

DLCM: Comment être innovant en gestion des données de recherche

Pierre-Yves BURGI

Directeur adjoint DiSTIC et DLCM Project Director

Hugues Cazeaux

Responsable du pôle eResearch, DiSTIC et DLCM CTO



Agenda

1. What's the problem ?
2. Design methodology: « under the hood »
3. What next ?
4. Questions

1. What's the problem?

Four (good) reasons for practicing data governance (in research):

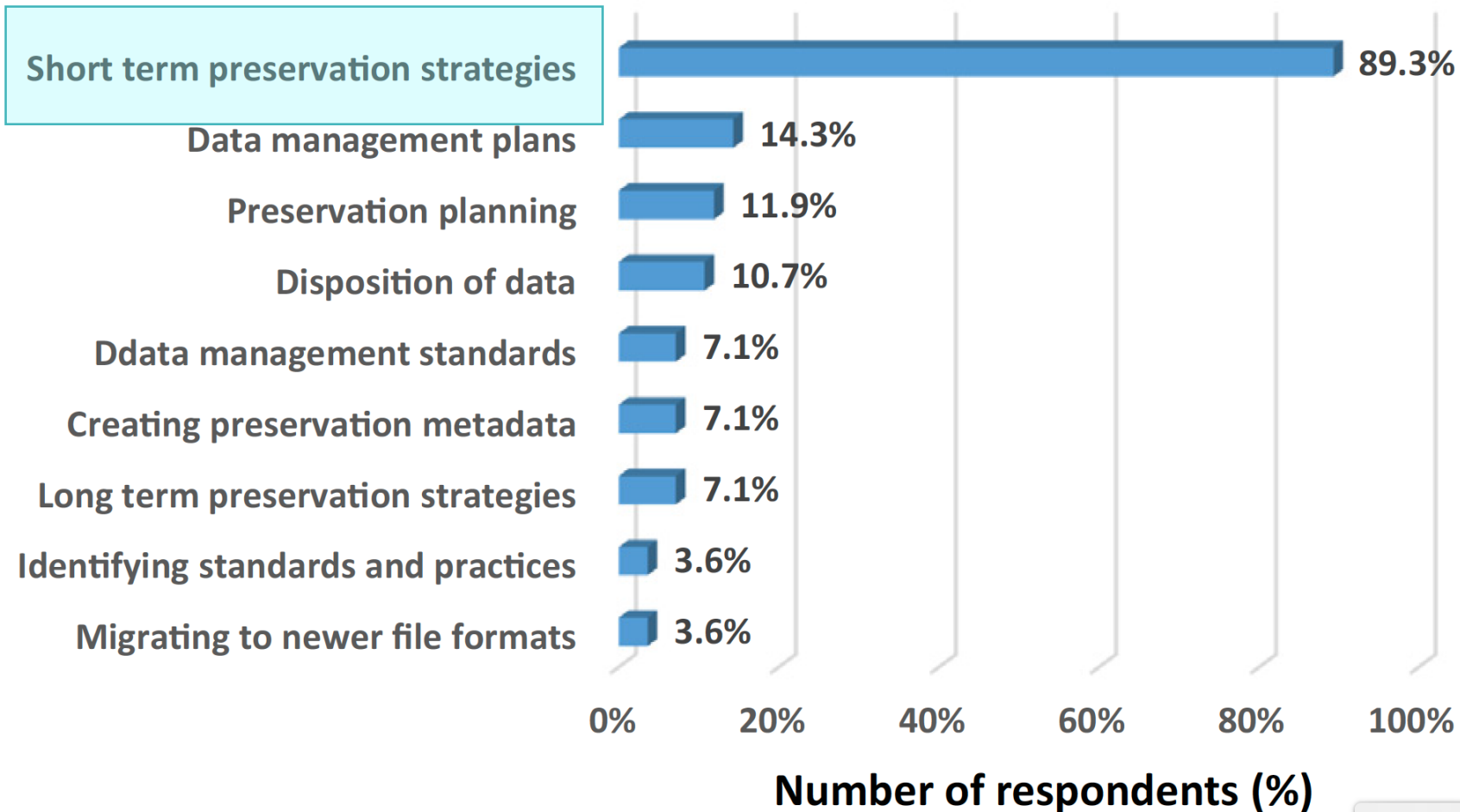
- 1) Avoid repeating experiments and thus reduces the costs
- 2) Contributes to global development by the study of phenomena other than those initially thought of
- 3) Guards against research fraud and improve reproducibility/replicability of results
- 4) Can be used by students in their learning activities

“**but**, researchers face challenges in organizing and sharing their data”

Factors affecting researchers' RDM (n = 84).

Challenge

- Lack of policy frameworks
- Lack of incentives
- Lack of skills to create metadata
- Lack of curation skills and training
- Lack of storage network infrastructure
- Lack of guidance and support
- Lack of curation tools and software
- Finding data produced by others
- Lack of support from the university
- Lack of incentives to share data
- Most data is not trustworthy
- Failure of re-users to cite my data
- Prohibitive institutional policies
- Obsolescence of technologies



Follow good practices !



Make Data FAIR

DATA

FAIR DATA PRINCIPLES

AH!



FINDABLE

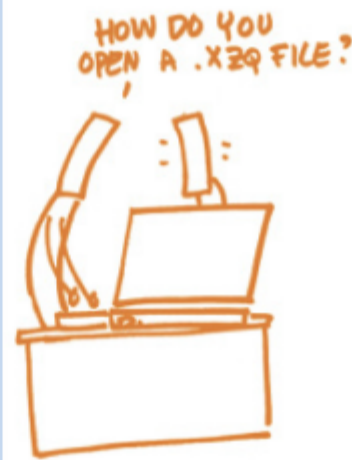
- ✓ persistent identifier
- ✓ enriched metadata
- ✓ Searchable and findable online



ACCESIBLE

- ✓ retrievable using standard communication protocols
- ✓ possibility to define access rights

