



ILLUMINATING DARK DATA

WHY YOU SHOULD STOP STORING DATA

ANTIQUES ROADSHOW



**CASH
IN THE
ATTIC**



Survey responses indicate The typical US researcher:



Publishes about
60% of their
publishable data



Has 10-24
unpublished
data sets



Has 300-500
unused
samples



Has \$29,000 in
unused resources

Unused resources at US Academic institutions:



Social Sciences
8%, \$0.5 B



Physical Sciences
10%, \$0.6 B



Engineering
19%, \$1.2 B



Life Sciences
63%, \$3.9 B



\$6.2
Billion

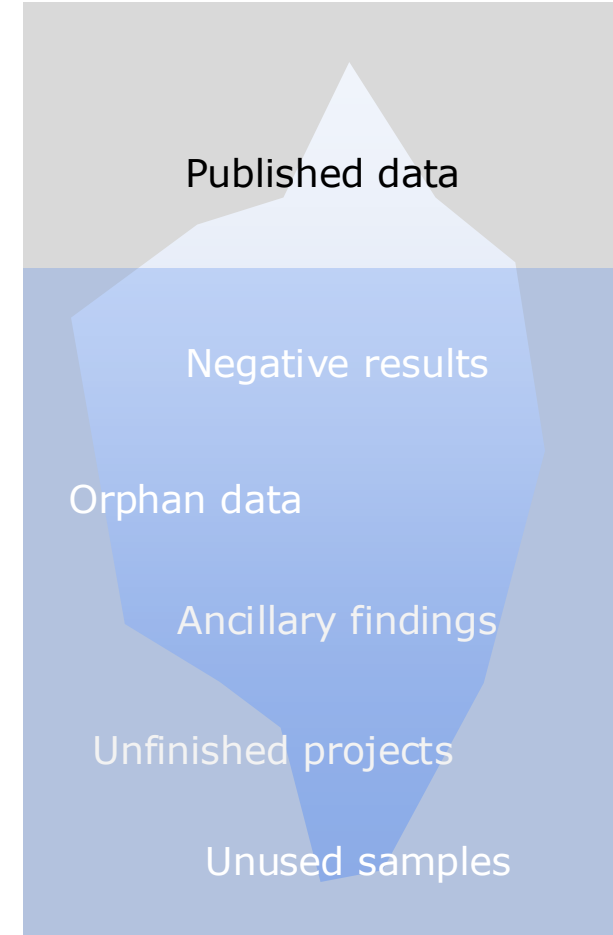
Highlights

Unpublished data and
unused samples represent
billions of dollars in
stranded assets

Unfinished projects and
orphan data were the
primary sources of
unpublished data

Unused samples are not
shared because others
don't know they exist

Addressing logistical
problems can increase the
efficiency of research



WHAT IS DARK DATA?

- Researchers need to ensure their data and metadata comply with existing standards to be useful.
 - Incompatible data is essentially useless for further research.
- Researchers must allocate time to standardize their data and learn the current standards of repositories.
- Academia and institutions often do not support the time needed for these tasks.
- There are few incentives for researchers to share, format, and standardize data.
- This lack of support leads to reduced professional collaboration and wasted research efforts.
- Valuable scientific data can become forgotten and underutilized, despite significant costs associated with its creation and storage.
- Such neglected data is referred to as "dark data" due to its lack of visibility in research.

Are you storing dark data?



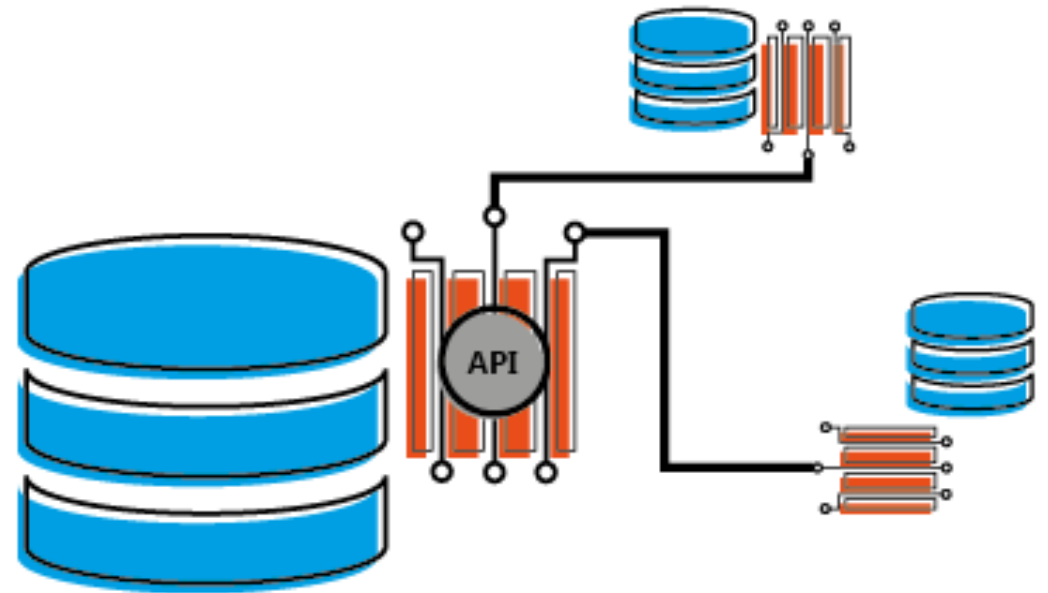
‘Dark data’ is killing the planet – we need digital decarbonisation

