

Three years of publishing data in ETH Zurich's Research Collection: Lessons learned and new developments

Barbara Hirschmann

Head E-Publishing, ETH Library

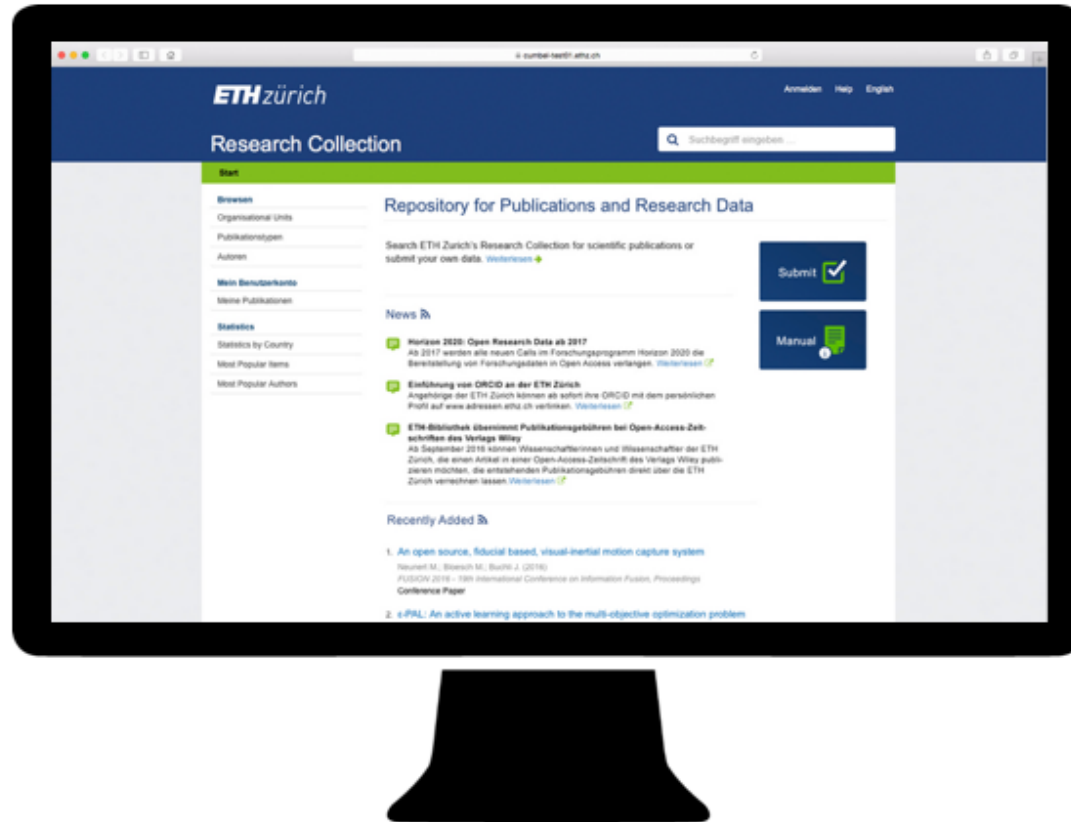
22 October 2020, Swiss Research Data Day

Agenda

1. Research Collection: Overview
2. How ETH researchers use the repository
3. Quality assurance and compliance checks
4. New developments

1. Research Collection: Overview

Research Collection: «3 in 1»