→ in a repository that ensures long-term preservation



INTEROPERABLE

- ✓ standard formats
- ✓ controlled vocabulary to describe data

→ data will be compatible and combinable with others

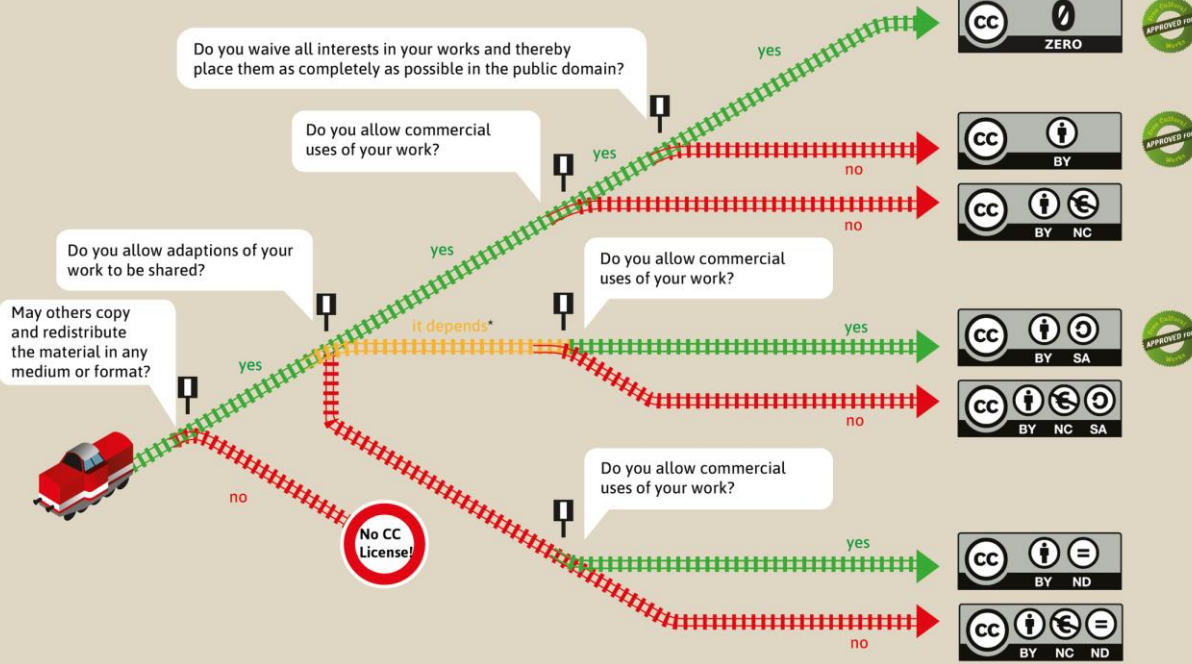


REUSABLE

- ✓ well-described & documented (e.g. in a README file)
- ✓ clear conditions to cite and reuse data (e.g. CC licenses)

→ allow data to be correctly interpreted and reused

"CHOO-CHOO-CHOOSE YOUR LICENSE!"



*it depends = yes, as long as others share alike (meaning: They use the same license as you.)



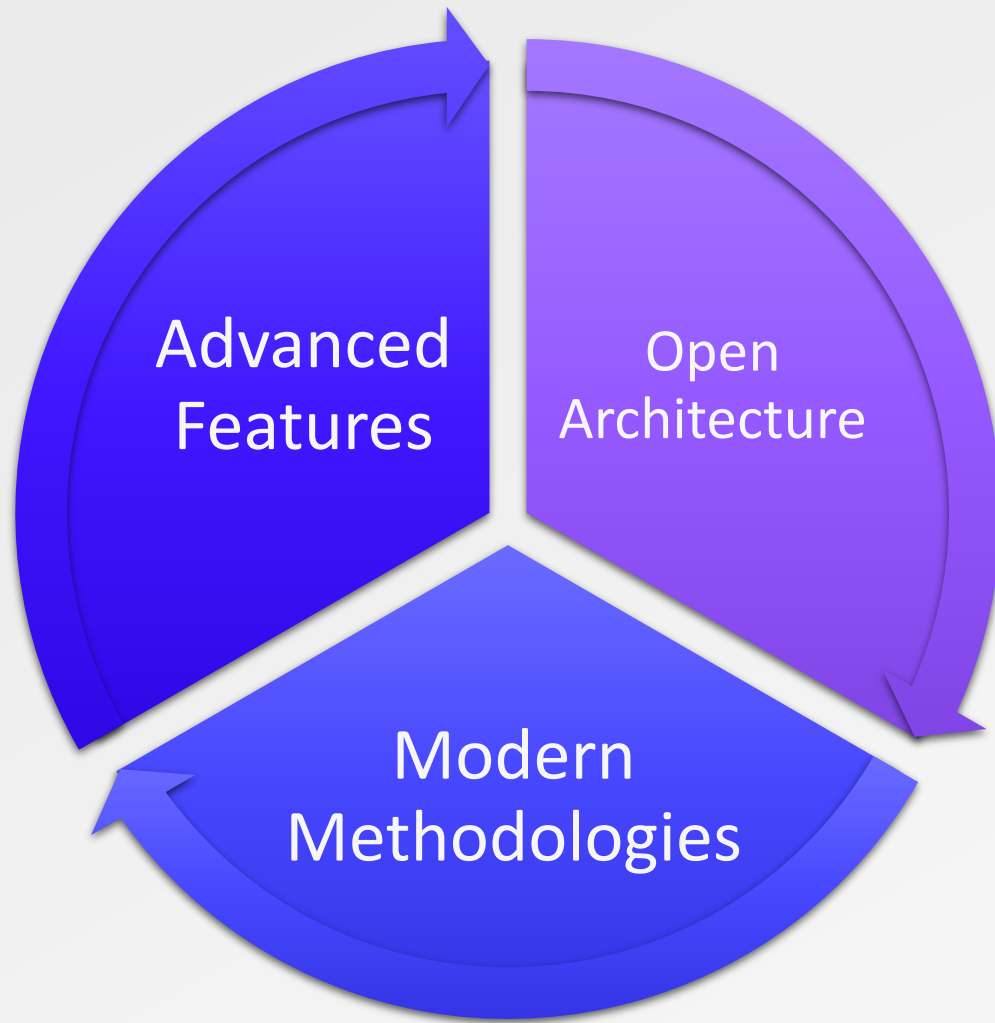
This Graphic "Choo-Choo-Choose your license!" is based on the work "Welche CC-Lizenz ist die richtige für mich?" by Barbara Klute und Jöran Muuß-Merholz für wb-web* unter CC BY-SA 3.0**. The English version is a translation and enhancement by Jöran Muuß-Merholz under the same licence.
 * <http://www.wb-web.de> | ** <https://creativecommons.org/licenses/by-sa/3.0/deed.de>

partage	DP	CC 0	CC BY
modifier		CC BY SA	
commercial		CC BY NC	CC BY SA
partage		CC BY NC	CC BY SA
modifier		CC BY NC SA	
partage		CC BY ND	
		CC BY NC ND	
tous droits réservés			