Published: September 29, 2022 6.12pm CEST

MAIN IDEA – DISTRIBUTE AND STANDARDISE

- Observation: Data will always reside in various places, and is usually slow to move
 - Also moving usually means copying, creating different versions
- Solution: leave data where it is, but make it accessible via *interfaces* (API's) so that users and applications can access them when needed
 - Improves control, findability and security

Stop storing data, and start *managing it!*



FAIR D



https://youtu.be/KWVCSwUNtBA?si=vqL-Q88JDGj_4pRS

FAIR PRINCIPLES

FINDABLE



Persistent unique identifiers enriched with metadata allow (meta) data to be searched for and found.

ACCESSIBLE



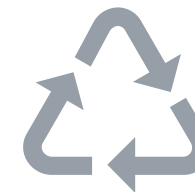
The (meta) data is accessible through standard communication (internet) protocols and available for research. Metadata is always accessible.

INTER-OPERABLE



Commonly used and open data formats, programming language, and vocabularies are used to allow easy exchange, integration, and reuse with other data and software.

REUSABLE



All data is clearly and elaborately documented so that the data can be correctly interpreted and reanalyzed (by others) under specified license.



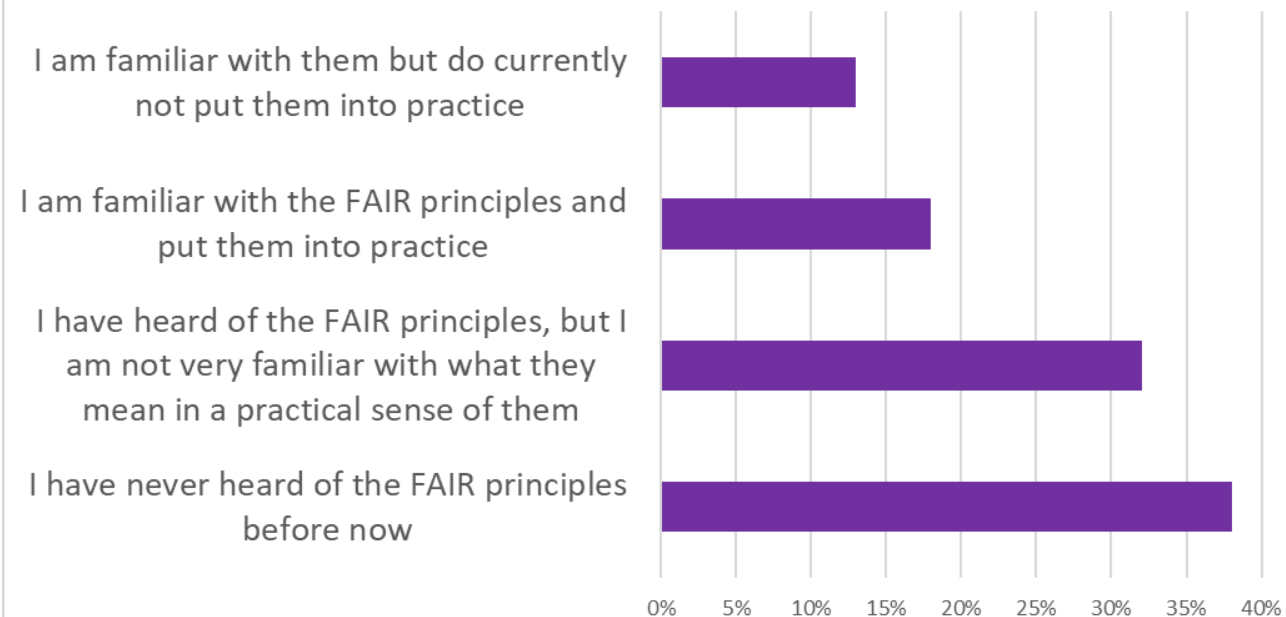
Expectation



Reality

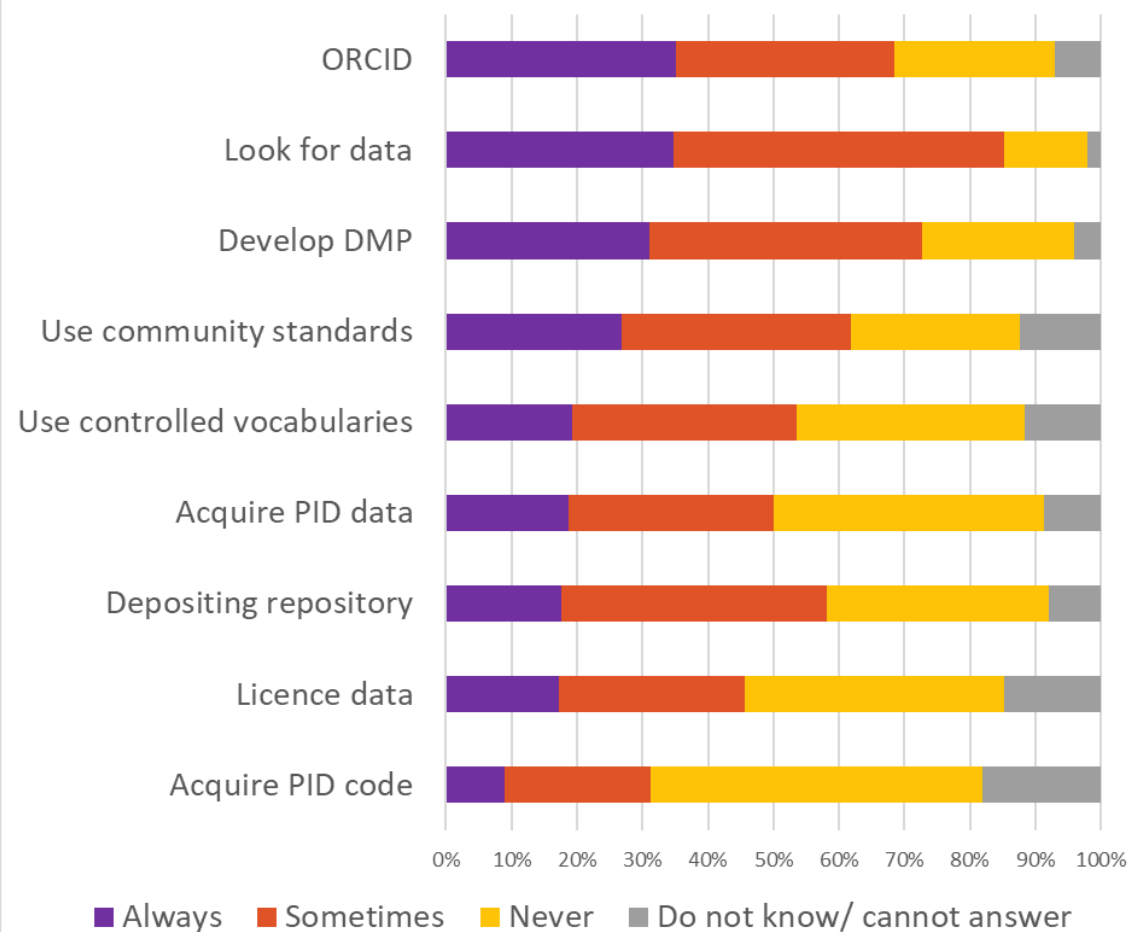
LACK OF ADOPTION

FAIR principles in EU researchers



<https://zenodo.org/doi/10.5281/zenodo.6778743>

FAIR practices in research



WHY A LACK OF ADOPTION?

