# Long Term Storage Service at CSCS – An Introduction

CSCS Webinar

Mario Valle, CSCS

June 23, 2021

# Webinar agenda

1. Why the Long Term Storage service is needed? (Mario Valle)

2. Technical structure of the Long Term Storage service (Giuseppe Lo Re)

3. Demo of the service (Stefano Schuppli)

4. Q&A

CSCS

ETH zürich

# Volume of scientific data is not the (only) issue
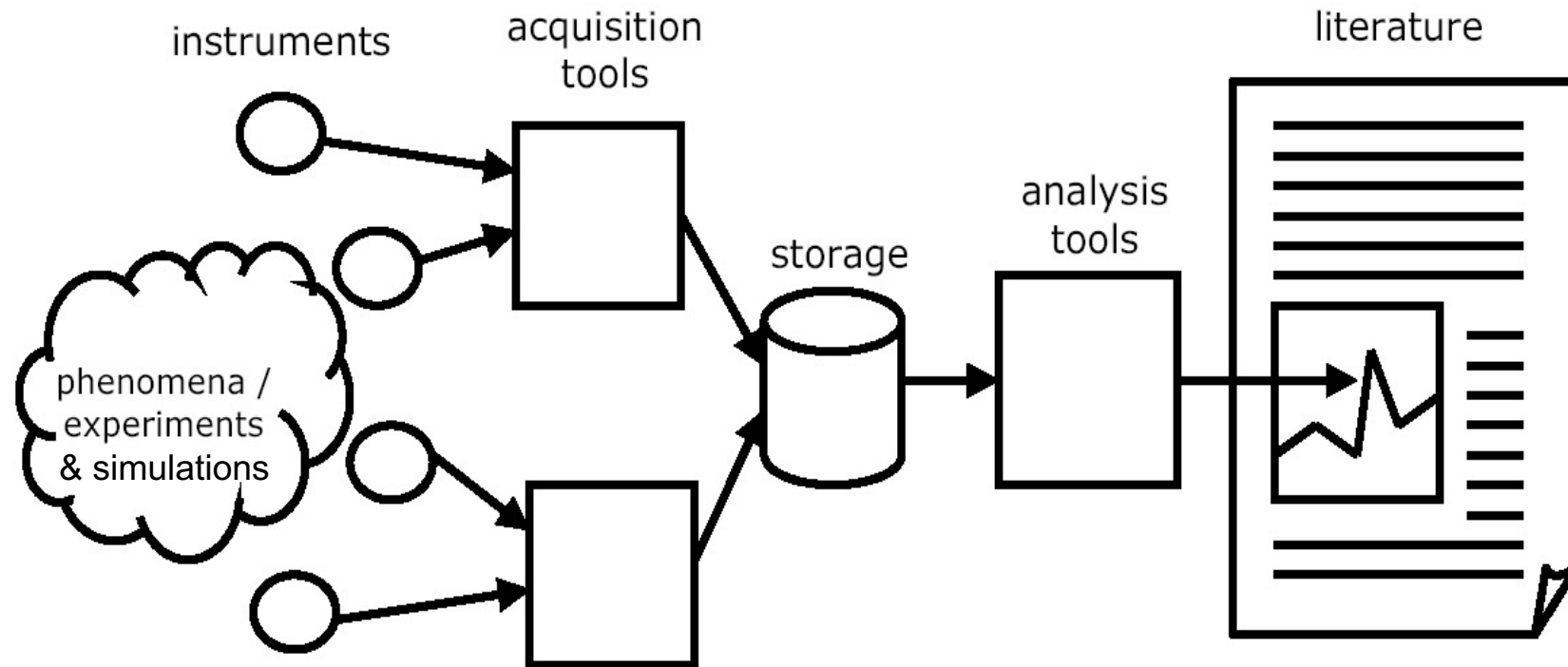


Inside the CSCS tape archive

Every CSCS user remembers the Richard Hamming's admonition:

**"Purpose of computation is insight, not numbers"**

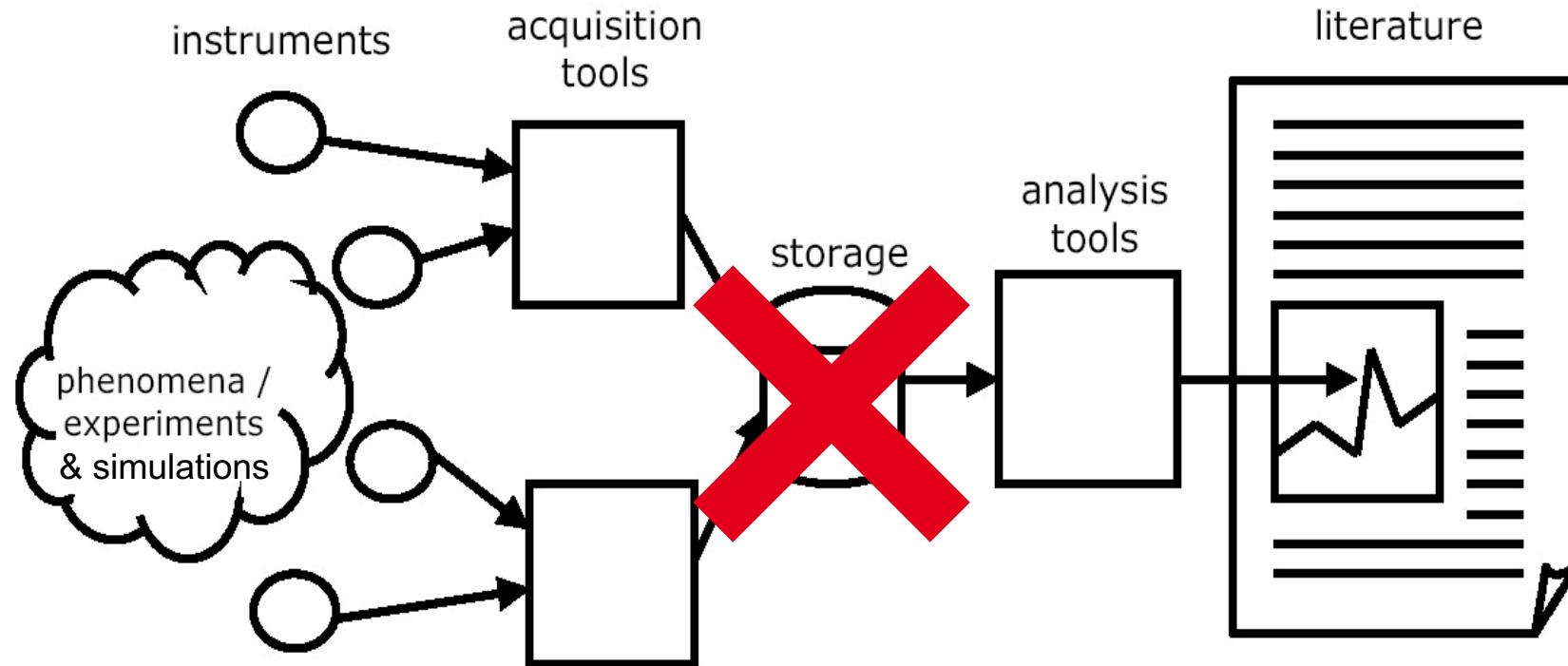We produce data and more data to help our insight, but often the result is that:

**"We are drowning in data, but starving of information"**

cscs

**ETH** *zürich*

# The usual scientific data lifecycle

# The usual scientific data lifecycle often ends with publication

Too often, after publications, data might consciously or unconsciously be forgotten and sometimes even throw away.
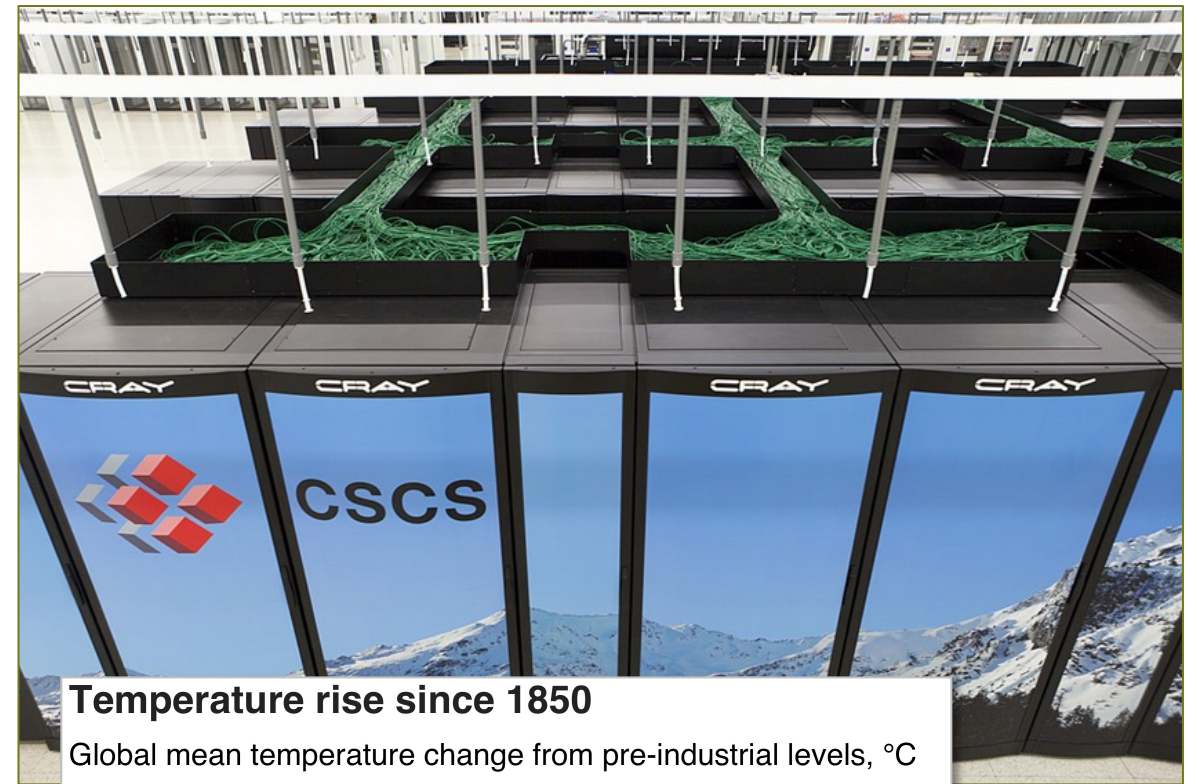
# Data creation is costly

Data are costly to collect or generate:

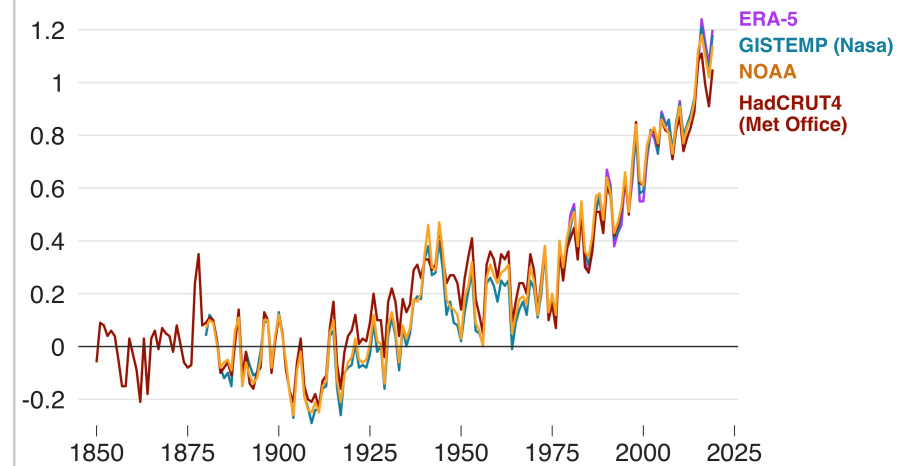- Compute/hours for simulation
- Telescope time
- Etc.

Some data cannot even be reproduced at will:

- Earthquake data
- Climate records
- Etc.



**Temperature rise since 1850**

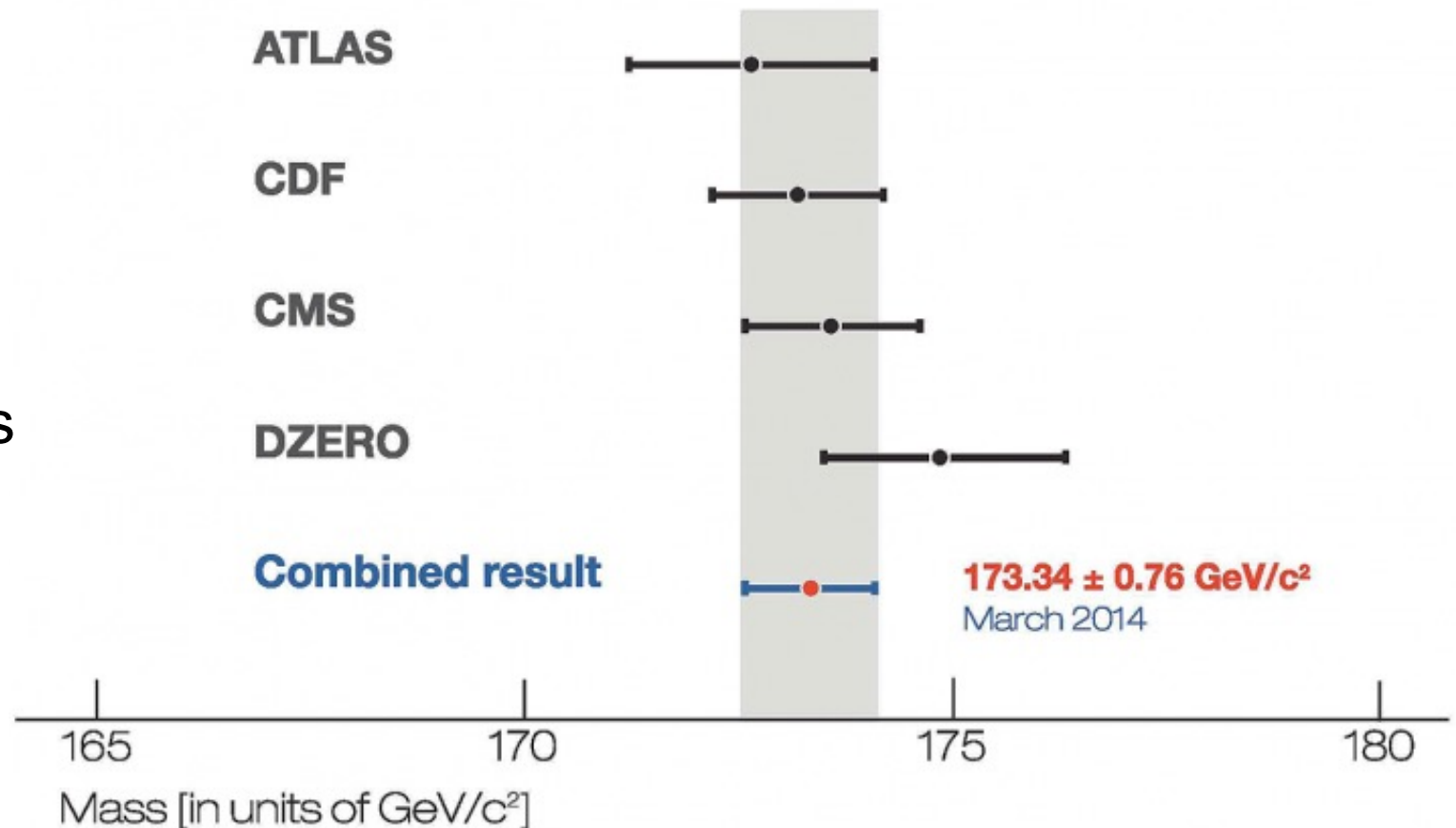Global mean temperature change from pre-industrial levels, °C

