ntuple_merged_11.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/ntuple_merged.	I_11.h5 /arkivum/higgs1/Hi ····			
		ntuple_merged_12.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/ntuple_merged.	I_12.h5 /arkivum/higgs1/Hi ····			
		ntuple_merged_13.h5 /arkivum/higgs1/HiggsToBBNtupleProducerTool/Train/ntuple_merged.	I_13.h5 /arkivum/higgs1/Hi ····			
		Entries per page:	: 15 ∨ < 1 2 3 () 7 > ≫			

#### InvenioRDM Workflow



#### InvenioRDM Workflow

https://vimeo.com/723439662/260a70b326

#### Xrootd Workflow

- Xrootd server hosted in GCP
- Transfer of datasets on-demand





#### Xrootd Workflow

https://vimeo.com/723493247/212b1f4de5

#### **Cloud-hosted Processing of Archived Research Datasets**

- GCP hosted applications
  - Cloud Functions (serverless code)
  - Cloud Run (containers)
  - GKE (Kubernetes)
  - Compute Engine (VMs)
- Arkivum Webhooks
  - Ingest, preservation, export
- Arkivum REST API
  - Including search, get metadata, export files
- Xrootd server
  - Easy integration of scientific apps



#### Webhooks

- Can be set per job (ingest, preservation, export)
- Default webhooks for a Datapool (ingest, preservation)
- Includes API tokens to allow calls to secured endpoints

Datapool Name	Datapool Path	Encrypted	Preservation Ena	AtoM Enabled	Location Set	Metadata Names	Webhook URL	
Arkivum	/arkivum	true	false	false	quickaccess	technical,http://www.arki		
Arkivum Preserved	/arkivum-preserved	true	true	false	quickaccess	technical,http://www.arki		
Deep Archive	/deep_archive	false	false	false	deeparchive	technical,http://www.arki		
Default		true	false	false	quickaccess	technical,http://www.arki		
Frequent Access	/frequent_access	Webbooks			×	technical,http://www.arki		
ILCDOCS	/ilcdocs					technical,http://www.arki		
InvenioRDM	/invenio-rdm	Ingest Webhook URL	.archiver.arkivum.net:8080/cre	ate_record_from_ingest		dataciteTypes,dataciteD	http://invenio-rdm-1.arc	
MAGIC1	/magic1					technical,magic_telesco		
				D	hiscard Apply			



≡	Google Cloud	cern-dedicated 👻	Q Search Pro	ducts, resources,	docs (/)		v) ii D	) <b>.</b> ()	: (
٨	Kubernetes Engine	Workloads	C REFRESH	+ DEPLOY	DELETE	6	G OPERATIONS	r 🖻 HELP A	SSISTAN
<b>(</b> ])	Clusters	Cluster	▼ Na	mespace	- RESE	T SAVE			
- 54	Workloads	Wedde de ere deel	Workloads are deployable units of computing that can be created and managed in a cluster.						
A	Services & Ingress	cluster.							
	Applications	OVERVIEW	COST OPTIMIZATI	ON					
⊞	Secrets & ConfigMaps	Ţ Filter (Is s	system object : False (	Filter workloads	3			×ø	III
	Storage	Name 1	`	Status	Туре	Pods	Namespace	Cluster	
:=	Object Browser	reana-ca	che	📀 ОК	Deployment	1/1	default	reana	
	object browser	reana-db	1	🖉 ОК	Deplo				
A	Migrate to containers	reana-m	essage-broker	🛇 ок	State	Cal	na		
0	Backup for GKE	reana-re	source-quota-update	🕗 ок	Cron				
۲	Config Management	reana-ru 47cf-b96	n-batch-c0c4997f-612d 2-eaac1d1c4c24	с- 🕑 ОК	Job				
<b></b>	Protoct	reana-se	rver	📀 ОК	Deplo		Yo	our workf	lows
۲	FIOLECT NEW	reana-tra	efik	📀 ОК	Deplo				
		reana-ui		📀 ОК	Deplo				
		reana-we	orkflow-controller	📀 ОК	Deplo		S	earch	

Your workflows	€ Refreshed at 16:06:06 U	тс
Search	Q	•
Status 🝷	Sort by 👻	
root6-roofit-yadage-kubernetes #1 Finished 11 minutes ago	<b>finished</b> in 1 min 21 sec step 2/2	
<ul> <li>root6-roofit-snakemake- kubernetes #1</li> <li>Finished 14 minutes ago</li> </ul>	<b>finished</b> in 32 seconds step 2/2	
♥ root6-roofit-serial-kubernetes #1 Finished 15 minutes ago	<b>finished</b> in 19 seconds step 2/2	
root6-roofit-cwl-kubernetes #1 Finished 16 minutes ago	<b>finished</b> in 1 min 17 sec step 2/2	



