



Data Preservation

Techniques, Traps and Processes

Shaun de Witt, United Kingdom Atomic Energy Authority

RECAP

- Why do we preserve data
 - Requirement of funders, support claims, speed research, protect against data loss
- Data Loss scenarios
 - Media loss
 - Format loss
 - Bad hats
 - Funding loss
 - Link loss

Media loss



*HDD Failure
Often get 'soft' warnings
Failure rates 1-8% p.a.
Data recovery often impossible*



*Human Failure
No warning, moderate
rate
Data recovery often
impossible*

*Tape Failure
Very reliable
Increasing chance of failure with number of reads
High chance of partial data recovery*



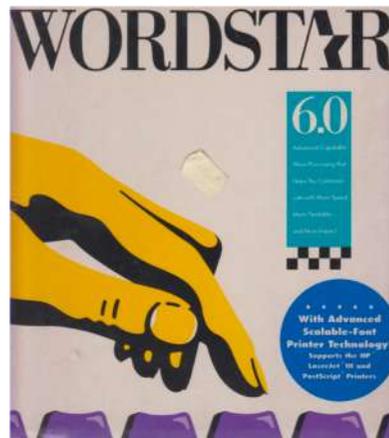
*Planning Failure
Data still on media that can't be
read
Data recovery often
impossible*

Format Loss



But what about Science?

- *In-house formats still used (ipx in fusion)*
 - *But more and more 'open' formats are being adopted*
- *Myriad of formats*
 - [http://justsolve.archiveteam.org/index.php/Scientific Data formats](http://justsolve.archiveteam.org/index.php/Scientific_Data_formats)
 - *13 just for Astronomy!*
- *What about databases?*
 - *Still direct access to databases and poorly understood database design*



Commercial Obsolescence



Its not just word processors!

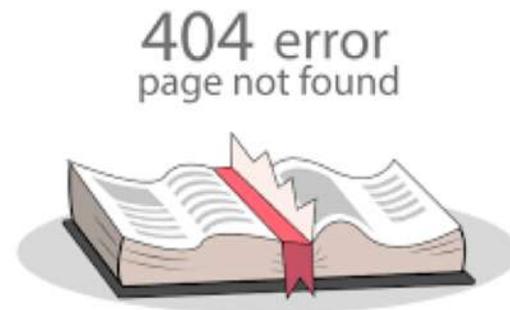
Bad Hats and accidents

- Bad hats can get in and delete, modify, or add to your data
 - In the 1990's Gigabytes of porn was dumped on University servers at Reading University
 - "Climategate"
- Sometimes people make mistakes
 - 600TB of particle physics data was accidentally deleted while trying to correct another problem (about 2/3 of all its data)
 - "I needed more space so I deleted this file called vmunix?"
 - Using `rm -rf *` as root (find directory ..)

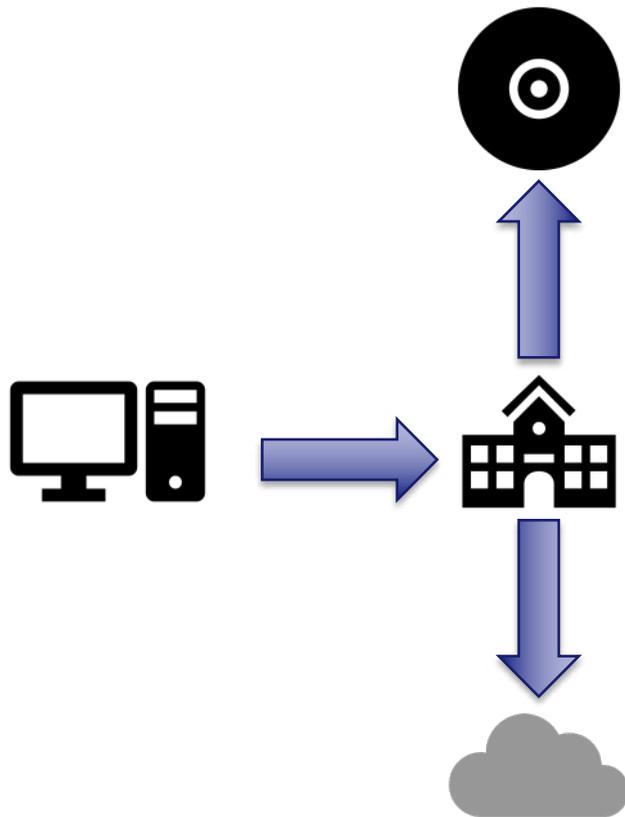


Link Rot

- 2014 study showed 50% of the URLs in U.S. Supreme Court opinions no longer link to the original information
 - doi:10.1017/1472669614000255
- August 2015 Weblock analyzed more than 180,000 links from references in the full-text corpora of three major open access publishers and found that overall 24.5% of links cited were no longer available.
 - <https://web.archive.org/web/20160304081204/https://weblock.io/report?id=all>
- 2016-17: Study of Yahoos! Directory links showed a half life of 2 years



Data Preservation Techniques - Backup



Traditional Backup to Media

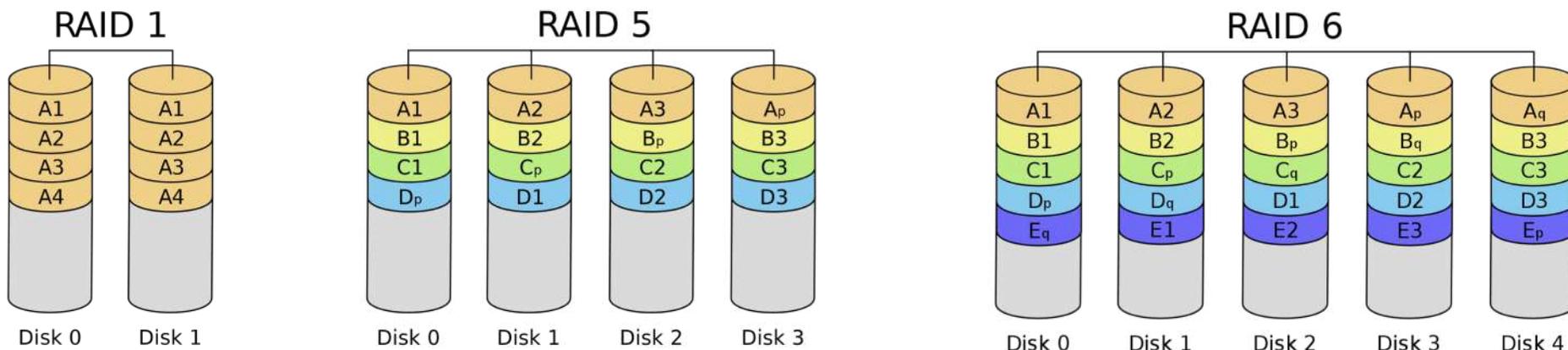
- *Still subject to media failure*
- *Only done once a day*
- *Effortful to recover*

Backup to 'cloud'

- ***Can be*** backed up more frequently
- Recovery ***can be*** costly and time consuming
- *Subject to commercial Terms and Conditions*
- *Expensive to use for long term preservation*
 - *E.g. UKAEA cloud backup is only 30 days*
 - *For longer term preservation we still use n in-house system*
- *In both cases it may be difficult to recover version history and is not suitable for large volumes of data*
 - *It can take > 1 day to back up scientific data sets*