- Publications directory / Institutional bibliography



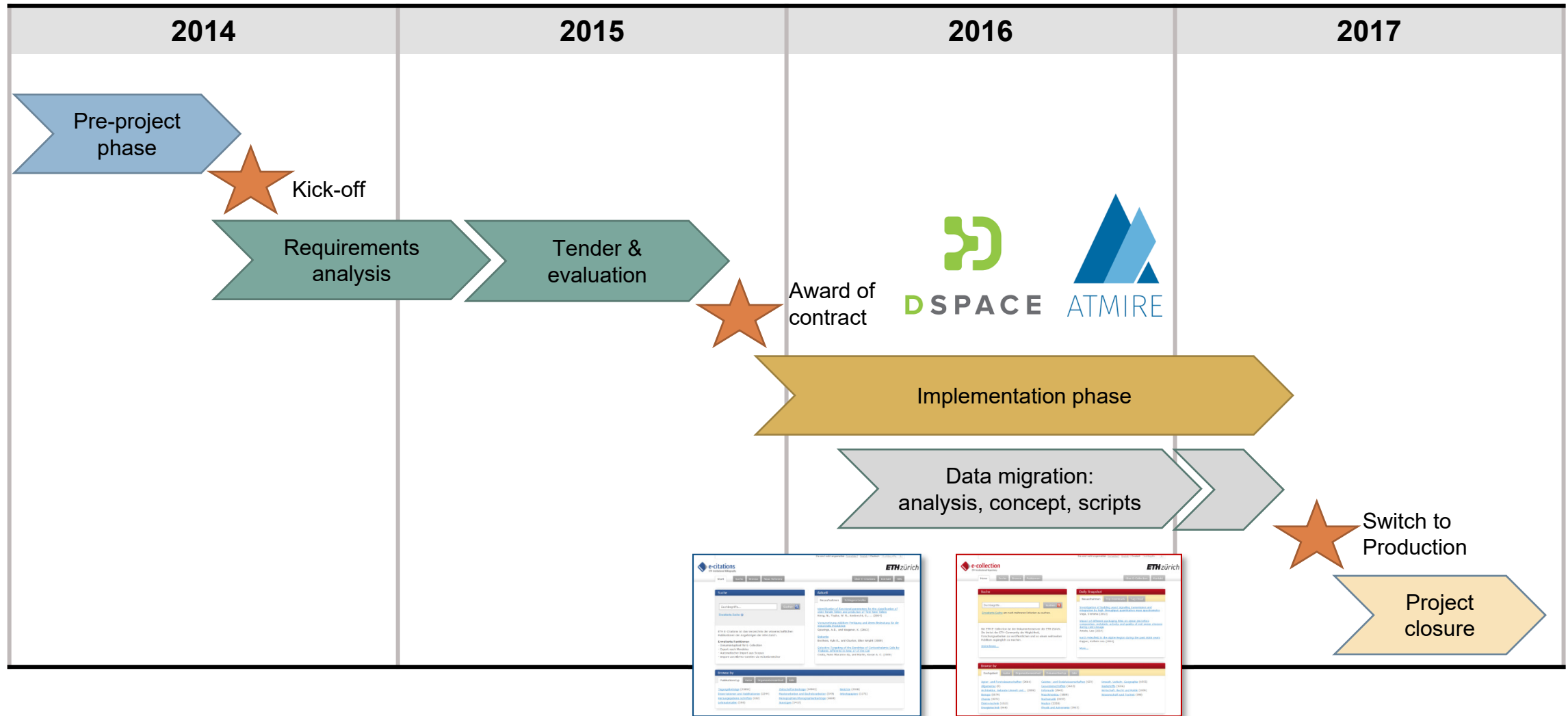
- Open access repository



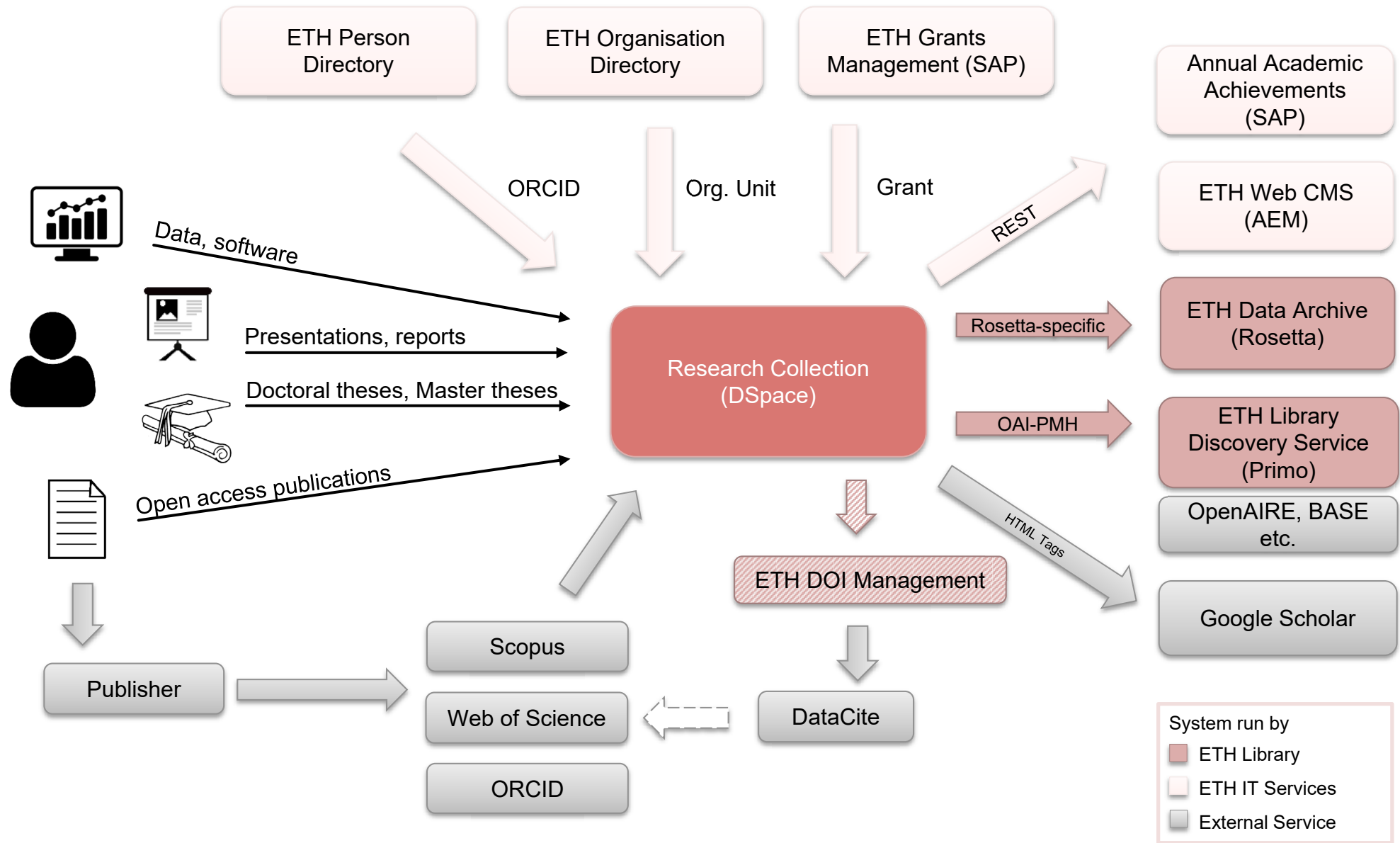
- Research data repository

www.research-collection.ethz.ch

Project timeline



System landscape



Features for publishing research data

- Deposit research data either as **supplementary material** or as **standalone** publication
- **Linking** between publications and research data
- Flexible **access rights**
- Option for **licensing** data under an open content license
- **DOIs** for research data; option to preview/reserve DOIs
- Usage **statistics** on file and item level
- ZIP- / TAR **preview**
- All **file formats** permitted
- Limited **retention period** possible
- Integration with **preservation system** ETH Data Archive (Rosetta)

The screenshot displays the ETH Zürich Research Collection interface. The top navigation bar includes the ETH Zürich logo, a search bar, and links for 'Login', 'Help', and 'Deutsch'. The main content area is titled 'Research Collection' and features a breadcrumb trail: 'Home → Research Data → Dataset → View Item'. On the left, there are several menu sections: 'Browse' (Organisational Units, Publication Types, Authors), 'Publish' (New Submission), and 