Source: Met Office
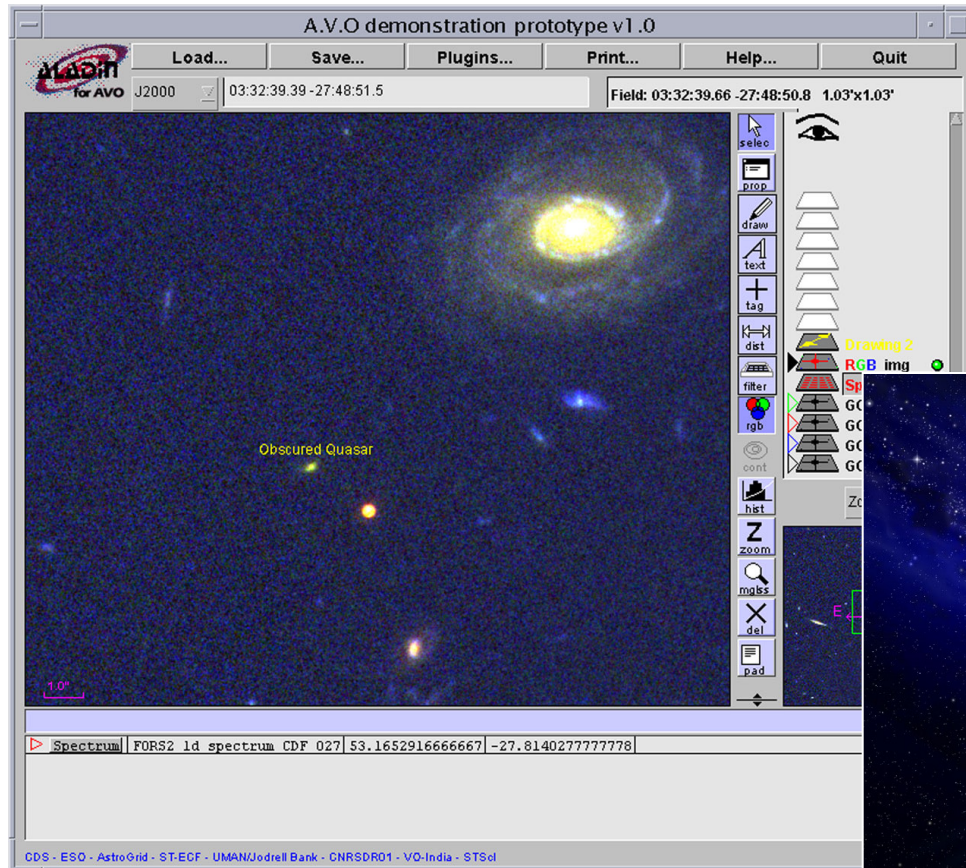
# Data could be reanalyzed or combined

It is routine at CERN, for example

Here data from CERN LHC and Fermilab DØ experiments are combined to give a better estimate of the top quark mass



## Top quark mass measurements

ATLAS

CDF

CMS

DZERO

**Combined result** — 173.34 ± 0.76 GeV/c² March 2014

165 · 170 · 175 · 180

Mass [in units of GeV/c²]

cscs

**ETH** zürich
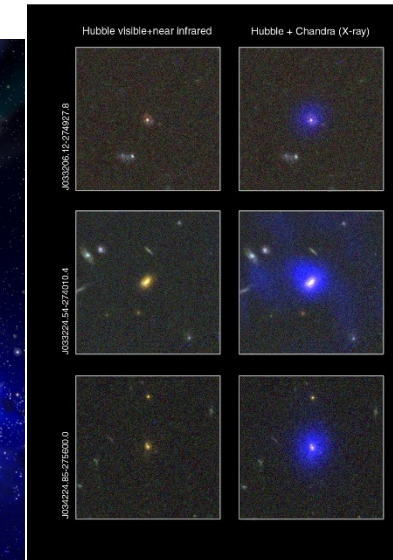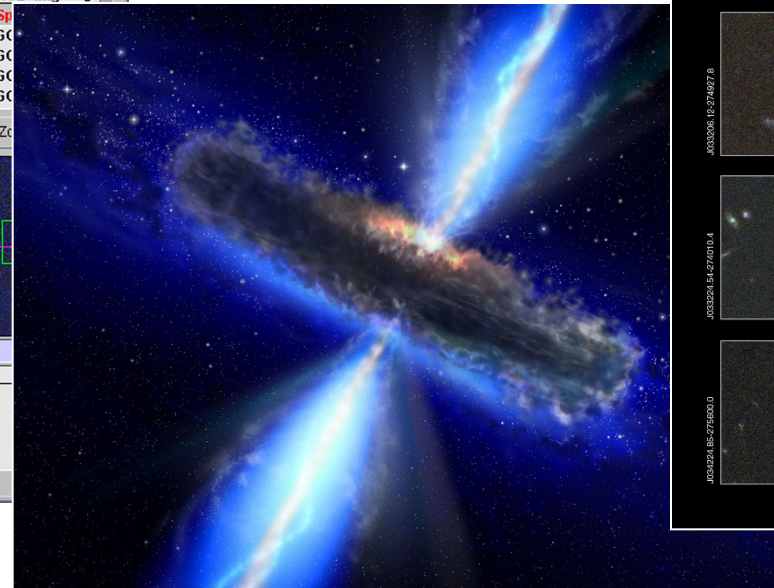
# Archived data could support "discovery by browsing"