All images from [en>User:Cburnett](#) - CC-BY-SQ3.0

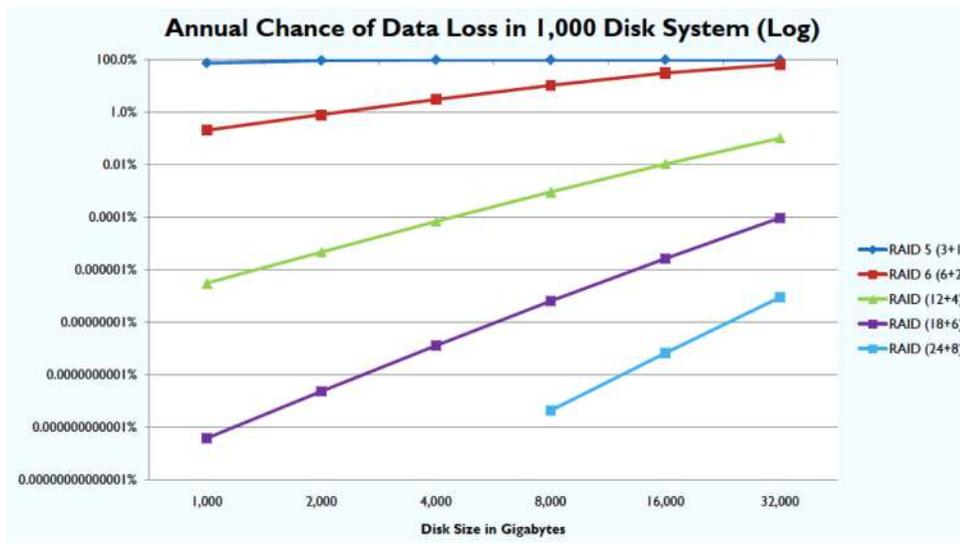
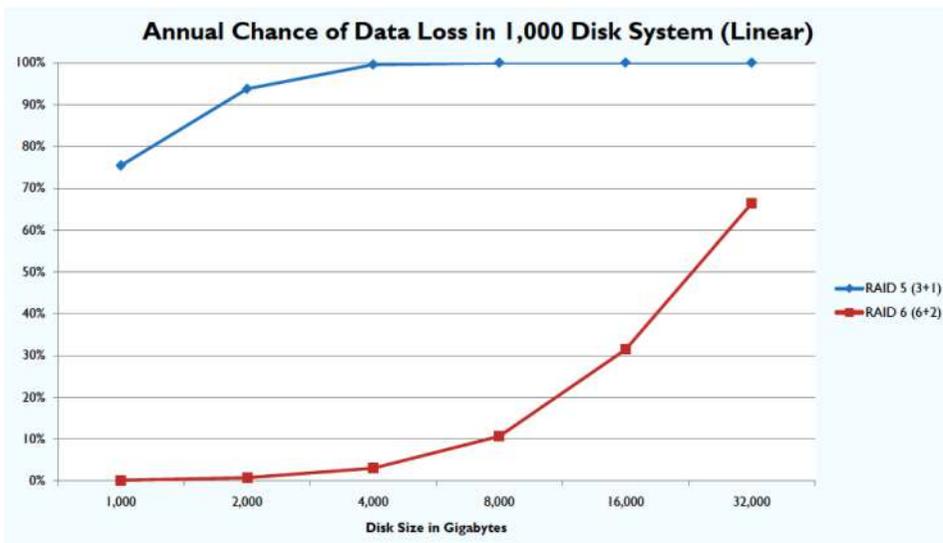
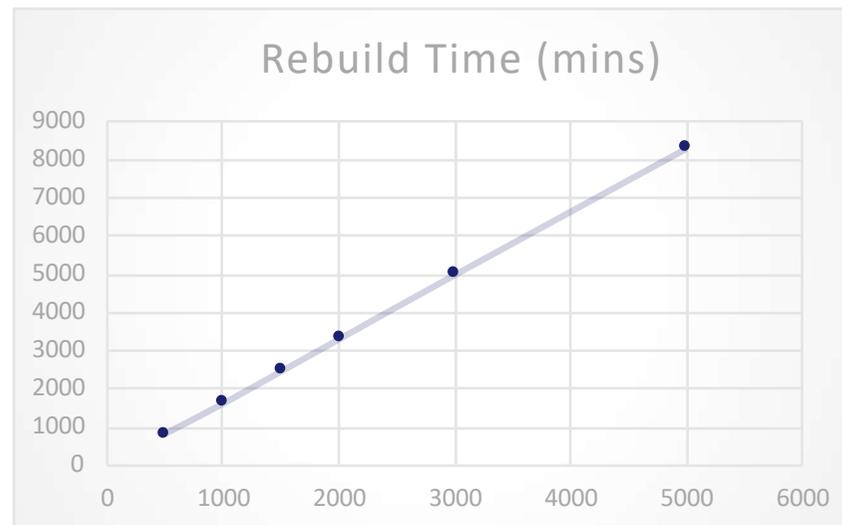
Data Preservation - RAID



RAID Level	Description	Min. Drives	Space Loss	Fault Tolerance
RAID 0	Striped	2	1	None
RAID 1	Mirroring	2	1/n	n-1 drive failures
RAID 1+0	Striped + mirrored	2	1/n	n-1 drive failures
RAID 5	Striped + 1 parity	3	1-1/n	1 drive failure
RAID 6	Striped + 2 parity	4	1-2/n	2 drive failures

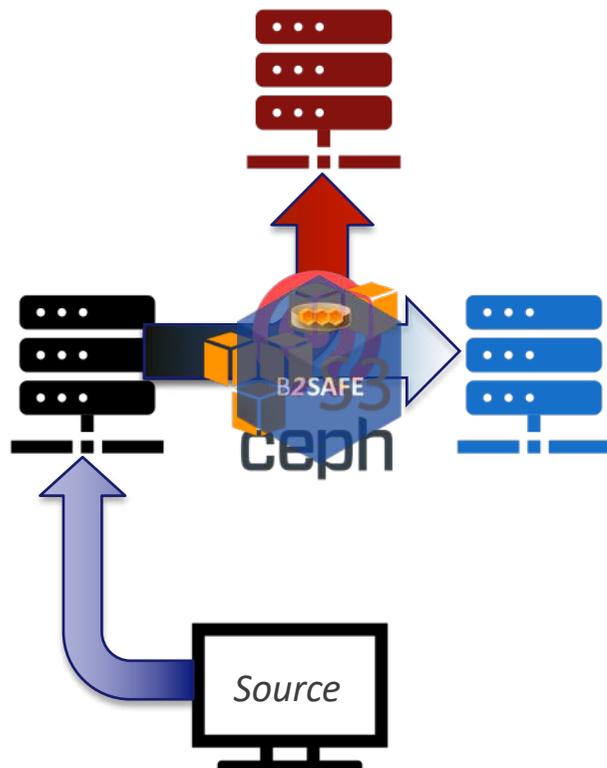
The problems with RAID

- Data rebuild times increase with disk size
- This graph is theoretical
 - Reality is double or triple this
- Typically 1 or two 'hot spares' are reserved to cover disk failures, so more space is lost than nominal
- And you still have the single site problem



© J. Resch, Cleversafe, Storage Developers Conference, 2010

Data Replication



- Commonly used within **object stores** with 3x replication
 - Normally replicas held at different places
- Also common on **cloud services** to give very high resilience
- Can be used distribute data to partner operations
 - E.g. data generated in UKAEA could be replicated to PSNC (another fusion research establishment) and CINECA (which hosts the community HPC facility)
- In most cases data integrity is maintained by **checking-summing** each object
 - Either as a background task or on retrieval
 - Many tools can self-heal in case of failure

The Good and the Bad of Checksums

- EXERCISE

Summary of Results



Original Image

MD5sum:

9bceae7b0b4c19df35052b9018c9e8d



Bit-flipped Image

MD5sum:

d46866f79200465f94baede9450928e0

<p>To Whom it May Concern:</p> <p>Alice Falbala fulfilled all the requirements of the Roman Empire intern position. She was excellent at translating roman into her gaul native language, learned very rapidly, and worked with considerable independence and confidence.</p> <p>Her basic work habits such as punctuality, interpersonal deperment, communication skills, and completing assigned and self-determined goals were all excellent.</p> <p>I recommend Alice for challenging positions in which creativity, reliability, and language skills are required.</p> <p>I highly recommend hiring her. If you'd like to discuss her attributes in more detail, please don't hesitate to contact me.</p> <p>Sincerely,</p> <p>Julius Caesar</p>	<p>May, 22, 2005</p> <p>Order:</p> <p>Alice Falbala is given full access to all confidential and secret information about GAUL.</p> <p>Sincerely,</p> <p>Julius Caesar</p>
---	--

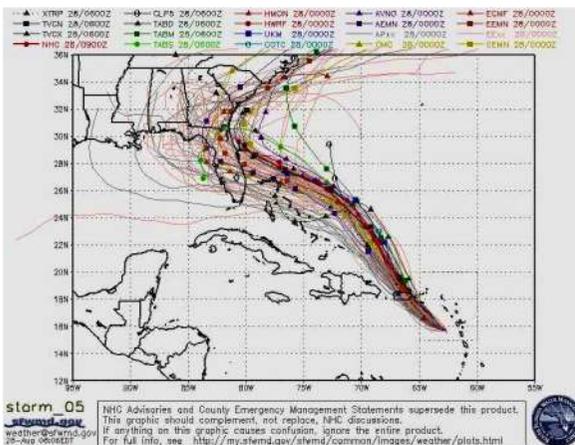


A Letter of Recommendation? Or Security Clearance?

A Panda? Or a Bucket?