<https://creativecommons.org/choose/>

2. Design Approach: « *under the hood* »

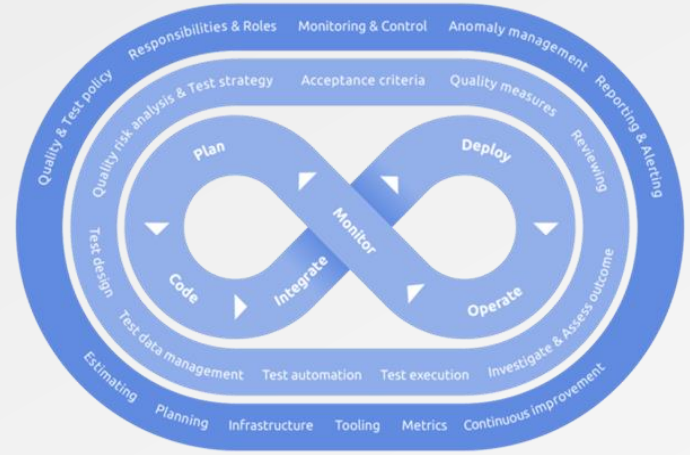
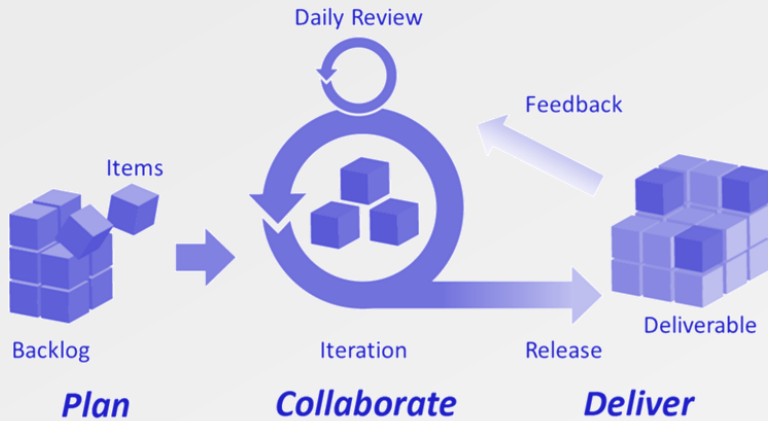
DLCM Project



2.1 Modern Methodologies

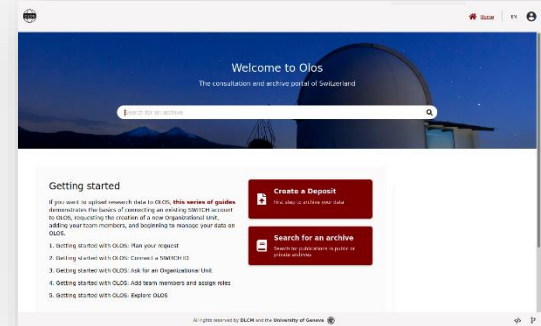
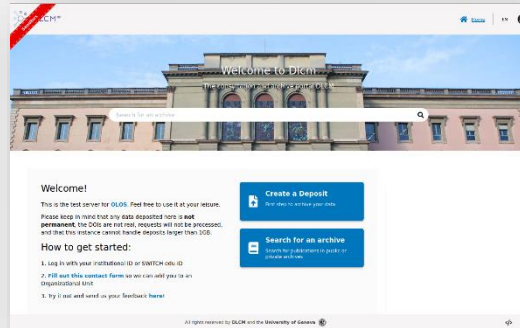
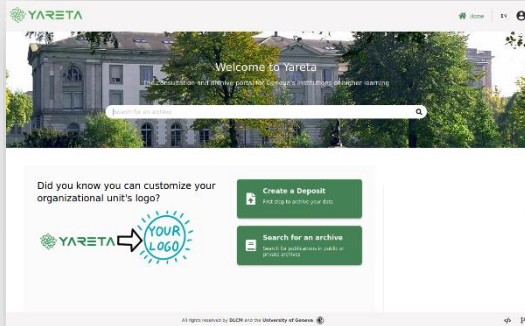
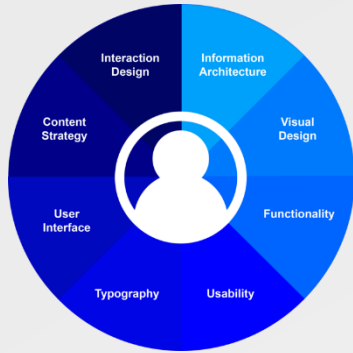
DLCM Project