*Paolo Padovani, "Discovering missing black holes: First Science from a Virtual Observatory", AVO (2004)*
*https://esahubble.org/news/heic0409/*

# Our data could become part of some scientific data collection

- There is already a growing list of managed data collections
- They make data accessible and "recycle" existing data

# New trend in scientific data lifecycle



**From** *publish and forget* **to** *use, reuse, recycle*

# So, be a data ecologist

Data is like a natural resource:

- Try to use it better
- Do not waste or pollute
- Recycle
- Assure quality
- Preserve for future generations

Scientific journals and funding agencies started this change by requiring data associated to papers/project to be preserved and made available

# FAIR principles

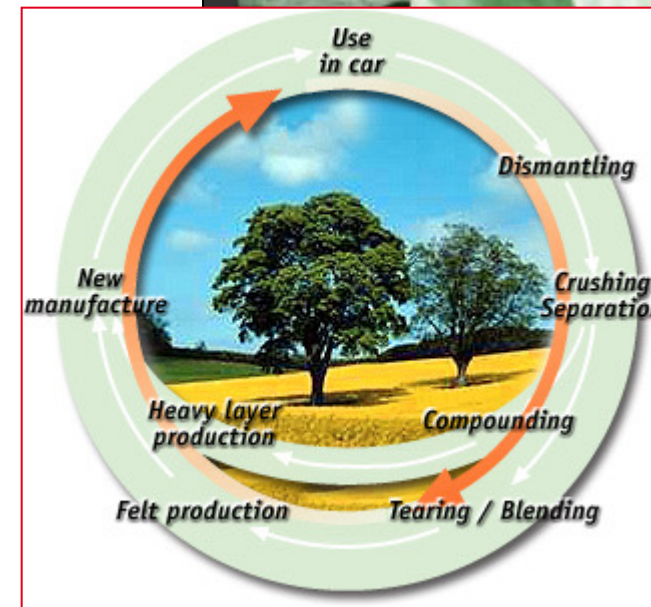The principles emphasize machine-actionability (i.e., the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention) because humans increasingly rely on computational support to deal with data as a result of the increase in volume, complexity, and creation speed of data.

https://www.go-fair.org/

Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.
**FINDABLE**

Metadata and data are understandable to humans and machines. Data is deposited in a trusted repository.
**ACCESSIBLE**

Metadata use a formal, accessible, shared, and broadly applicable language for knowledge representation.
**INTEROPERABLE**

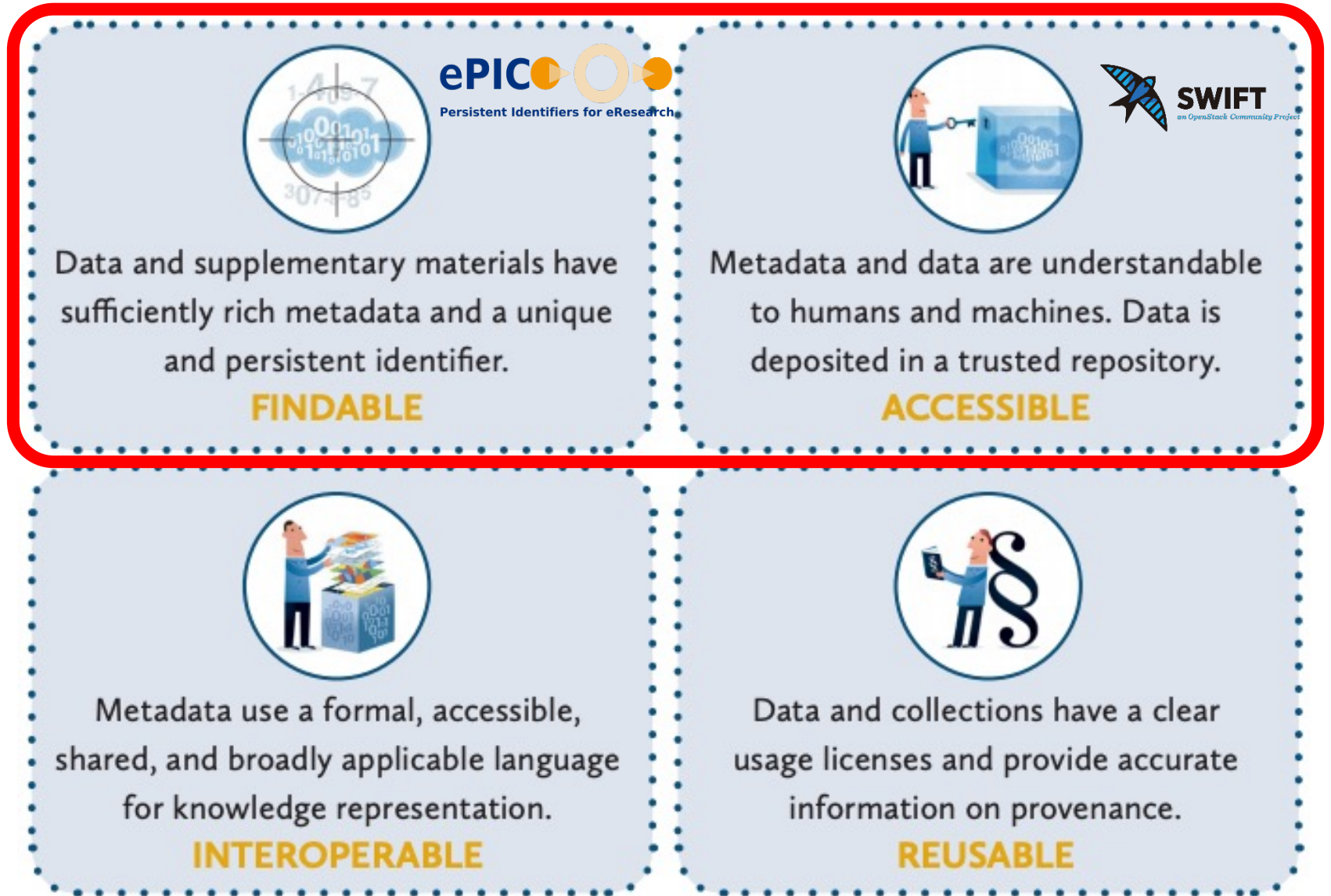Data and collections have a clear usage licenses and provide accurate information on provenance.
**REUSABLE**