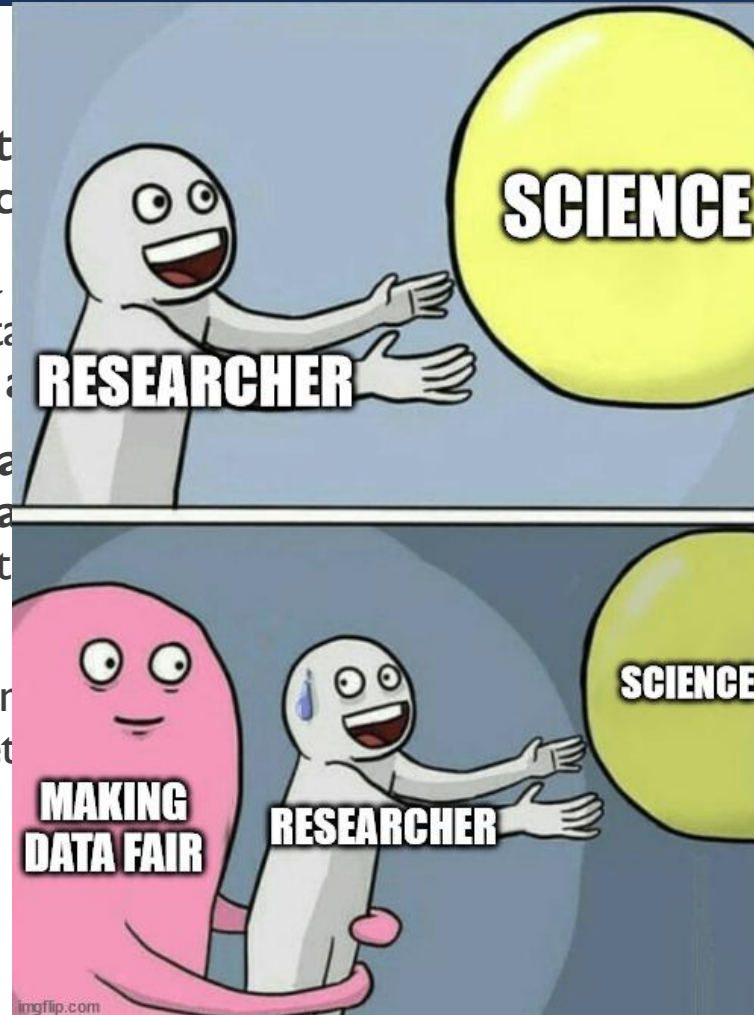
Schembera, B., Durán, J.M. (2019), “**Dark Data as the New Challenge for Big Data Science and the Introduction of the Scientific Data Officer**”:

- Researchers who wish to share their data need to make sure that both the data and the metadata they are submitting are compatible with existing standards. **Incompatible data might be as good as noise, for it cannot be processed for further use and is therefore useless for research.**
- **Researchers must find time in their already busy agenda not only for standardizing their data and metadata but also for acquiring updated knowledge on the standards used by a given repository.** Unfortunately, academia and other institutions are not always receptive to their researchers spending qualitative time on such endeavors.
- As a result, these unacknowledged—but implicitly required—efforts give researchers very few incentives to share, format, and standardize their data and metadata for further use, with the subsequent loss of professional collaboration and research efforts.

WHY A LACK OF ADOPTION?

Schembera, B., Durán, J.M. (2019), “Dark Data: The Introduction of the Scientific Data Office”

- Researchers who wish to share their data and the metadata they are submitting are compatible with existing standards, but their data cannot be processed for further use.
- Researchers must find time in their schedules to standardize their data and metadata but also for acquiring updates. Unfortunately, academia and other institutions do not spend time on such endeavors.
- As a result, these unacknowledged—but important—barriers to data format, and standardize their data and metadata, leading to a loss of collaboration and research efforts.



data Science and the

and the metadata they are
as good as noise, for it
ch.

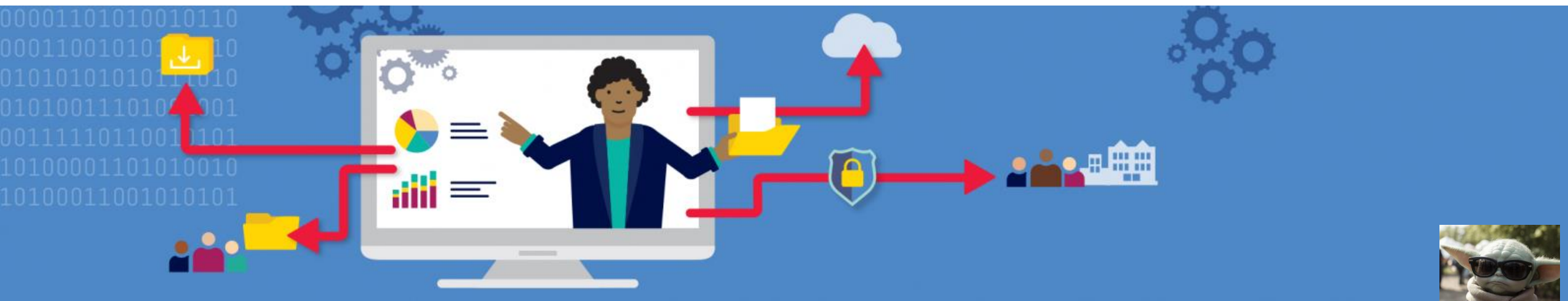
standardizing their data and
used by a given repository.
researchers spending qualitative

ers very few incentives to share,
ent loss of professional

WHAT IS YODA*?