Software Agile Development

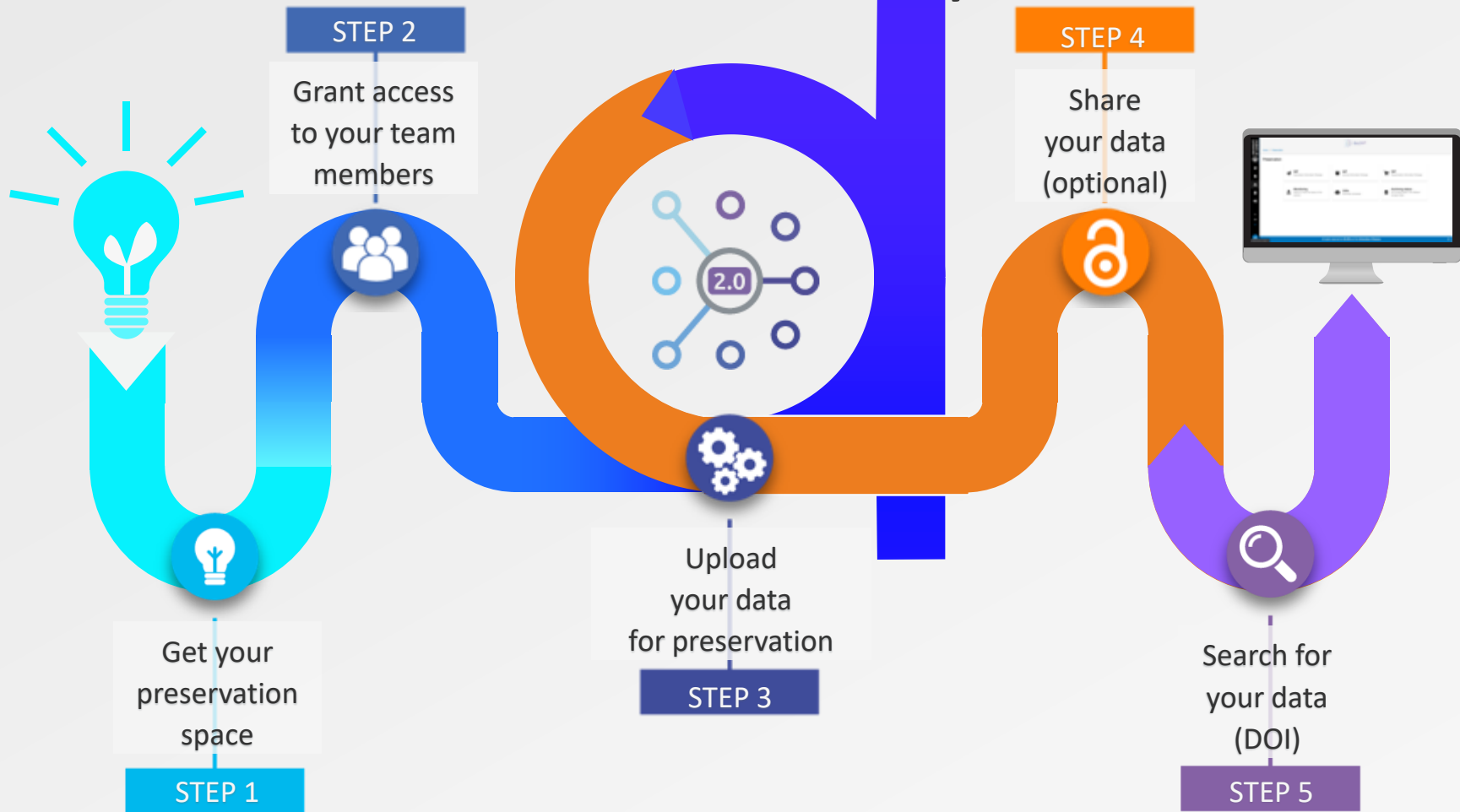


DevOps

UX & UI



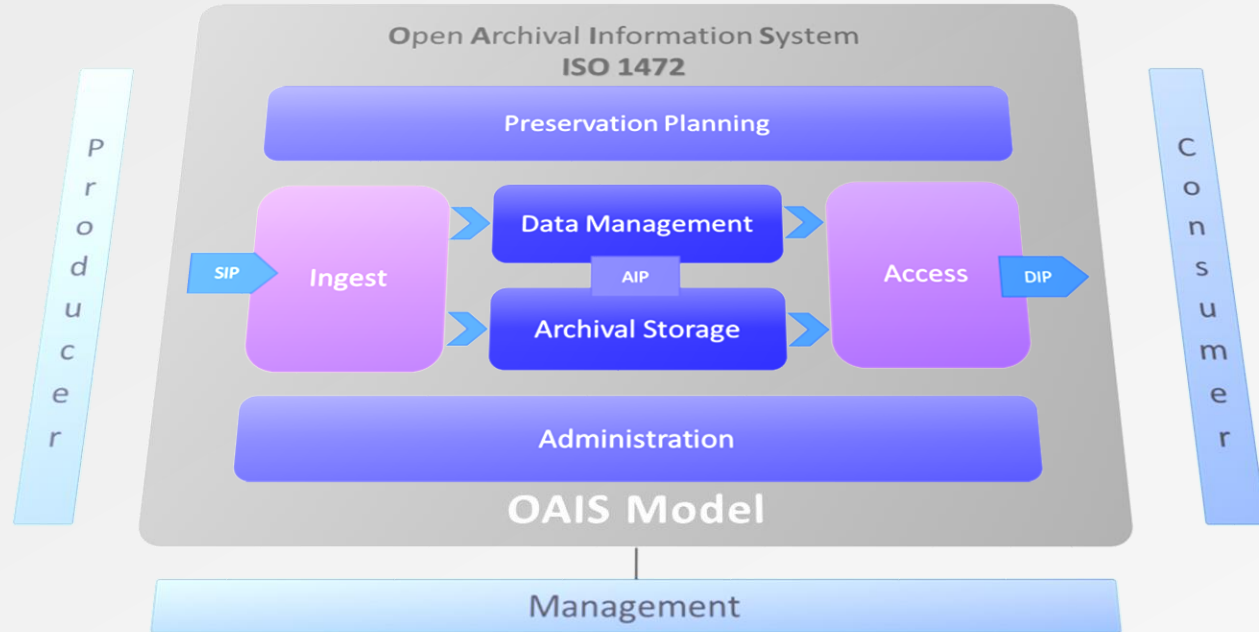
Researcher's Journey



2.2 Open Architecture

DLCM Project

Innovative Design



Innovative Design

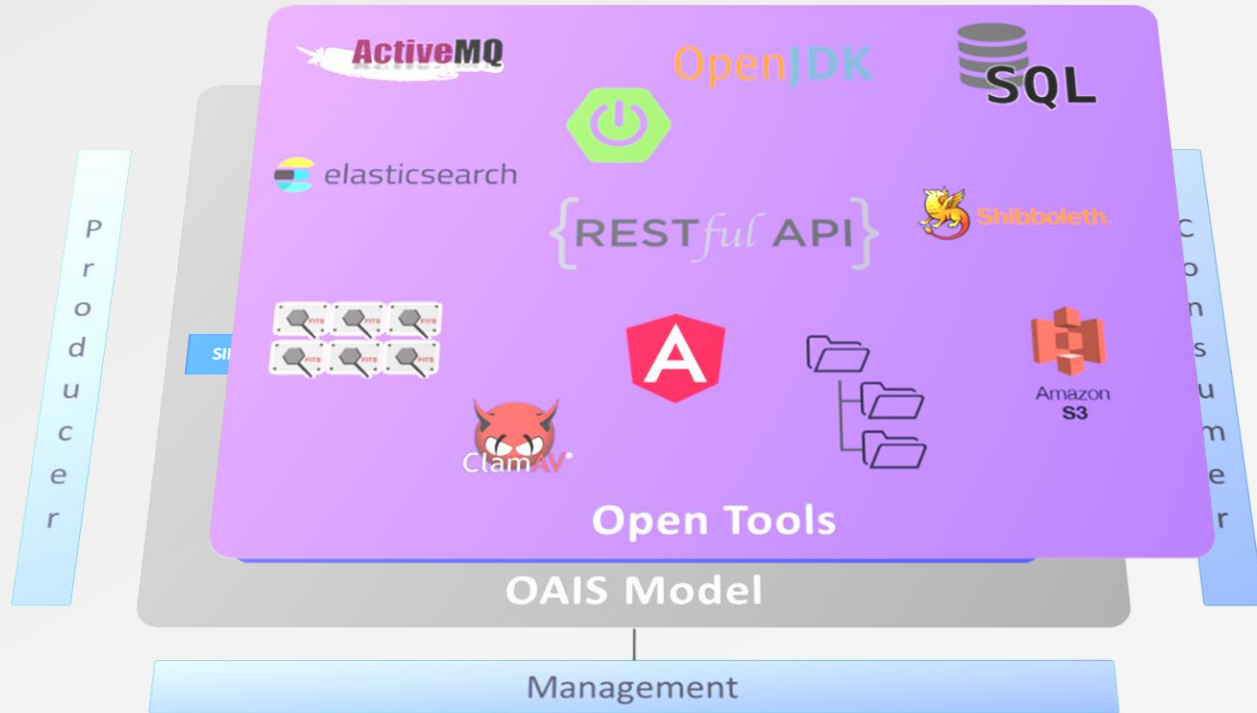
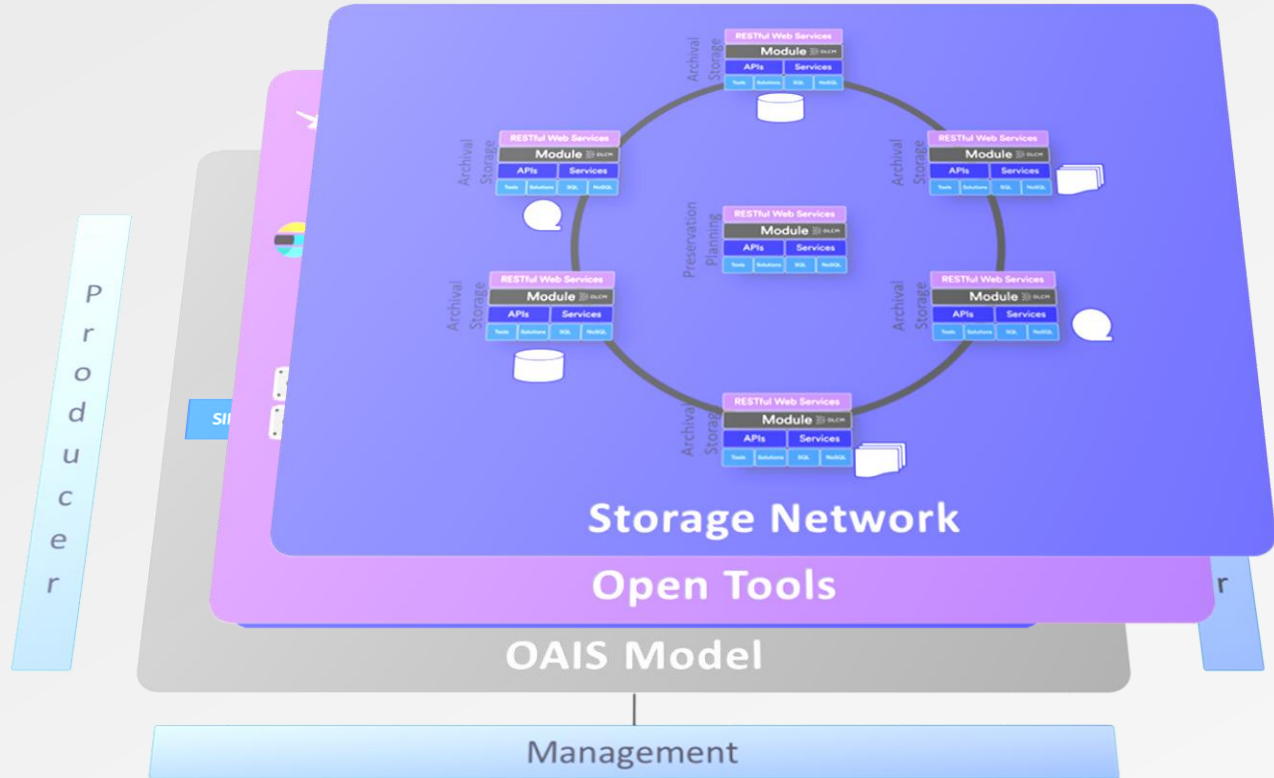