# FAIR and CSCS data services

The four FAIR Principles, together with increasing requirement from funding agencies to have the data produced by their funded research publicly available, has motivated CSCS to offer its:

**Long-Term Storage (LTS) Service**



ePIC
Persistent Identifiers for eResearch

Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.
**FINDABLE**

SWIFT
an OpenStack Community Project

Metadata and data are understandable to humans and machines. Data is deposited in a trusted repository.
**ACCESSIBLE**

Metadata use a formal, accessible, shared, and broadly applicable language for knowledge representation.
**INTEROPERABLE**

Data and collections have a clear usage licenses and provide accurate information on provenance.
**REUSABLE**

# Findability by the CSCS PID service

To make LTS data findable, CSCS provides a service to generate and manage a certain range of **Persistent Identifiers (PID)** assigned to Switzerland by the ePIC consortium and to resolve them.
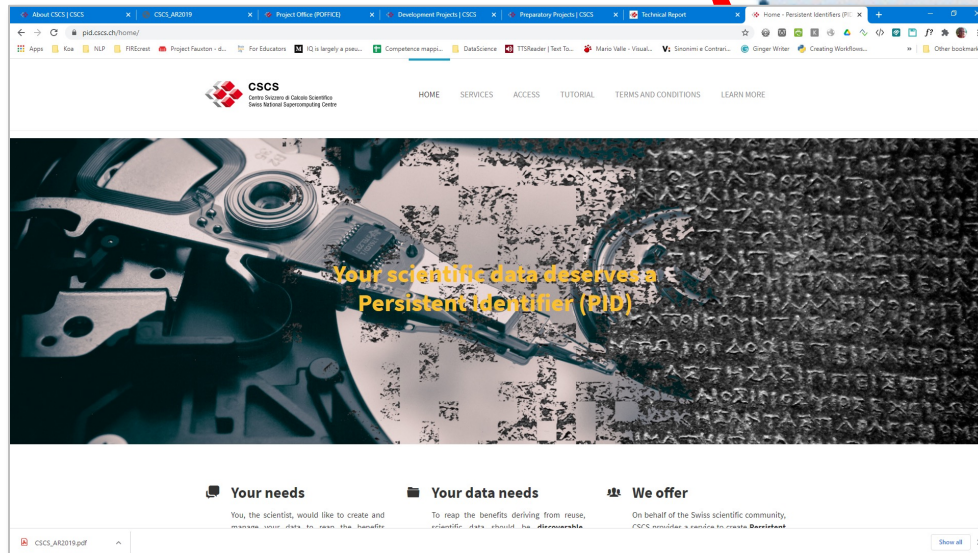


Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.

**FINDABLE**

Metadata and data are understandable to humans and machines. Data is deposited in a trusted repository.