- System for FAIR research data management
- Based on iRODS
 - High performance scientific storage system
- YODA is a RDM focused GUI *on top* of iRODS

*Yoda is a work of data engineering. Any similarity to specific space opera characters, living or dead, or galactic events, is purely coincidental.



FEATURES

- Based on iRODS architecture and language.
- Yoda is the user-friendly interface on top of iRODS.
- Store and work with data in the Yoda Research environment.
- Archive/secure data in the Yoda Vault environment.
- Add metadata to folders.
- Easily share data within and outside the insititute.

The screenshot displays the Yoda Portal interface. At the top, there is a navigation bar with "Research" and "research-death-star" (with "Go to group manager" and "Go to vault" buttons). Below this, the folder "research-death-star" is shown with a "Secured" badge and action buttons for "Metadata", "Create Folder", "Upload", and "Actions". A table lists files and folders:

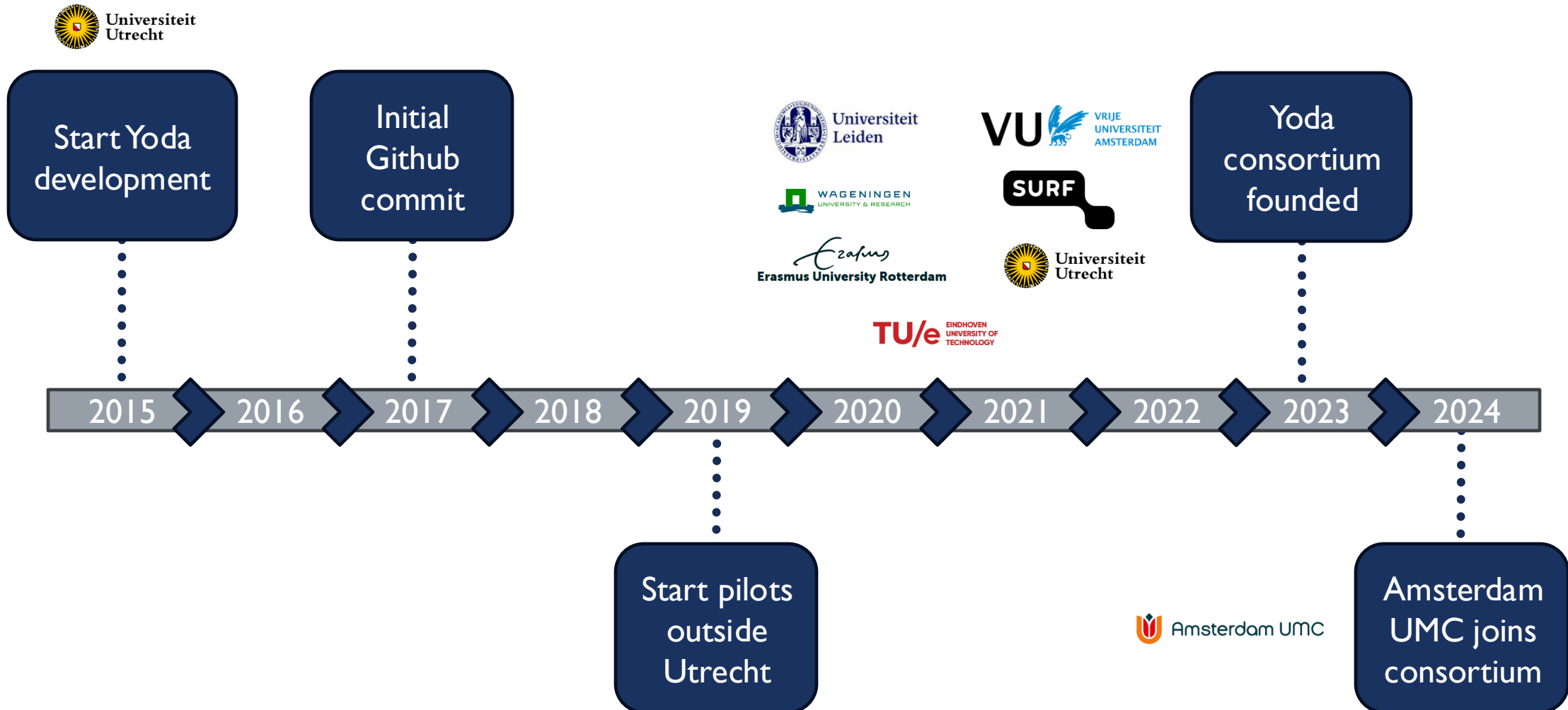
<input type="checkbox"/>	Name	Size	Modified date	
<input type="checkbox"/>	Cafeteria menus		2022-10-07	⋮
<input type="checkbox"/>	Personnel files			
<input type="checkbox"/>	Schematics			
<input type="checkbox"/>	blueprint.jpg			
<input type="checkbox"/>	transformation-back			
<input type="checkbox"/>	yoda-metadata.json			