Image by [Painting Valley](#) CC BY-NC 4.0

Innovative Design



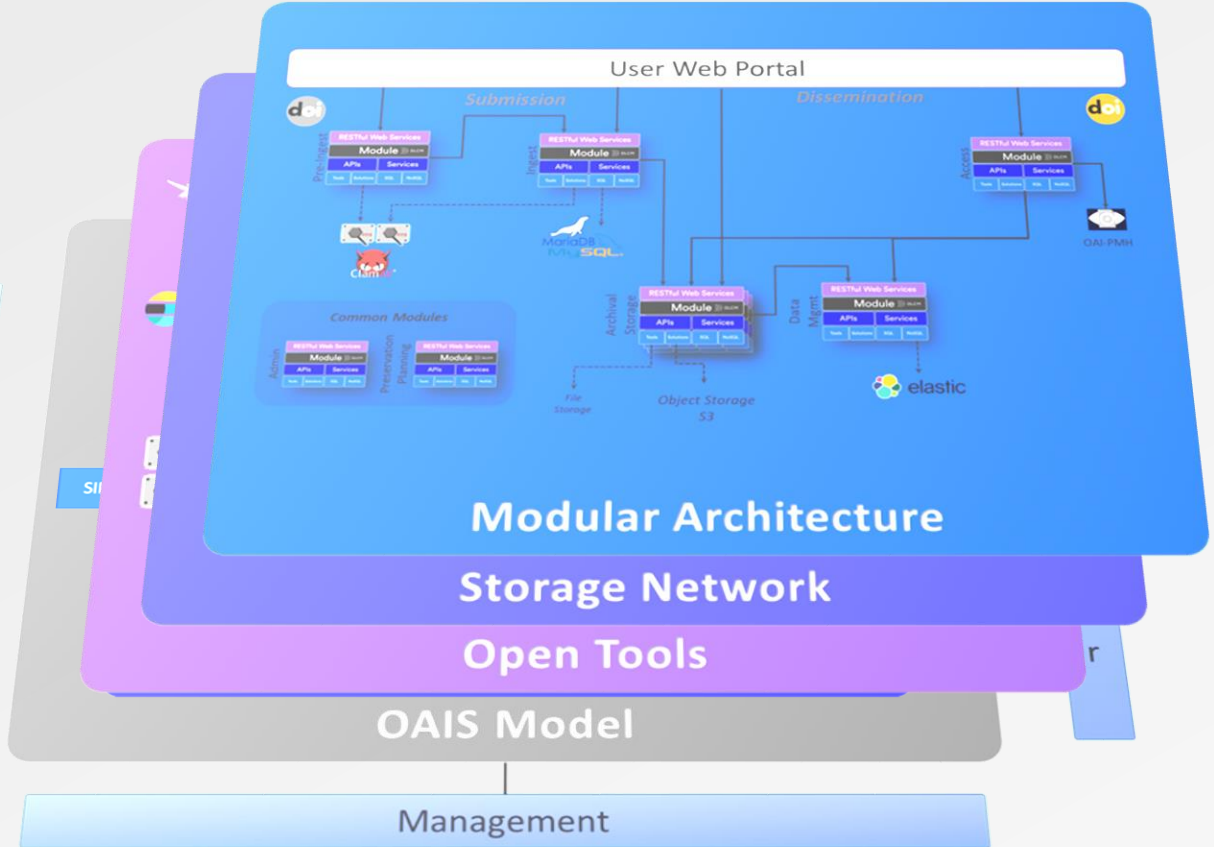
Image by [Painting Valley](#) CC BY-NC 4.0