**ACCESSIBLE**

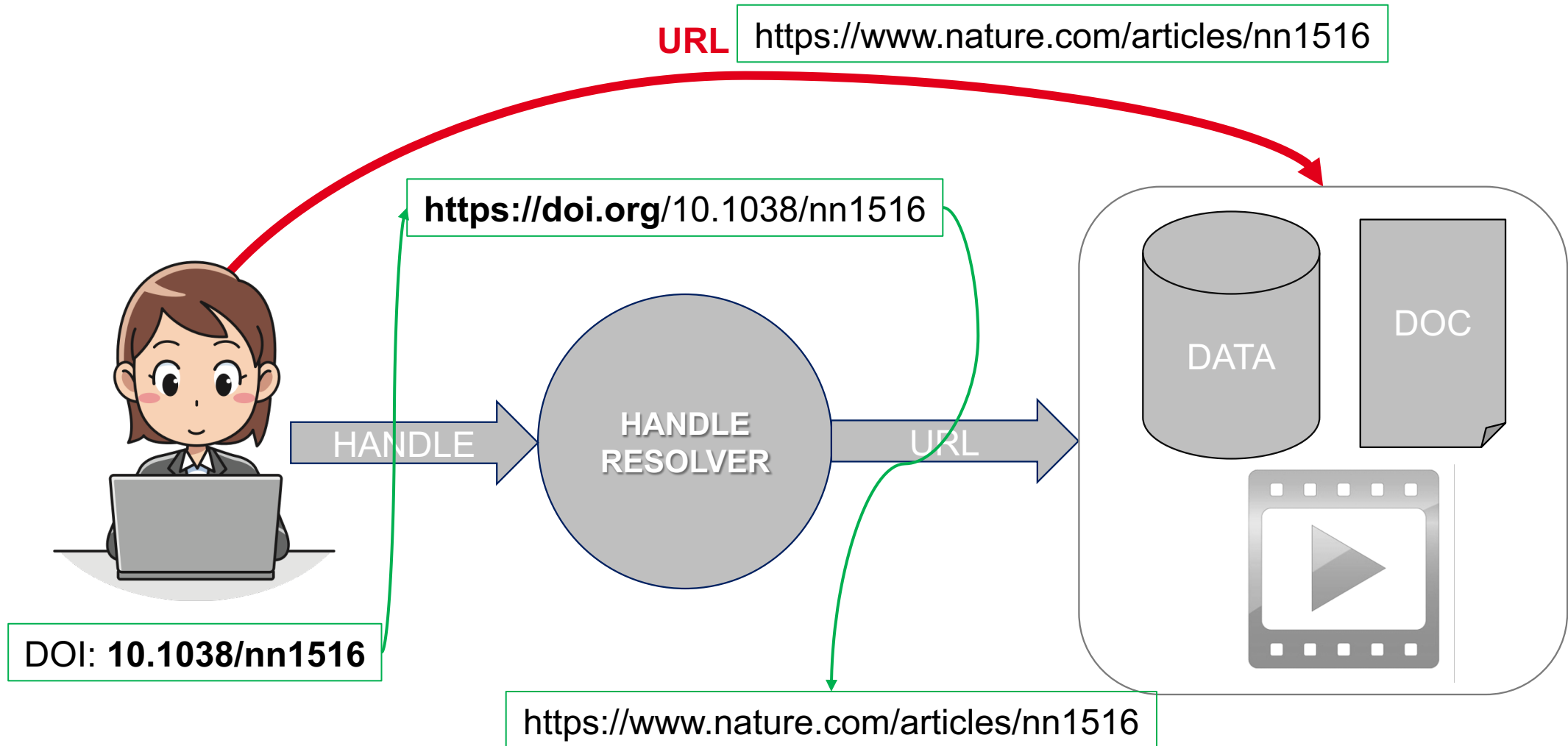use a formal, accessible, broadly applicable language wledge representation.

**TEROPERABLE**

Data and collections have a clear usage licenses and provide accurate information on provenance.