A "Yoda Portal" overlay is visible, containing a "Metadata form - /research-death-star" window. The form includes a "Save" button, a "Required for the vault" progress indicator, and fields for "Title*" and "Description *".

Title*
Plans for a novel spacecraft with planetary destructive capabilities
The title of your data package

Description *
The Death Star is a mobile space station and galactic superweapon designed and operated by the galactic empire. Its dimensions are 160 kilometers (99 mi) in diameter, and it is crewed by an estimated 1.7 million military personnel and 400,000 droids.
A description of your data package

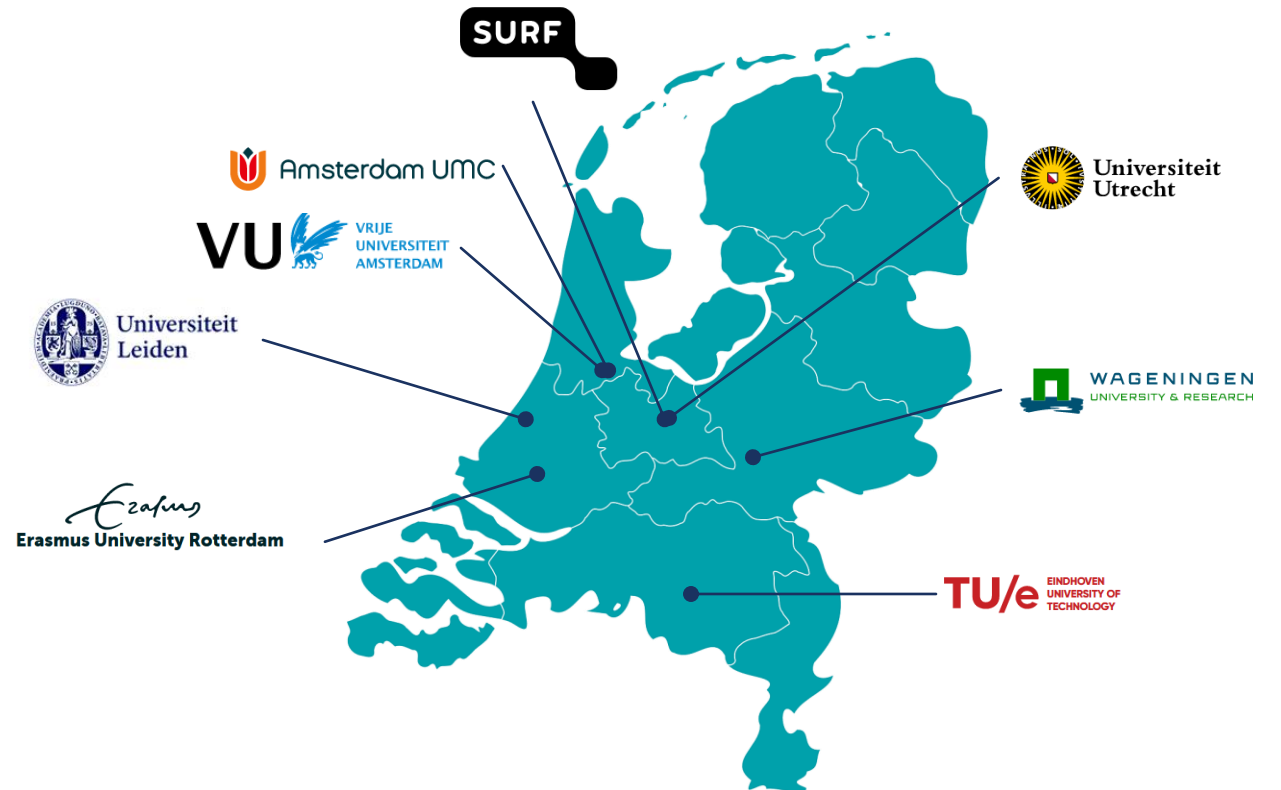
TIMELINE YODA & CONSORTIUM



YODA CONSORTIUM

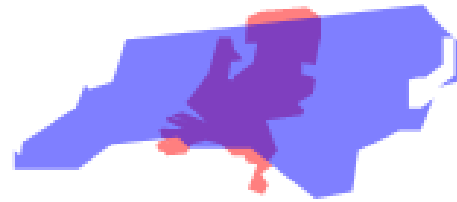
- 6 out of 13 major universities in the Netherlands & 1 UMC