Innovative Design



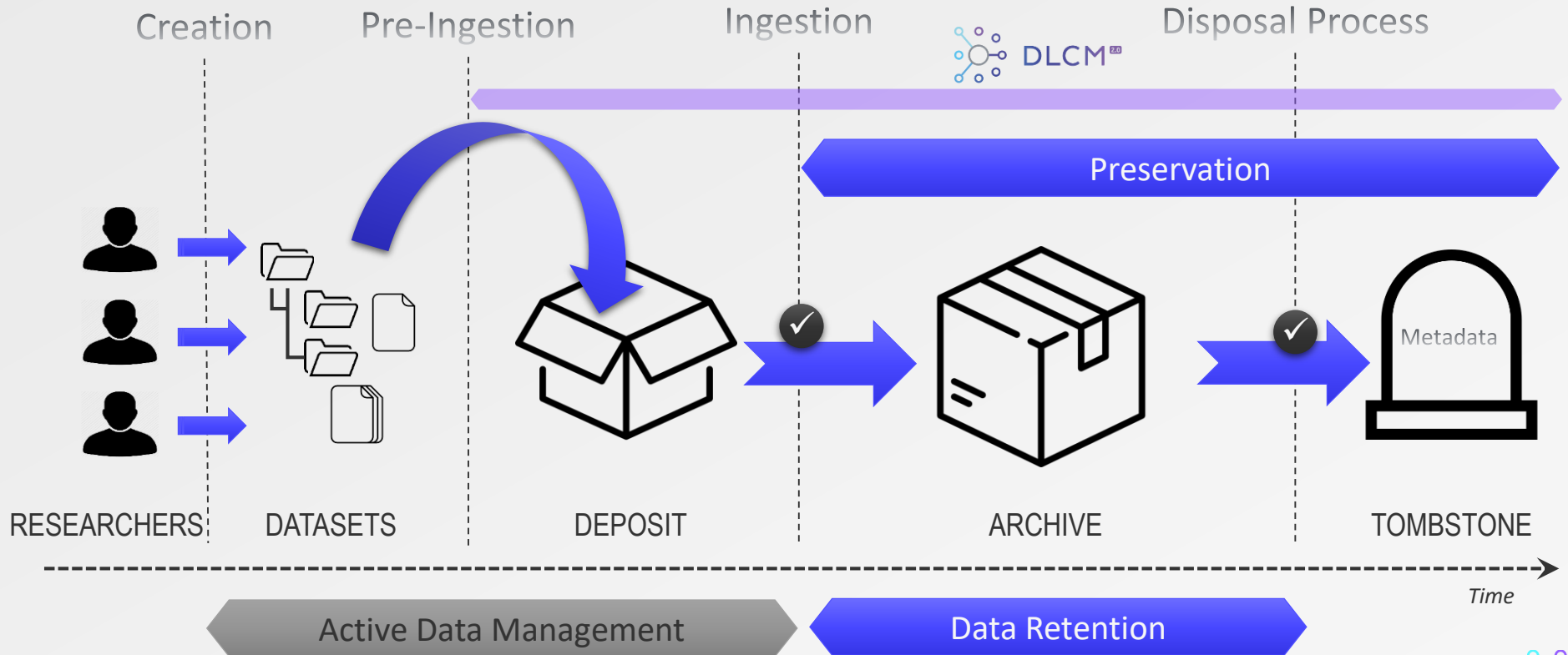
producer



2.3 Advanced Features

DLCM Project

What is Preserving Data?



Backup vs. Archiving



ARCHIVE

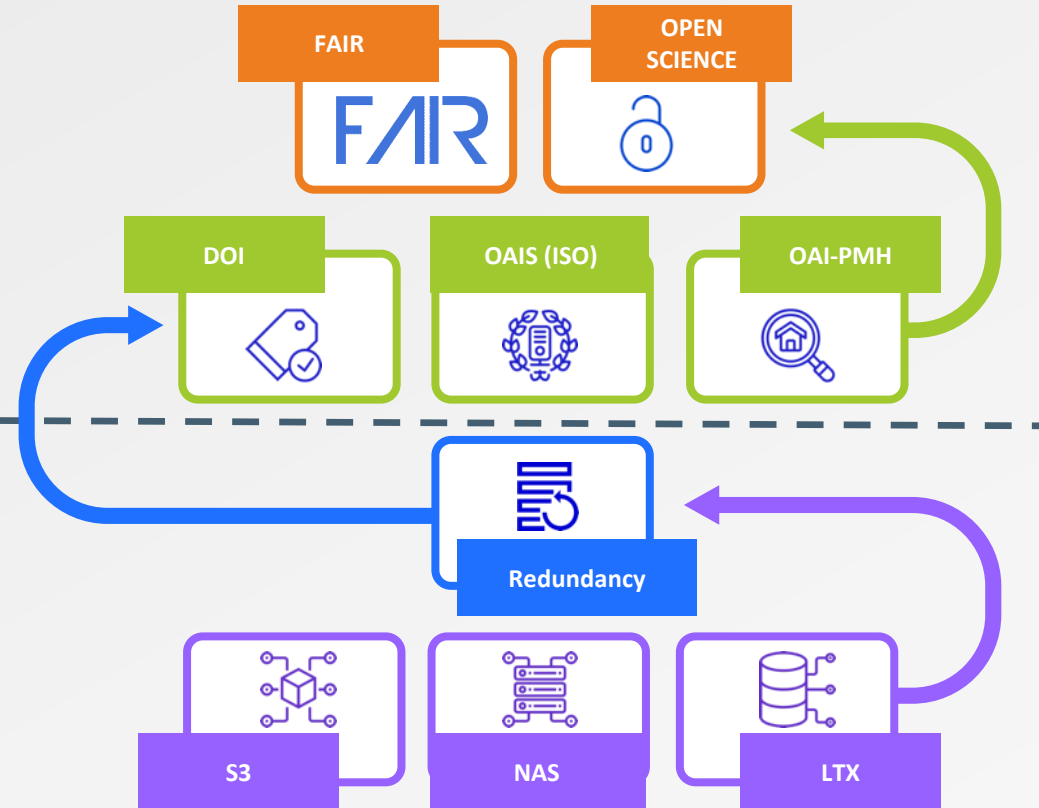
Preserve information as required by regulations and institutional policies

- Auditable
- Follows a life cycle
- Self-described
- Data integrity

BACKUP

Insurance policy against unforeseen system failures

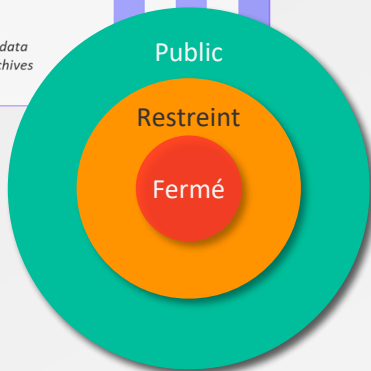
- Incremental
- Multiple snapshots
- Retained on short periods of time
- Not searchable