Extending FAIR

- Making data **open** or **FAIR** is generally a good thing
 - Faster Science, More verifiable results, Better impact assessment, driving the data and digital economies, etc, etc...
- BUT... There are **risks**
 - **Accidental and Malicious**
 - Lets look at some examples



More Examples – The Reproducibility Crisis

Gene name errors are widespread in the scientific literature

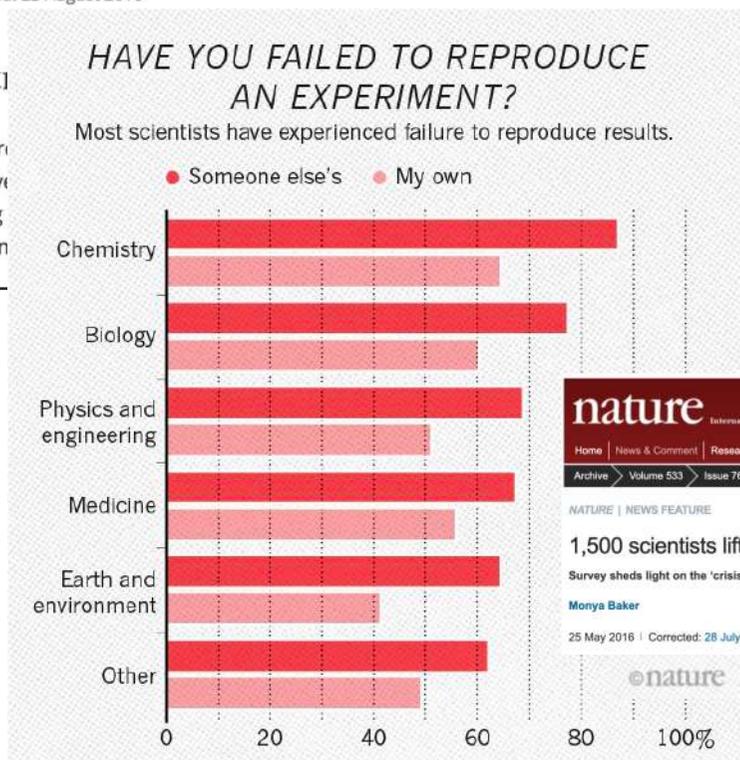
Mark Ziemann, Yotam Eren and Assam El-Osta

Genome Biology 2016 17:177 | DOI: 10.1186/s13059-016-1044-7 | © The Author(s). 2016

Published: 23 August 2016

Abstract

The spring to convert leading supplement



Ioannidis et al. Nature Genetics, 41, 2010
doi:10.1038/ng.295

ANALYSIS



56% of analyses could not be repeated, of which 30% were because of software issues. 50% did not state software version, 39% did not provide raw data. Only 11% could be reproduced satisfactorily.

Repeatability of published microarray gene expression analyses

John P A Ioannidis¹⁻³, David B Allison⁴, Catherine A Ball⁵, Issa Coulibaly⁴, Xiangqin Cui⁴, Aedin C Culhane^{6,7}, Mario Falchi^{8,9}, Cesare Furlanello¹⁰, Laurence Game¹¹, Giuseppe Jurman¹⁰, Jon Mangion¹¹, Tapan Mehta⁴, Michael Nitzberg⁵, Grier P Page^{4,12}, Enrico Petretto^{11,13} & Vera van Noort¹⁴

Given the complexity of microarray-based gene expression studies, guidelines encourage transparent design and public

CAMBRIDGE JOURNAL OF ECONOMICS

Issues | JEL | More Content | Submit | Purchase | About



Volume 38, Issue 2
March 2014

Does high public debt consistently stifle economic growth? A critique of Reinhart and Rogoff

Thomas Herndon, Michael Ash, Robert Pollin

Cambridge Journal of Economics, Volume 38, Issue 2, March 2014, Pages 257-279, <https://doi.org/10.1093/cje/bet075>

Published: 24 December 2013 | Article history

Cite | Permissions | Share

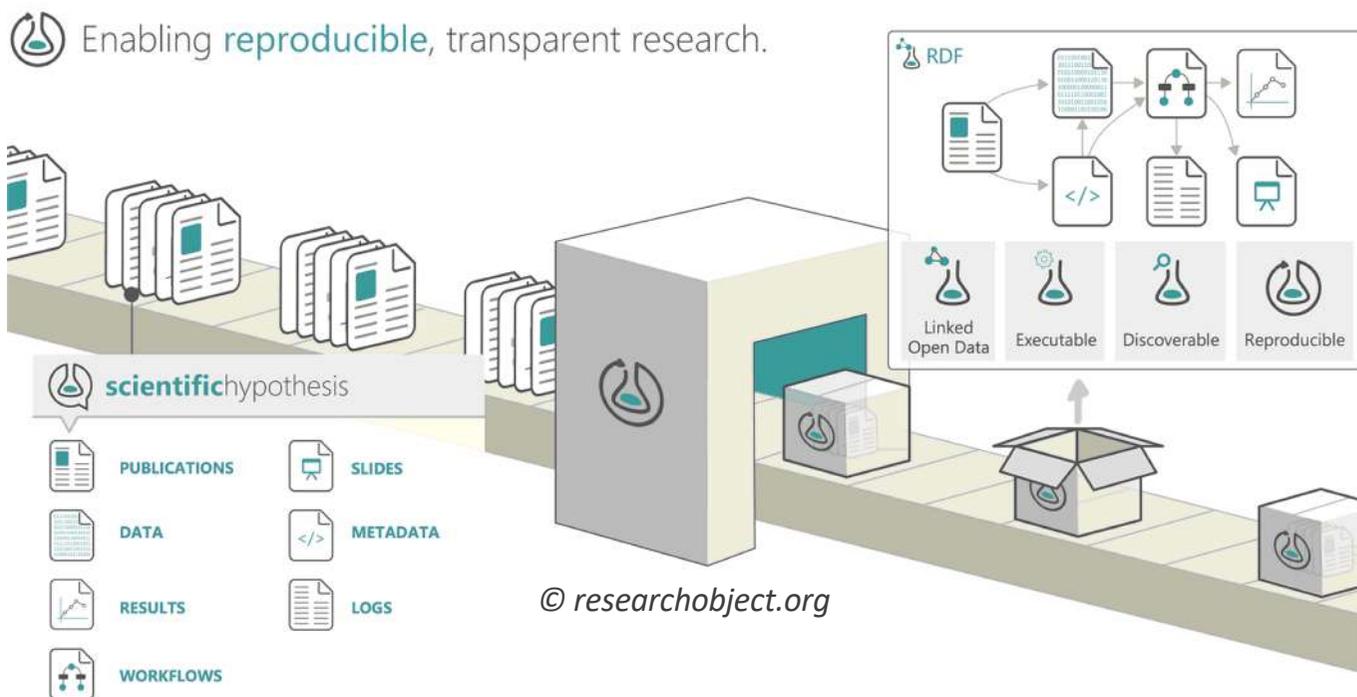
Abstract

We replicate Reinhart and Rogoff (2010A and 2010B) and find that selective exclusion of available data, coding errors and inappropriate weighting of summary statistics lead to serious miscalculations that inaccurately represent the relationship between public debt and GDP growth among 20 advanced economies. Over 1946-2009, countries with public debt/GDP ratios above 90% averaged 2.2% real annual GDP growth, not -0.1% as published. The published results for (i) median GDP growth rates for the 1946-2009 period and (ii) mean and median GDP growth figures over 1790-2009 are all distorted by similar methodological errors, although the magnitudes of the distortions are somewhat smaller than with the mean figures for 1946-2009. Contrary to Reinhart and Rogoff's broader contentions, both mean and median GDP growth when public debt levels exceed 90% of GDP are not dramatically different from when the public debt/GDP ratios are lower. The relationship between public debt and GDP growth varies significantly by period and country. Our overall evidence refutes RR's claim that public debt/GDP ratios above 90% consistently reduce a country's GDP growth.

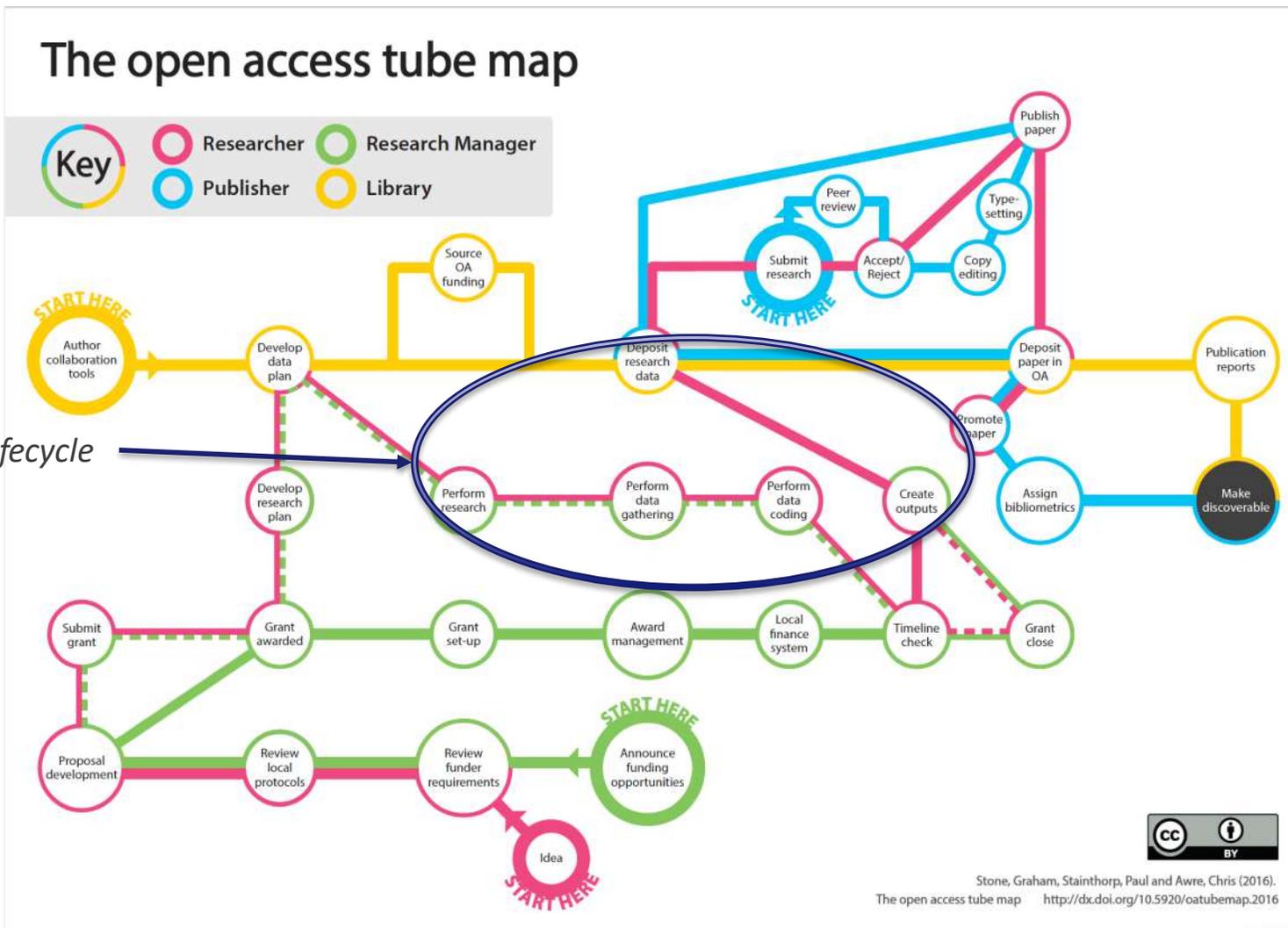
JEL: E60 - General, E62 - Fiscal Policy, E65 - Studies of Particular Policy Episodes