#### Long-Term Digital Preservation of Research Data



#### Self-Assessment using the CTS+FAIR Capability Maturity Model

Area	Short Requirement Name	FAIR Principle ( <u>Nature</u> )	RDA Indicator	FAIRsFAIR Metric	Evidence Links	
Digital Object Management	13. Data discovery and identification	Discovery F2. data are described with rich metadata (defined by R1 below)	RDA-F2-01M Rich metadata is provided to allow discovery (Essential)	FsF-F2-01M Metadata includes descriptive core elements (creator, title, data identifier, publisher, publication date, summary and keywords) to support data findability.	Metadata support in the Arkivum solution includes DublinCore, DataCite and custom fields. Metadata can be ingested and stored in its native format (json and XML), e.g. domain specific schemas. Metadata fields in domain specific metadata can be mapped to DC or DataCite and then made searchable. DataCite metadata can be exported to invenioRDM.	
		<b>Discovery</b> F4. (meta)data are registered or indexed in a searchable resource	RDA-F4-01M Metadata is offered in such a way that it can be harvested and indexed (Essential)	FsF-F4-01M Metadata is offered in such a way that it can be retrieved by machines.	InvenioRDM supports publication of metadata landing pages. InvenioRDM supports metadata harvesting using OAI-PMH. https://inveniordm.docs.cern.ch/reference/oai_pmh/ InvenioRDM supports PID services, e.g. DOI registration using DataCite. https://inveniordm.docs.cern.ch/customize/dois/ Metadata In the Arkivum solution can be searched and retrieved though an API (authenticated users only)	
		<b>Identification</b> F1. (meta)data are assigned a globally unique and persistent identifier	RDA-F1-01M Metadata is identified by a persistent identifier (Essential) RDA-F1-01D Data is identified by a persistent identifier (Essential) RDA-F1-02M Metadata is identified by a globally unique identifier (Essential) RDA-F1-02D Data is identified by a globally unique identifier (Essential)	FsF-F1-02D Data is assigned a persistent identifier. FsF-F1-01D Data is assigned a globally unique identifier.	All entities in the Arkivum solution (PCDM Collections, objects, Files) are assigned globally unique UUIDs. Users can provide one or more of their own identifiers as metadata, e.g. PIDs such as DOIs or handles. The Arkivum metadata schema supports identifier type and indentifier value. URLs are possible to entities in the Arkivum solution (authentication is required for access). Records in InvenioRDM can have one or more PIDs, including those registered with external PID systems or as Alternative Identifiers. Identifier schemes include: ARK, arXiv, Bibcode, DOI, EAN13, EISSN, Handle, IGSN, ISBN, ISSN, ISTC, LISSN, LSID, PubMed ID, PURL, UPC, URL, URN, W3ID. https://inveniordm.docs.cern.ch/reference/metadata/#external-pids https://inveniordm.docs.cern.ch/reference/metadata/#alternate-identifiers-0-n	
		Identification F3. metadata clearly and explicitly include the identifier of the data it describes	RDA-F3-01M Metadata includes the identifier for the data (Essential)	FsF-F3-01M Metadata includes the identifier of the data it describes.	Metadata in InvenioRDM includes identifiers (see above). Likewise, all records in the Arkivum solution have identfilers. Records in InvenioRDM can include files. The InvenioRDM roadmap includes support for referencing data in external systems. https://inveniosoftware.org/products/rdm/roadmap/	
## Further Materials and Summary

0

0

 $\bigcirc$ 

•

0

## **Further Materials**

- User Guide
- API Guide
- Cookbooks and examples for specific activities
- Demos and pilot projects
- Onboarding process, including detailed training
- Zendesk support portal
- Mappings to NDSA preservation levels, DPC RAM, CTS, FAIR metrics
- iPRES 2021 paper and iPRES 2022 panel
- ARCHIVER deliverables

## Summary

- The Arkivum solution enables LTDP of research data:
  - Ingest, Preserve, Safeguard, Search, Navigate, Download, Export, Publish
- Support for reuse of scientific datasets:
  - xrootd, InvenioRDM, webhooks, reana, snakemake, scripts, k8s, VMs
- Flexible deployment
  - GCP, AWS, Open Stack
- Authentication and Authorisation
  - eduGAIN, request/approve workflows, control over actions and data
- Support for good practice, assessment and certification
  - NDSA preservation levels, DPC RAM, CoreTrustSeal, FAIR principles and metrics
- Available for use today
  - Listed on EOSC Marketplace, contact Arkivum or Google













## QUESTIONS?



matthew.addis@arkivum.com



orcid.org/0000-0002-3837-2526



www.arkivum.com

www.arkivum.com

hello@arkivum.com