Features & Best Practices



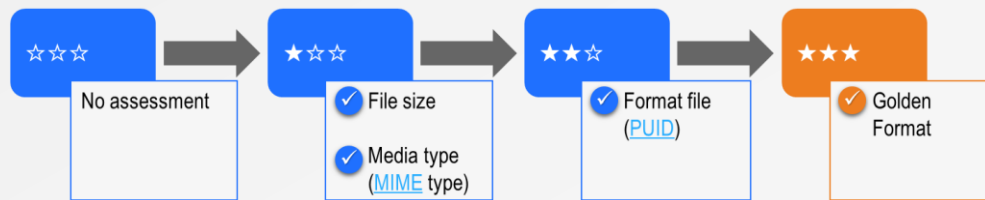
USER ROLES



	Access Level		
	Public	Restricted	Closed
Compatibility	Public	Restricted	Closed
Not defined	✓	✓	✓
Public	✓	✓	✓
Controlled public	✗	✓	✓
Accountable	✗	✗	✓
More accountable	✗	✗	✓
Fully accountable	✗	✗	✓
Maximum restricted	✗	✗	✓

Data Sensitivity (Data tag)

Compliance Level



PUID = PRONOM Unique Identifiers
The National Archive, UK
<https://www.nationalarchives.gov.uk/PRONOM>

Golden Format = recommended format for long term preservation
Library of Congress, USA
<https://www.loc.gov/preservation/conserv/ffu/>

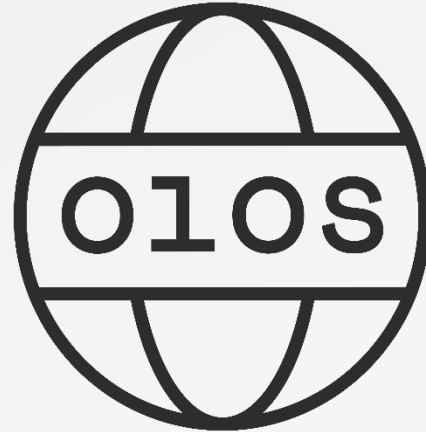


DLCM^{2.0}



YARETA

<https://yareta.unige.ch/>



<https://olos.swiss/portal/>

3. What Next ?

DLCM Project

Adding Redundancy: Byzantine Agreement Problem

- First studied by Lamport et al. in 1982
- Distributed system with n replicas, which could tolerate f faults, $n = 3f + 1$
- LOCKSS* $n \gg 3f + 1 \rightarrow$ from elections to opinion polls
- LOCKSS' design principles:
 - Cheap storage (unreliable), diverse technologies
 - Hash comparisons (and no stored checksums)
 - No long-term secret for identities (just for polling \rightarrow sha1 enough)
 - Use inertia - running very slowly limits attacks and damage
 - Integrate intrusion detection intrinsically and assume a strong adversary



* Lot of Copies Keep Stuff Safe

SAFE PLN

SAFE Archive FEderation

mun.ca



ucl.ac.be



ulb.ac.be



ugent.be



Université
de Montréal



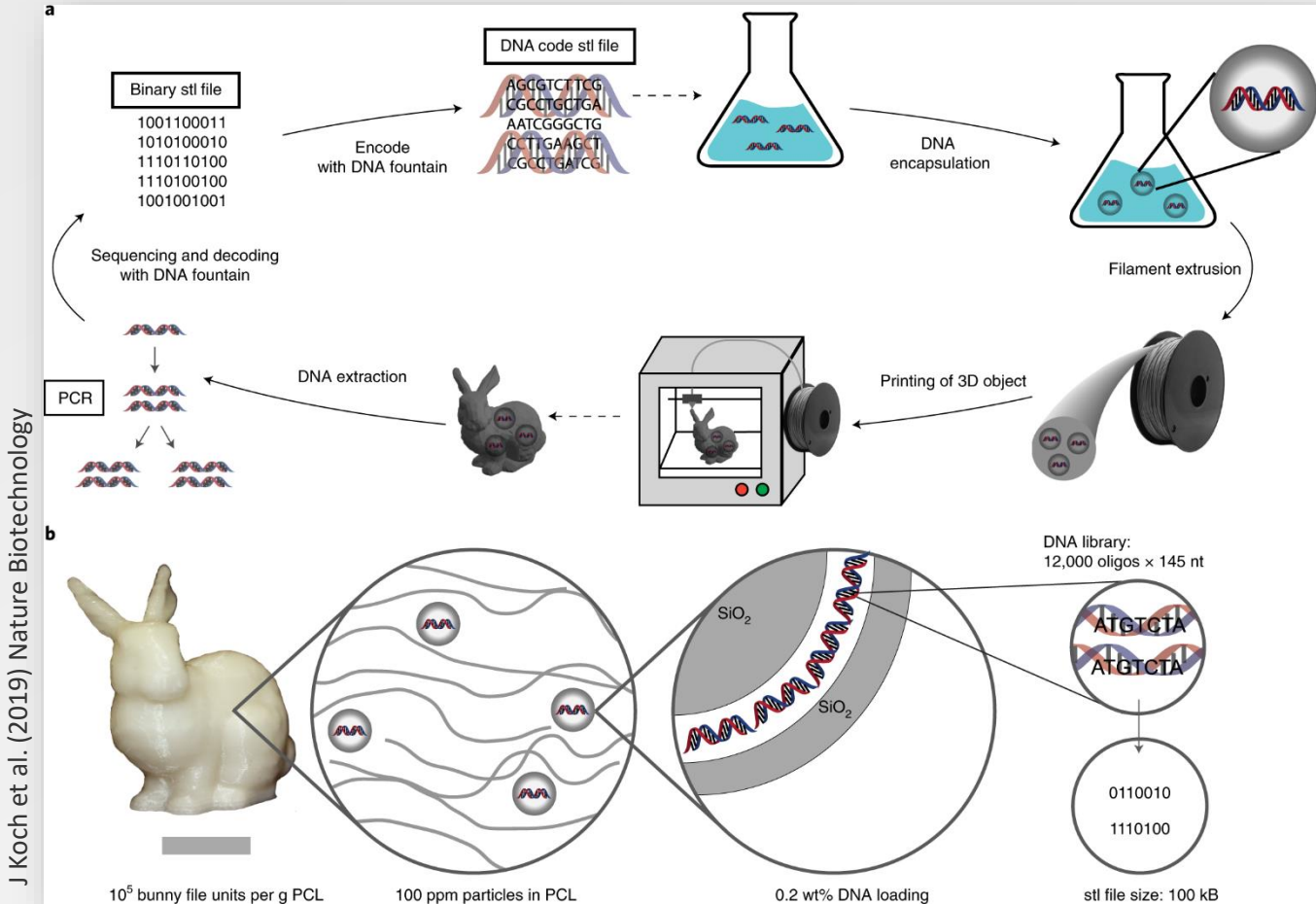
**UNIVERSITÉ
DE GENÈVE**

uni-bielefeld.de



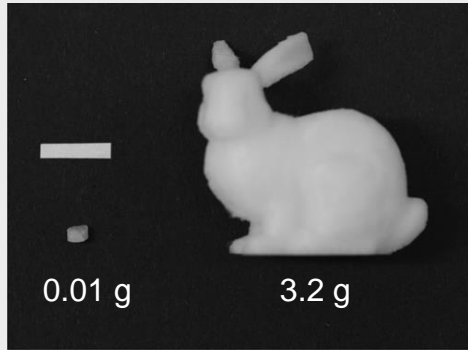
Used by UNIGE to archive electronic thesis