Issue Section: Article

Overcoming the Problem – 1. The Research Object Paradigm

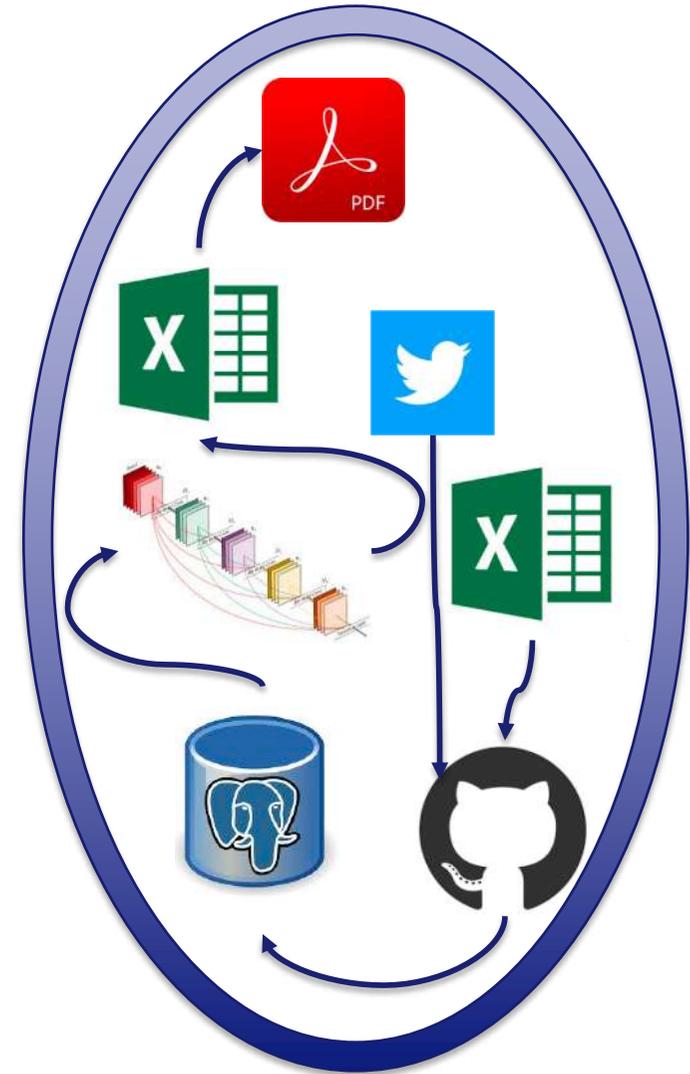


Taking the RO to the Next Level



Making the RO Usable

- Linked Open Data and RDF == **PROVENANCE**
- Open Questions...
 - Q: How much detail should I go into
 - A: How reproducible do you want your results to be?
 - Q: This is a lot of effort, why should I do it?
 - A: How else can a user trust your data?
 - Q: My observational outputs are not reproducible – do I still need provenance
 - A: Yes – need to know how you gathered your results
- Exposing Provenance
 - **Make it a part of your metadata**



Metadata – What is it?

- Data about Data
 - => Metadata IS data
- A way of describing a research artefact in a structured, machine readable way
- Something my community or funder insist I have
- “A love note to the future” ❤️
- A way of tracking your every movement...

The 'Me' in Metadata

Almost every digital file we generate carries invisible tags.

METADATA EXAMPLES

- Focal length
- Camera type
- Date & time taken
- Exposure
- Flash setting
- Preferred language
- Home location
- Place ID
- Internet provider
- Mail client
- IP address

From the tweet:

- COORDINATES: 42.59640 -114.4012
- LANGUAGE: EN

From the email:

- DATE/TIME: 2013-04-22 15:57:33

From the photo:

- CAMERA TYPE: HTC ONE X
- FLASH SETTING: AUTO

- 1 Geoff poses by a waterfall and snaps a self-portrait, which he immediately tweets, then emails to his grandmother.
- 2 Geoff's text, photo, and email ascend to a series of remote servers, each dragging their own trails of metadata.
- 3 Once there, the metadata may be extracted and interpreted by any interested party with access.

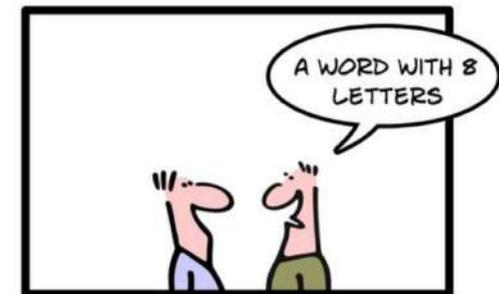
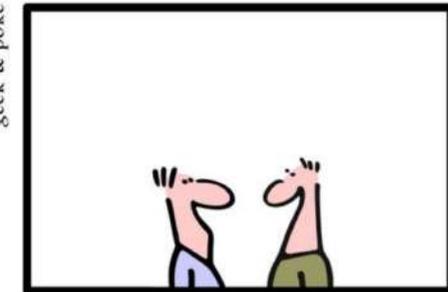
Thanks to the above metadata, without ever having met Geoff, we know he was at Shoshone Falls near Twin Falls, Idaho, at 3:57 p.m. on April 22, that he has an HTC One X smartphone, and that he is an English language speaker.

Note: Actual metadata code modified for readability. Many metadata values, such as time, can be extracted from multiple sources, and values may differ slightly. Source: staff reporting The Wall Street Journal