'Statistics' (Downloads by Country, Most Popular Items, Most Popular Authors). The main content area displays the title 'Data for: Molecular Tracing of Riverine Soil Organic Matter From the Central Himalaya' with a database icon. Below the title, there are 'Open access' and 'CC BY' icons. The 'View/Open' section shows two files: 'tableS2.xlsx (MS Excel XML, 16.12Kb)' and 'tableS3.xlsx (MS Excel XML, 12.34Kb)'. The 'Rights / license' section indicates 'Creative Commons Attribution 4.0 International'. The 'Permanent link' is 'https://doi.org/10.3929/ethz-b-000431464'. The 'External links' section shows 'https://doi.org/10.1029/2020GL087403'. The 'Contributors' section lists 'Contact person: Märki, Lena' and 'Data collector: Märki, Lena'. The 'Publisher' is 'ETH Zurich'. The 'Organisational unit' is '02704 - Geologisches Institut / Geological Institute, 03868 - Eglinton, Timothy I. / Eglinton, Timothy I.'. The 'Related publications and datasets' section shows 'Is supplement to: https://doi.org/10.3929/ethz-b-000431470'. The 'More' section has a 'Show all metadata' link. The 'Atmetrics' section shows '1' tweet and a 'See more details' link.

Research data as independent publication type

Browse by Publication Type

Select a publication type to list all items of this type.