A DNA-of-things storage architecture to create materials with embedded memory

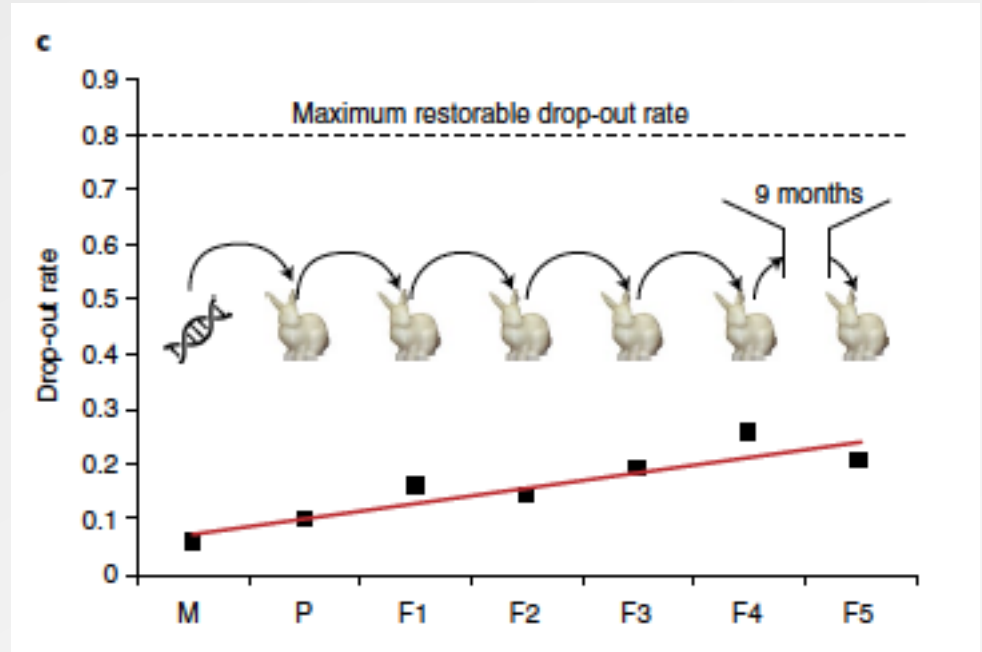


J Koch et al. (2019) Nature Biotechnology

... some numbers



0.01 g → 14'000 copies of the encoded file + redundancy

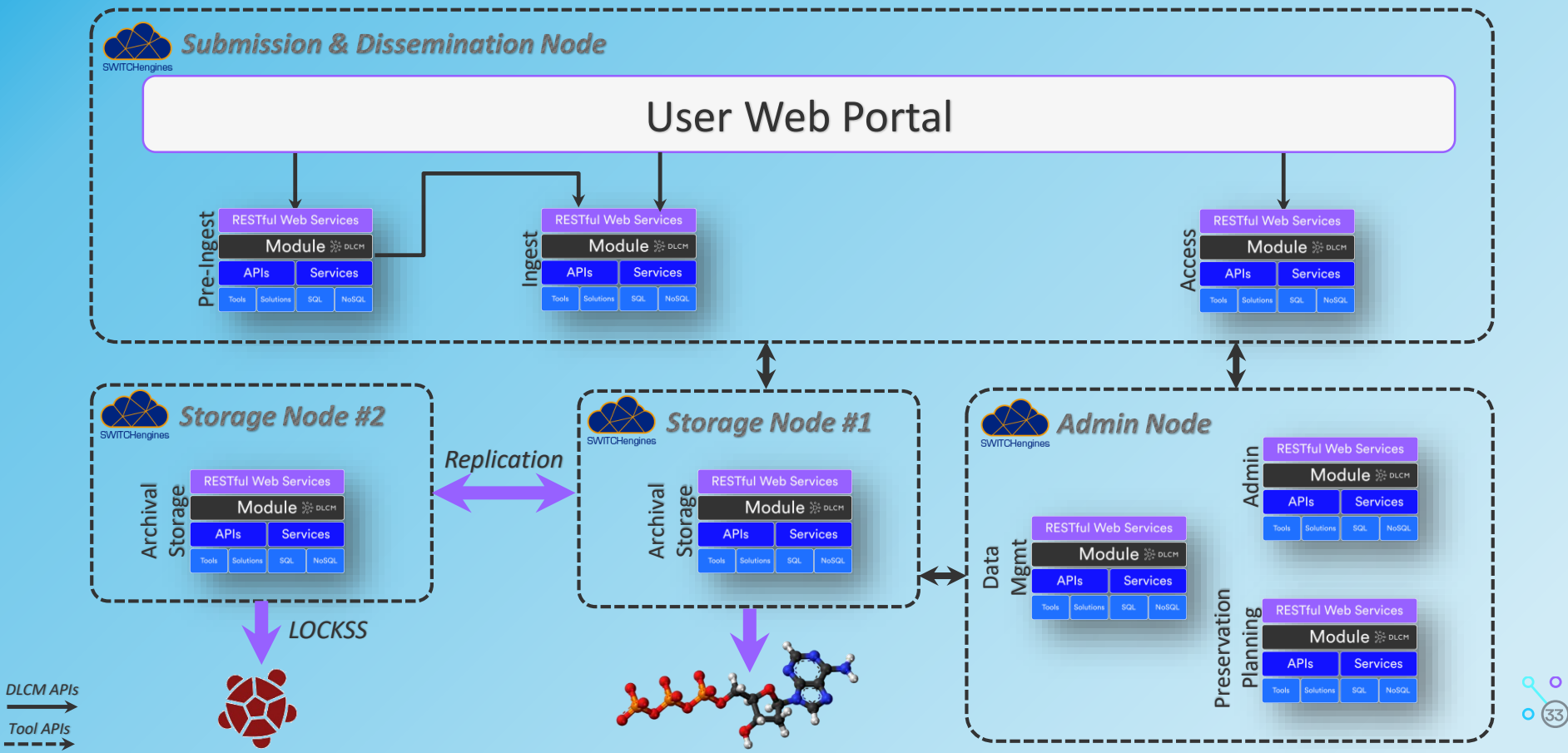


J Koch et al. (2019) Nature Biotechnology

DNA storage capacity: 215 PB per gram

Half-life of several thousands years

Future Cloud Deployment





Questions?

info@dlcm.ch

Go to <https://sandbox.dlcm.ch/>



Projet DLCM, 2021

This document is under license Creative Commons Attribution – Share Alike 4.0 International: <http://creativecommons.org/licenses/by-sa/4.0/deed.fr>