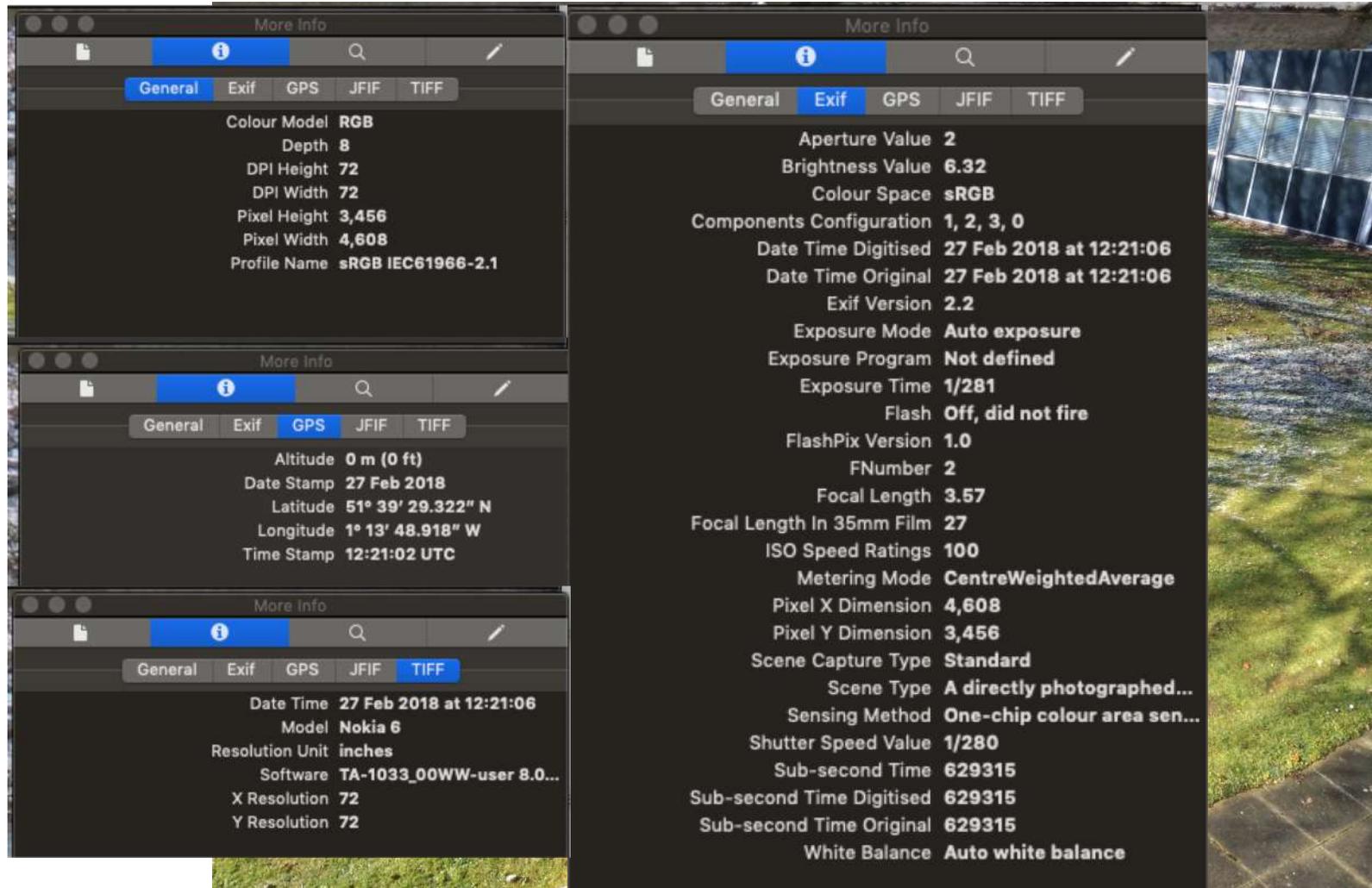
SIMPLY EXPLAINED: METADATA



geek & poke



Metadata in Life



The image shows a photo viewer interface with three metadata panels on the left and a large photo on the right. The photo shows a building facade with a grid of windows and a grassy area in the foreground.

General Metadata:

- Colour Model: **RGB**
- Depth: **8**
- DPI Height: **72**
- DPI Width: **72**
- Pixel Height: **3,456**
- Pixel Width: **4,608**
- Profile Name: **sRGB IEC61966-2.1**

GPS Metadata:

- Altitude: **0 m (0 ft)**
- Date Stamp: **27 Feb 2018**
- Latitude: **51° 39' 29.322" N**
- Longitude: **1° 13' 48.918" W**
- Time Stamp: **12:21:02 UTC**

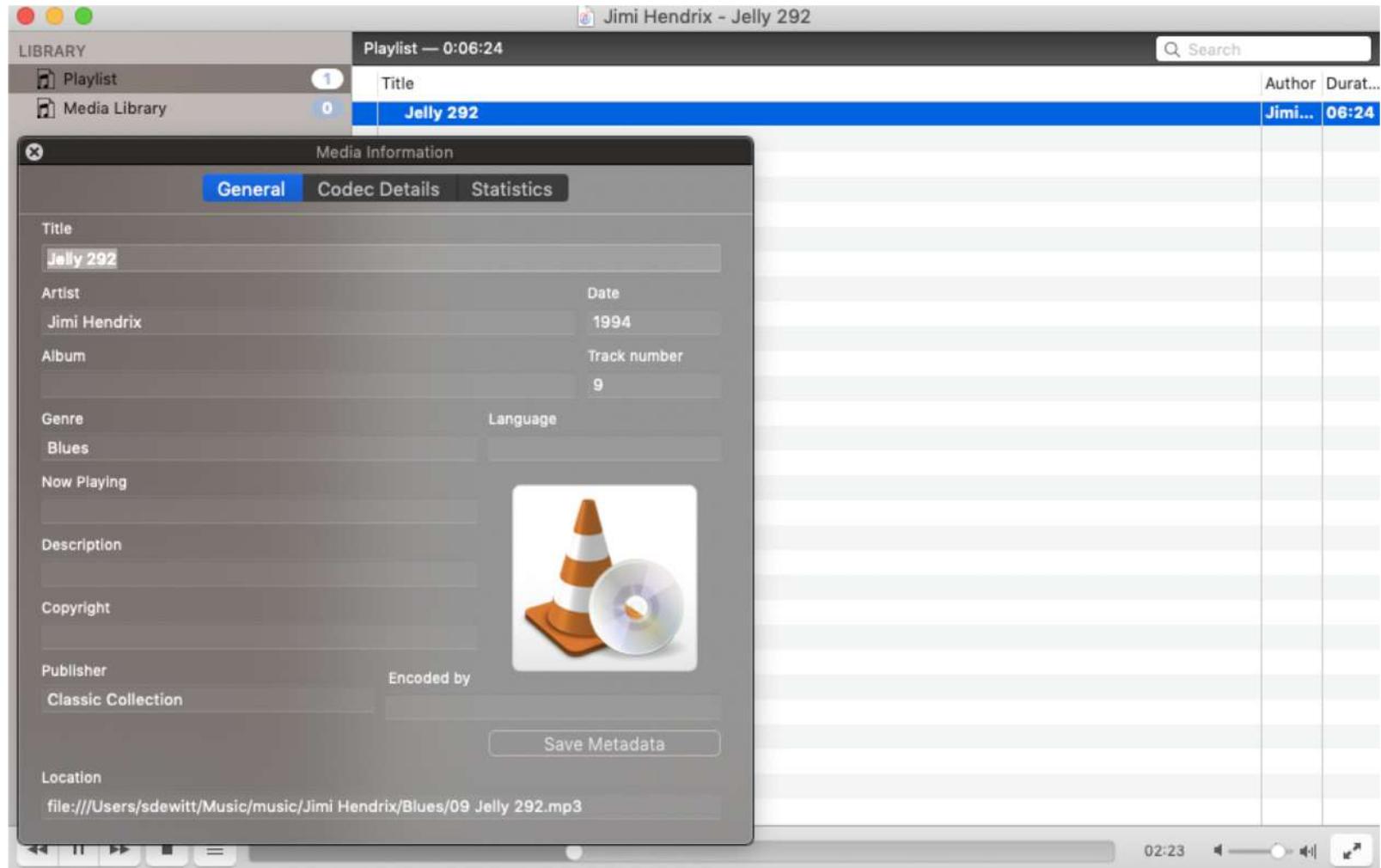
TIFF Metadata:

- Date Time: **27 Feb 2018 at 12:21:06**
- Model: **Nokia 6**
- Resolution Unit: **inches**
- Software: **TA-1033_00WW-user 8.0...**
- X Resolution: **72**
- Y Resolution: **72**

Exif Metadata:

- Aperture Value: **2**
- Brightness Value: **6.32**
- Colour Space: **sRGB**
- Components Configuration: **1, 2, 3, 0**
- Date Time Digitised: **27 Feb 2018 at 12:21:06**
- Date Time Original: **27 Feb 2018 at 12:21:06**
- Exif Version: **2.2**
- Exposure Mode: **Auto exposure**
- Exposure Program: **Not defined**
- Exposure Time: **1/281**
- Flash: **Off, did not fire**
- FlashPix Version: **1.0**
- FNumber: **2**
- Focal Length: **3.57**
- Focal Length In 35mm Film: **27**
- ISO Speed Ratings: **100**
- Metering Mode: **CentreWeightedAverage**
- Pixel X Dimension: **4,608**
- Pixel Y Dimension: **3,456**
- Scene Capture Type: **Standard**
- Scene Type: **A directly photographed...**
- Sensing Method: **One-chip colour area sen...**
- Shutter Speed Value: **1/280**
- Sub-second Time: **629315**
- Sub-second Time Digitised: **629315**
- Sub-second Time Original: **629315**
- White Balance: **Auto white balance**

Metadata in Life



Metadata Standards - Simple

Dublin Core

- Title
- Subject
- Description
- Creator
- Publisher
- Contributor
- Date
- Type
- Format
- Identifier
- Source
- Language
- Relation
- Coverage
- Rights