 - Additionally, 3 uni's use iRODS based solutions
- Partially funds and governs feature development of the Yoda software
- Representation for hosting of Yoda by SURF

- 15,000+ users
- 3,5+ PB of data



YODA CONSORTIUM

- 6 out of 13 major universities in the N & I UMC
 - Additionally, 3 uni's use iRODS based sc
- Partially funds and governs feature dev of the Yoda software
- Representation for hosting of Yoda by :
- 15,000+ users
- 3,5+ PB of data



North Carolina
(US) is **3.36**
times as big as
Netherlands

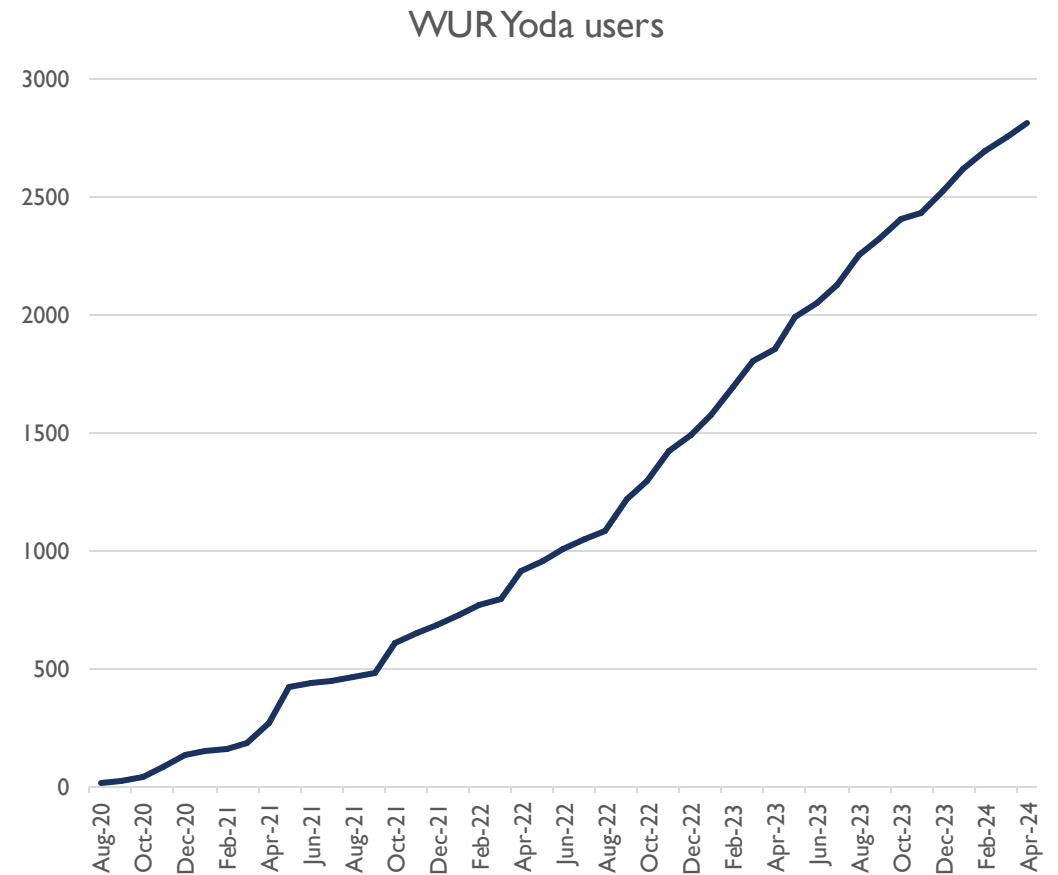
<https://mapfight.xyz/map/nl/>



BETTER TOOLS

- Growth of around 1000 users per year (within WUR)
- Similar growth in other consortium members
- User growth, without significant promotion or policy, indicates clear need
- Is part of their daily research process & environment