**REUSABLE**
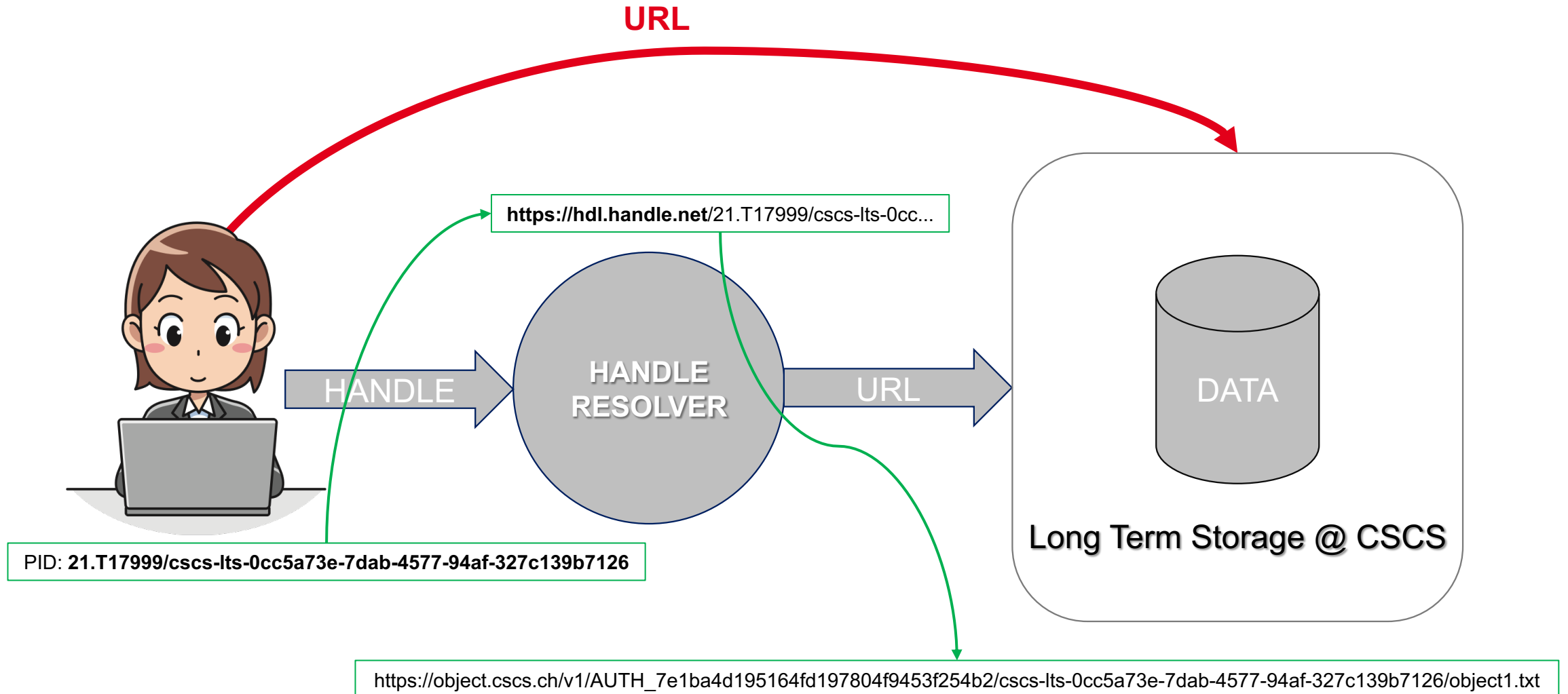
cscs

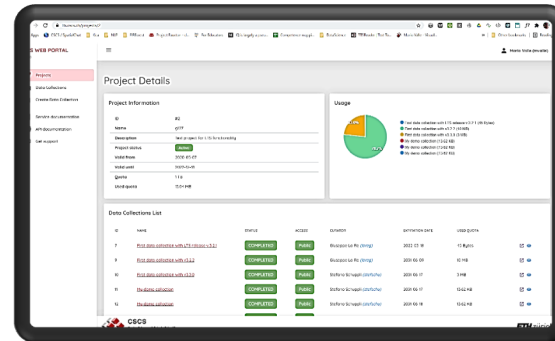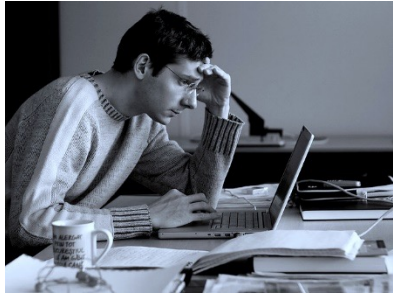ETH *zürich*

# Base of every PID system (DOI, Handle, URN, ARK, PURL, ISBN…)



URL https://www.nature.com/articles/nn1516

**https://doi.org**/10.1038/nn1516

HANDLE RESOLVER

HANDLE

URL

DATA

DOC

DOI: **10.1038/nn1516**

https://www.nature.com/articles/nn1516

cscs

**ETH** *zürich*

# Base of every PID system (DOI, Handle, URN, ARK, PURL, ISBN…)

URL

https://hdl.handle.net/21.T17999/cscs-lts-0cc...

HANDLE

HANDLE RESOLVER

URL

DATA

Long Term Storage @ CSCS

PID: **21.T17999/cscs-lts-0cc5a73e-7dab-4577-94af-327c139b7126**

https://object.cscs.ch/v1/AUTH_7e1ba4d195164fd197804f9453f254b2/cscs-lts-0cc5a73e-7dab-4577-94af-327c139b7126/object1.txt
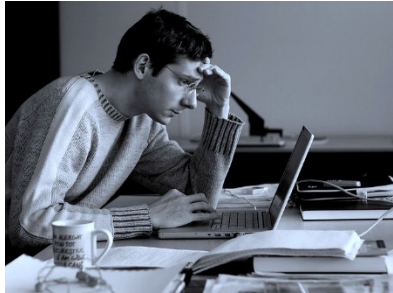
CSCS

ETH zürich

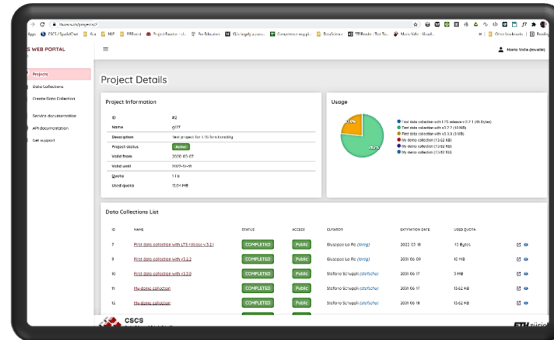# Long term storage at CSCS access through PID



CSCS Long-term
storage portal

Access simulation
results using PID

CSCS

ETH zürich

# Long term storage at CSCS shields users from technology changes



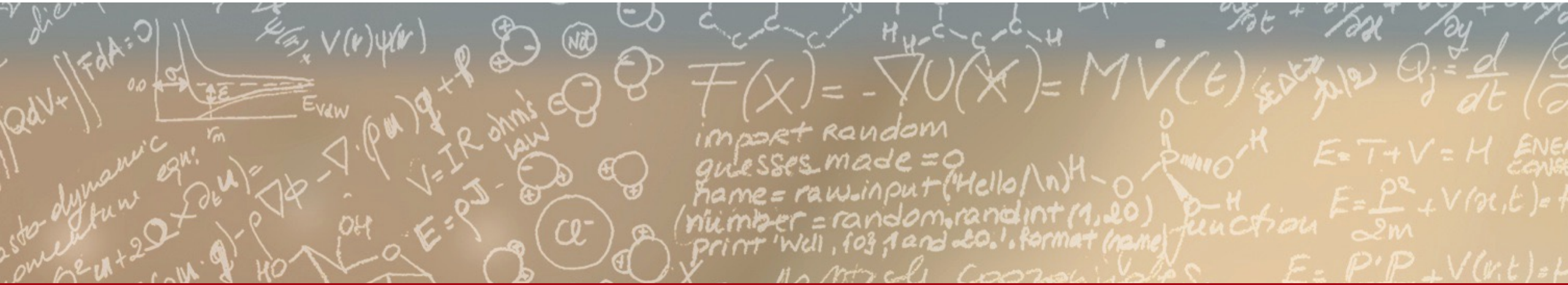Access simulation results using the same PID

CSCS Long-term storage portal

Data migration

**Thank you for your attention.**