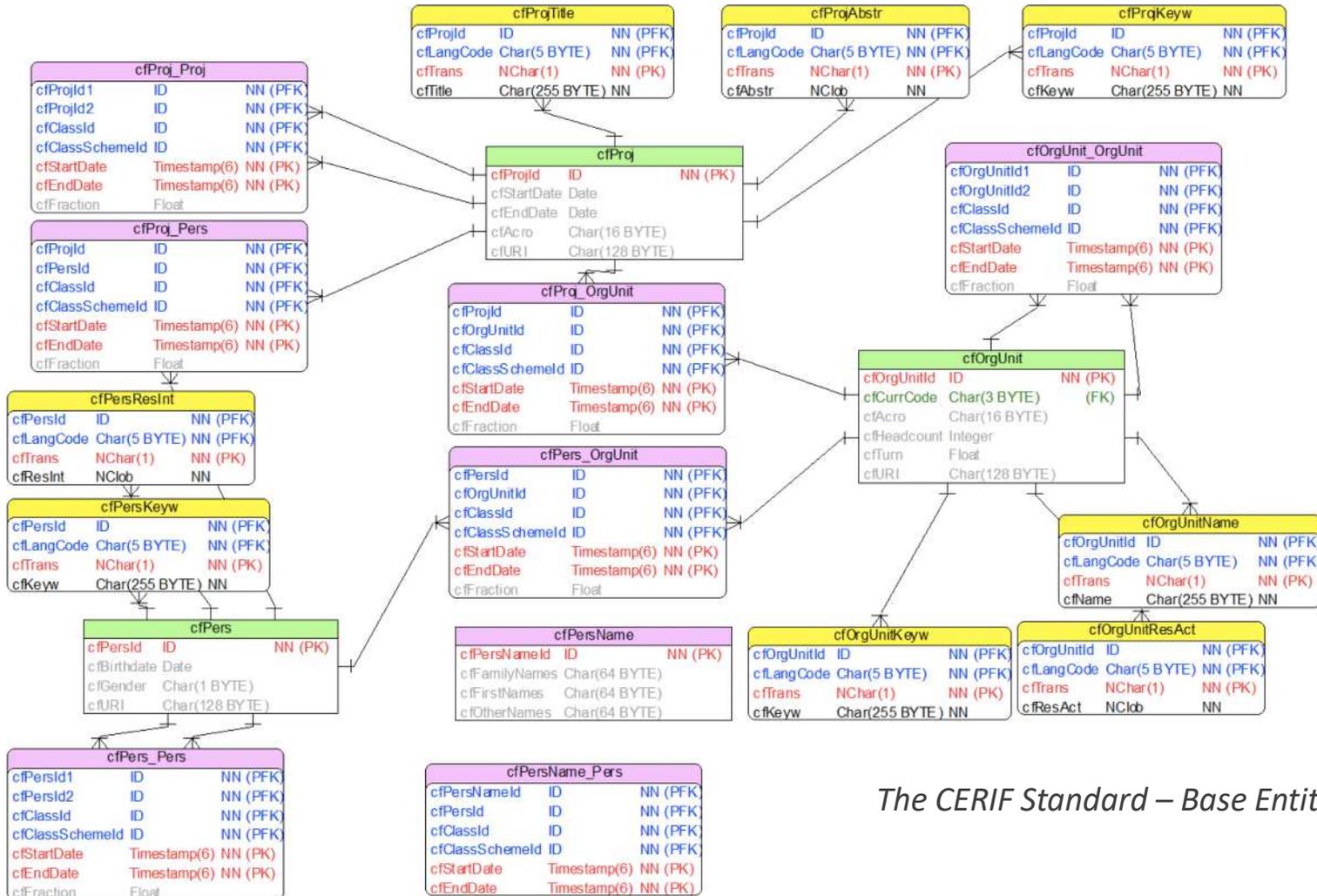
DataCite

- Title
- Creator
- Publisher
- Identifier
- Publication Year
- Resource Type
- Subject
- Contributor
- Date
- Related identifier
- Description
- Geolocation
- Language
- Alternate identifier
- Size
- Format
- Version
- Rights
- Funding Reference

EDMI

- Name
- Description
- Identifier
- url
- Creator
- Date Created
- license
- Data Standard
- Date Modified
- Access URL
- Access Interface
- Structure
- Included In
- Measurement Technique
- Keywords
- Variable Measured
- Format
- Scientific Type
- Includes
- Content Type
- Size
- Authentications
- Version
- Metric
- Same as
- Spatial Coverage
- Temporal coverage
- Citation
- Reference citation
- compression

Metadata Standards - Complex



The CERIF Standard – Base Entities

Exercises in metadata – Guess the film!

<i>Subject</i>	<i>None</i>
<i>Description</i>	<i>Three men searching for stolen gold</i>
<i>Creator</i>	<i>Age & Scarpelli, Luciano Vincenzoni</i>
<i>Publisher</i>	<i>Produzioni Europee Associate</i>
<i>Date</i>	<i>23 December 1966</i>
<i>Type</i>	<i>Spaghetti Western</i>
<i>Format</i>	<i>35mm anamorphic</i>
<i>Identifier</i>	<i>None</i>
<i>Source</i>	<i>Original Work</i>
<i>Language</i>	<i>Italian, English, Spanish</i>
<i>Relation</i>	<i>Part of trilogy</i>
<i>Coverage</i>	<i>US Civil War</i>
<i>Rights</i>	<i>Produzioni Europee Associate</i>



Exercises in metadata – Guess the film!

<i>Series</i>	<i>None</i>
<i>Cast</i>	<i>Marlon Brando, Charlie Sheen, Robert Duvali, Dennis Hopper, Harrison Ford</i>
<i>Credits</i>	<i>Director: Francis Ford Coppola, Writer: John Millus</i>
<i>Country</i>	<i>USA</i>
<i>Format</i>	<i>35mm film</i>
<i>Length</i>	<i>4205m</i>
<i>Duration</i>	<i>2hrs 27 mins</i>
<i>Language</i>	<i>English, French, Vietnamese</i>
<i>Year</i>	<i>1979</i>
<i>Identifier</i>	<i>Italian, English, Spanish</i>
<i>Genre</i>	<i>Part of trilogy</i>
<i>Relation</i>	<i>US Civil War</i>
<i>Source</i>	<i>Produzioni Europee Associate</i>



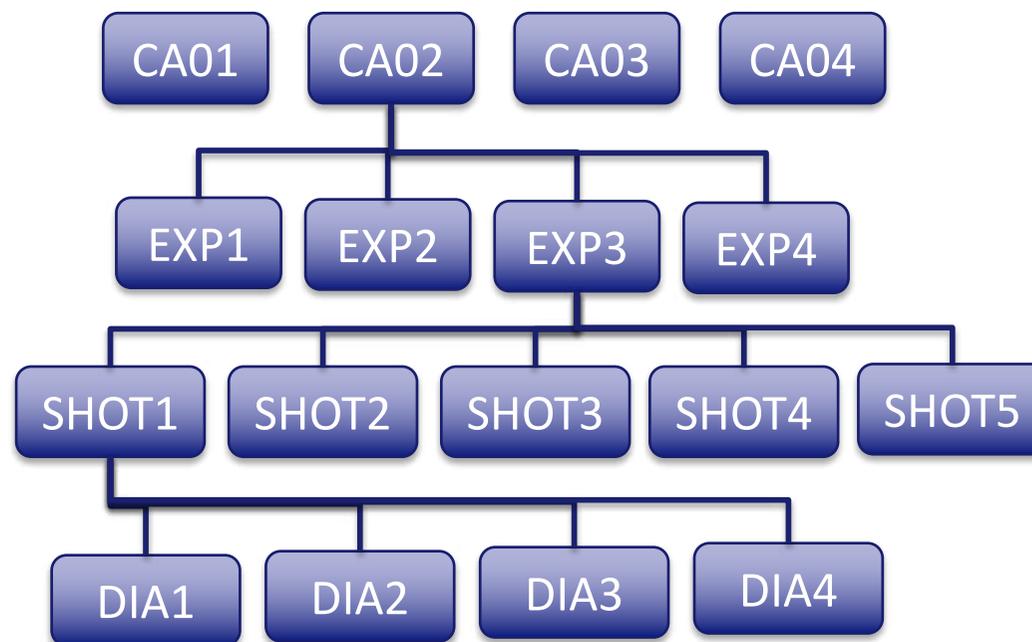
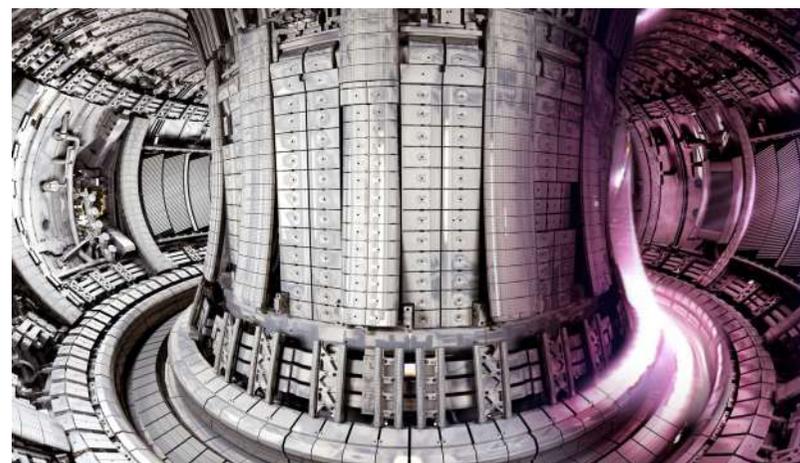
Exercises in Metadata

- http://nsteffel.github.io/dublin_core_generator/
- Use the Simple metadata generator to start with and to add metadata about yourself!
- If you get time – have a look at Advanced Generator but not essential



Metadata – the Last Word(s)

- FAIR
 - 13 of the 15 FAIR principles mention metadata
 - **It is as important as the data**
- Problems
 - Defining granularity
 - Multi-lingualism
 - Integrating existing schemas
 - Older Standards
- Metadata from different sources can be aggregated
 - Improve interdisciplinary science



Persistent Identifiers

- A way of giving your data a location independent link
- Means if someone cites data, the citation remains the same even if the data underneath moves
- Can be used to link the data with the metadata (in a PK-FK type relationship)
- Various forms....
 - Archival Resource Key (ARK):
 - <http://bnf.fr/ark:/13030/tf5p30086k>
 - Persistent Uniform Resource Locators (PURL) – not to be confused with personalised URL
 - <https://archive.org/services/purl/purl/Redford-Physics-of-God>
 - DOI
 - <https://doi.org/10.1109/5.771073>
 - International Standard Name Identifier
 - <http://isni.org/isni/000000012146438X>
- But they all have the same purpose – give a **unique**, **consistent**, **permanent** and **resolvable** id to something... e-mails, addresses etc are ephemeral

Last name of author(s), initial(s) of first name(s) with period. Put & before last author. Separate authors with a comma.