Yoda saves time and effort – it is practical RDM tooling

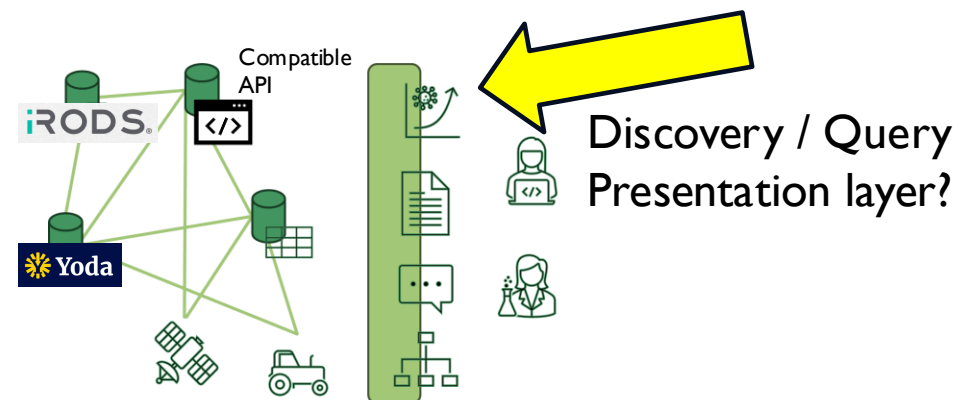
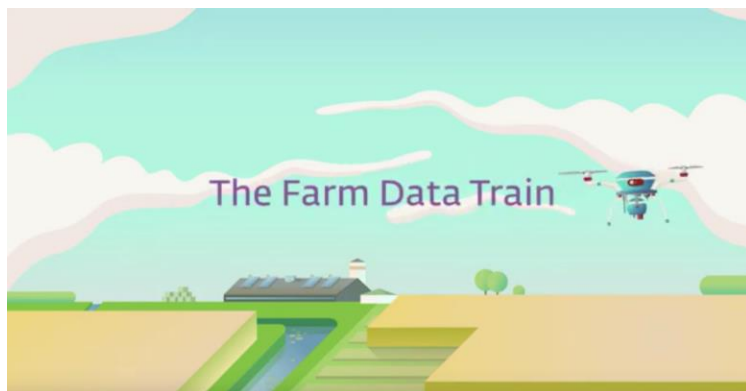


BETTER COLLABORATION



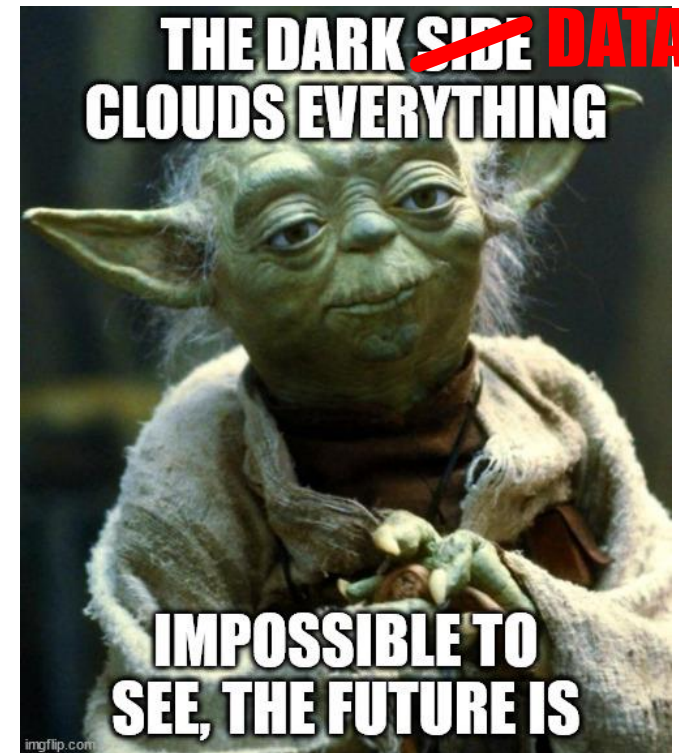
=

Dutch FAIR data train



RECAP & FUTURE OUTLOOK

- Dark data is a significant problem in –most- (Dutch) scientific institutions
- Yoda is a tool that supports researchers in their day-to-day data management activities, as well as enabling FAIRification of data
- Creates a large and active iRODS user base
- The Yoda consortium model creates a federated data network, without a single authority
- Expectation is that this model will work for most Dutch institutions
 - But also in a EU context (GaiaX, European Open Science Cloud and EU national activities)





THANKS!