Publication year in brackets.

Use sentence-case for title of article: capitalize first word of title and subtitle only. Capitalize proper nouns, such as companies or place names.

Use italics and capitalize all long words in the journal/source title.

Great – I can find the paper

Spitz, D., & Hunter, S. (2005). Contested codes: The social construction of Napster. *Information Society*,

21(3), 169-180. doi: 10.1080/01972240490951890

Italics for volume number, issue number in brackets.

First-page to last-page of article

Digital object identifier (doi) labelled with doi: and placed at end of citation. If needed, remove blue hyperlink.



Taylor & Francis Online
Access provided by Marie Curie Library - International
Journal
The Information Society
An International Journal
Volume 21, 2005 - Issue 3
Submit an article | Journal homepage

Oh – there's one
And they don't



nature immunology
nature.com > journal home > advance online publication > article > abstract
ARTICLE PREVIEW
view full access options >
NATURE IMMUNOLOGY | ARTICLE
The significance of T cell inflammation
Chen Dong
Department of Immunology and Center for Inflammation and Cancer, MD Anderson Cancer Center, Houston, Texas, USA.
Tsinghua University School of Medicine, Beijing, China.
Contact Chen Dong
Search for this author in:
NPG journals . PubMed . Google Scholar
orcid.org/0000-0002-0084-9130
Lecturer / Senior Lecturer
Molecular Biology / Cell Biology
University of Southampton

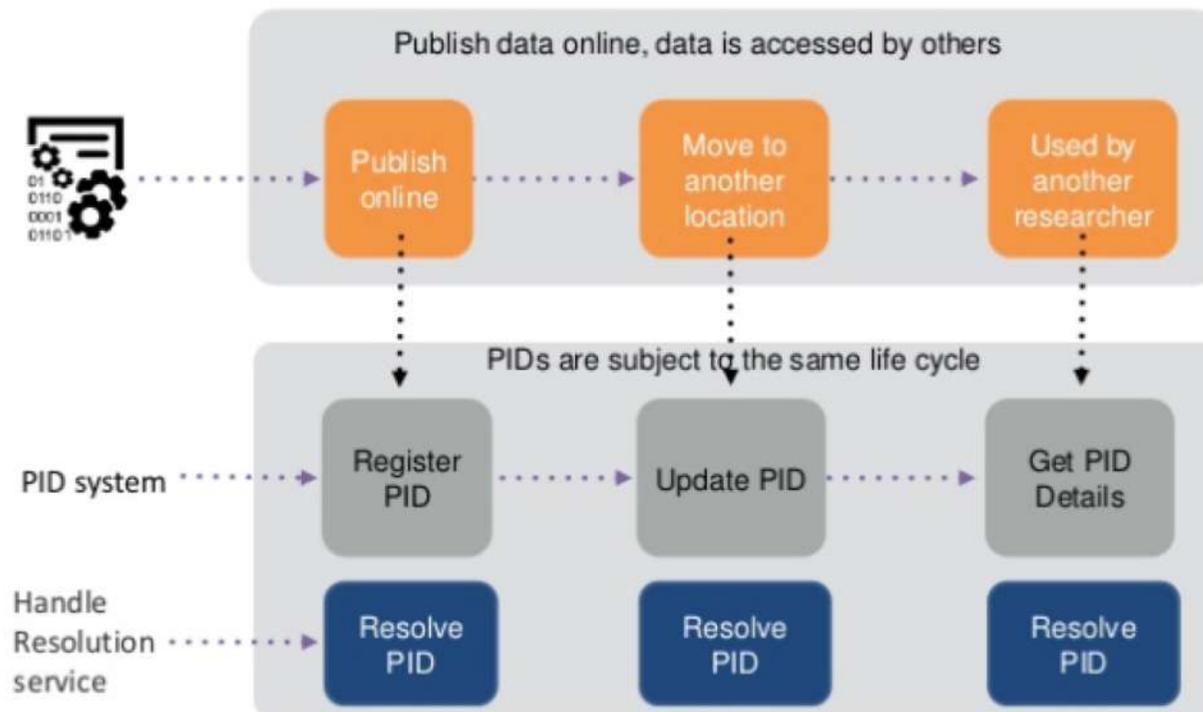
512 Views
17 CrossRef citations to date
0 Altmetric

Original Articles
Contested Codes: The Social Construction of Napster
David Spitz & Starling D. Hunter
Pages 169-180 | Received 04 Feb 2004, Accepted 20 Jun 2004, Published online: 23 Aug 2006

Download citation | <https://doi.org/10.1080/01972240490951890>

PIDs and Indirections

- All technologies rely on a resolver, and the resolution has to be kept updated – archivist/data manager/ data steward



Licensing – the last but most important thing you will do

- Why license
 - Legal protection
 - Allow users to understand what they can and con not do with the data
 - Make sure you get credit
- But which license?
 - B2SHARE can help – it has a tool to help you select your license



Examples of licenses commonly used for open data

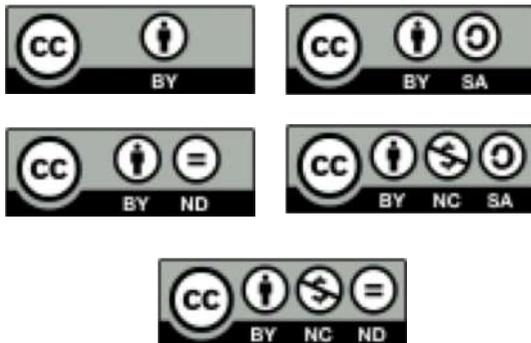
- Open licenses

- CC-0
- WTFPL
- Unlicense
- PDDL



- Permissive Licenses

- BSD
- MIT
- Creative Commons
- EUPL v1.2

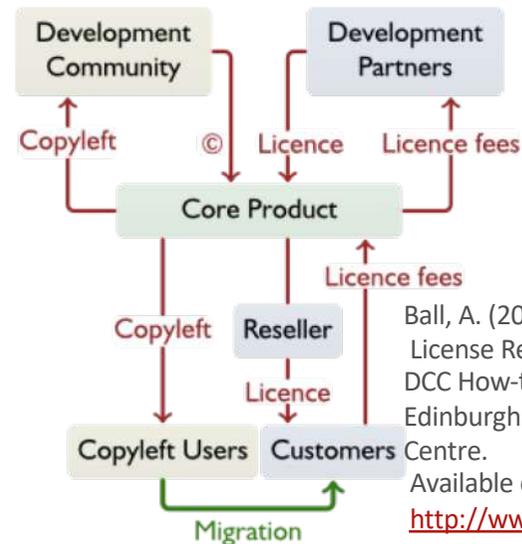


- Copyleft license

- GPL, LGPL, AGPL
- Common Development & Distribution License

- Specific Licenses

- IES Restricted Data License
- MetaShare no Redistribution License

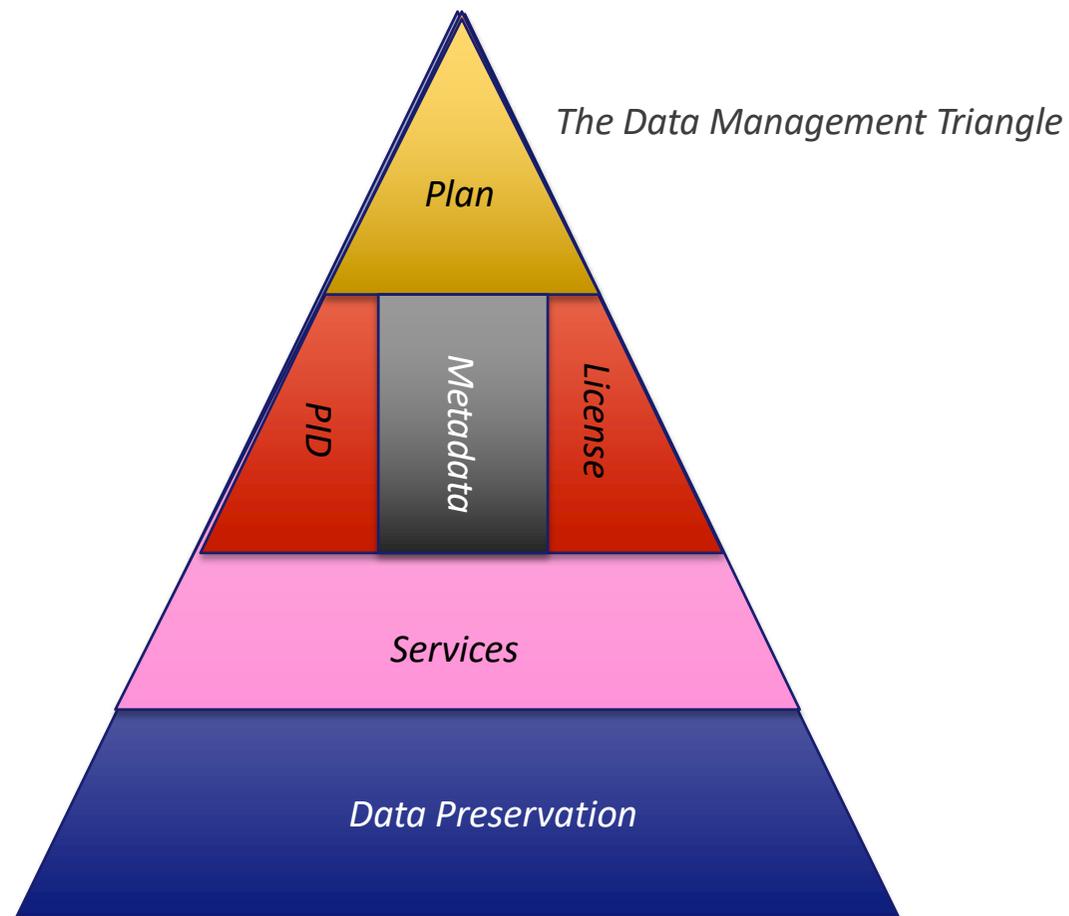
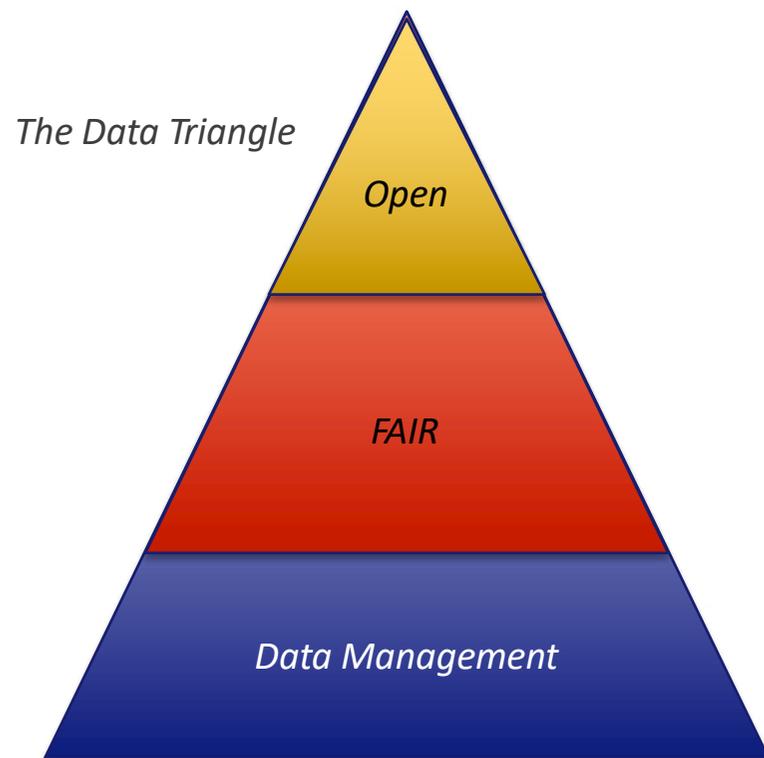


Ball, A. (2014). 'How to License Research Data'. DCC How-to Guides. Edinburgh: Digital Curation Centre. Available online: <http://www.dcc.ac.uk/resources/how-guides>

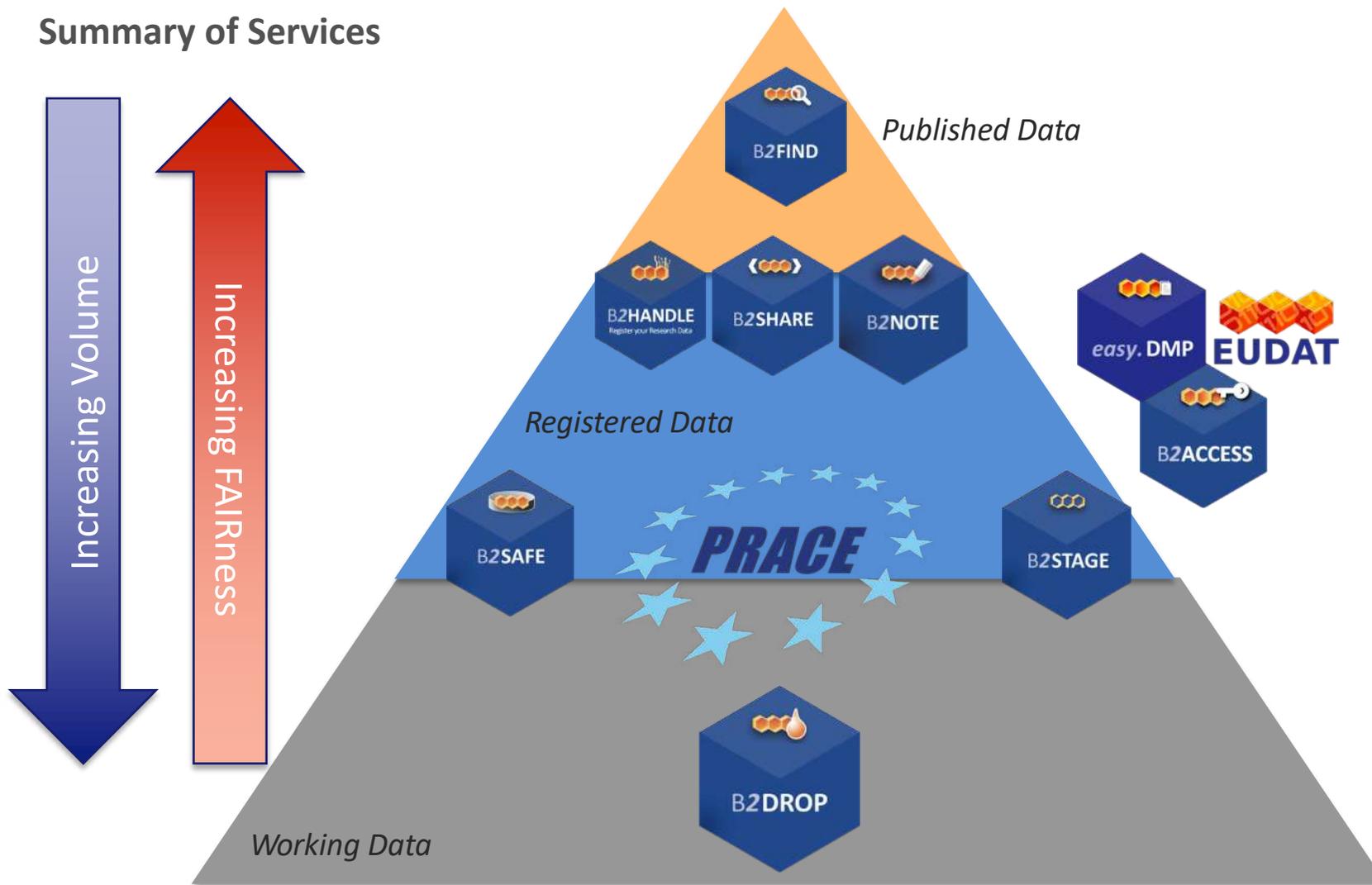
Summary

- **It is ridiculously easy to lose data**
- **Data Preservation** is a subset of **Information Preservation**
 - With overlap when **common standards** are used
- Various **Technologies** and **Procedures** can help you with data preservation
 - **None** are foolproof
 - AWS offer **99.999% availability** annual failure rate (AFR) of between **0.1%-0.2%**.
 - “We have no liability whatsoever for any damages, liabilities, losses (including any corruption, deletion, or destruction or loss of data, applications or profits)...”
- **Checksumming** (a.k.a. hashing or fixity checking) used to detect changes in files
 - Very good but can be subject to malicious attack
- **Metadata** and **Persistent Identifiers** needed for both data preservation and information preservation
 - **Provenance** for repeatability and trust

Summary...



Summary of Services



The Rest of The Week

- Introduction to some of the services...
 - B2DROP and B2SHARE for sharing your data and making it public
 - B2FIND and making your data discoverable
 - GridFTP – the interface between EUDAT and PRACE
- Taking you through a real use case from ENES Climate Analysis Service (ECAS)