- Books [9812]
- Conference Contributions [40960]
- Educational Material [412]
- Journal Contributions [93550]
- Other Publications [3062]
- Patents [26]
- Presentations [1573]
- Research Data [275]
- Theses [26595]
- Working Papers/Reports [8590]



Research Data [275]

- Data Collection [31]
- Dataset [184]
- Image [13]
- Model [15]
- Other Research Data [5]
- Software [11]
- Sound [2]
- Video [14]

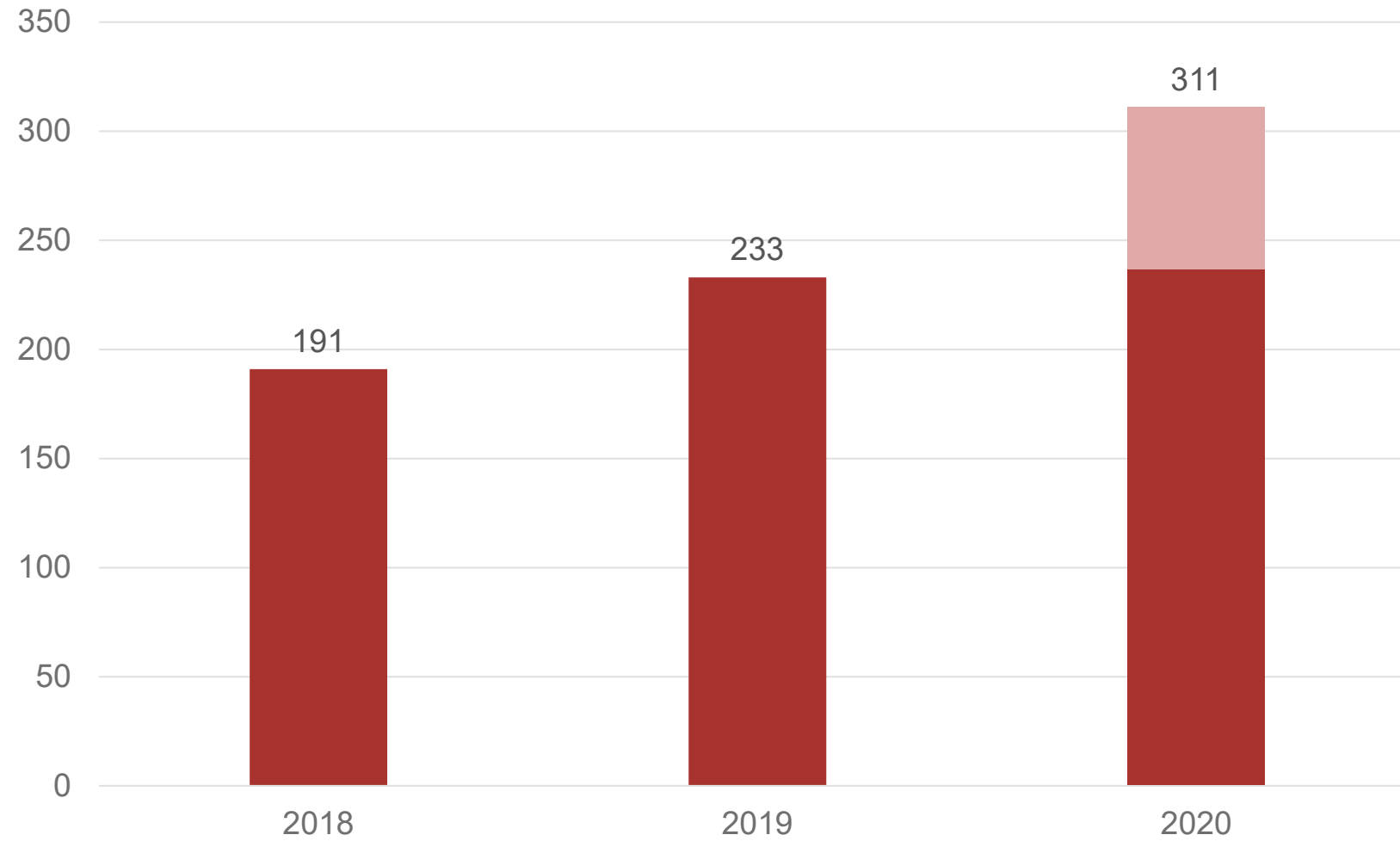
Access rights

	Open access	Embar- goed	ETHZ users	Selected users	Closed access
Publications	✓	✓			
Research data	✓	✓	✓	✓	✓

- Metadata are always freely accessible (via UI search and interfaces)
- Access to restricted datasets can be requested via a form

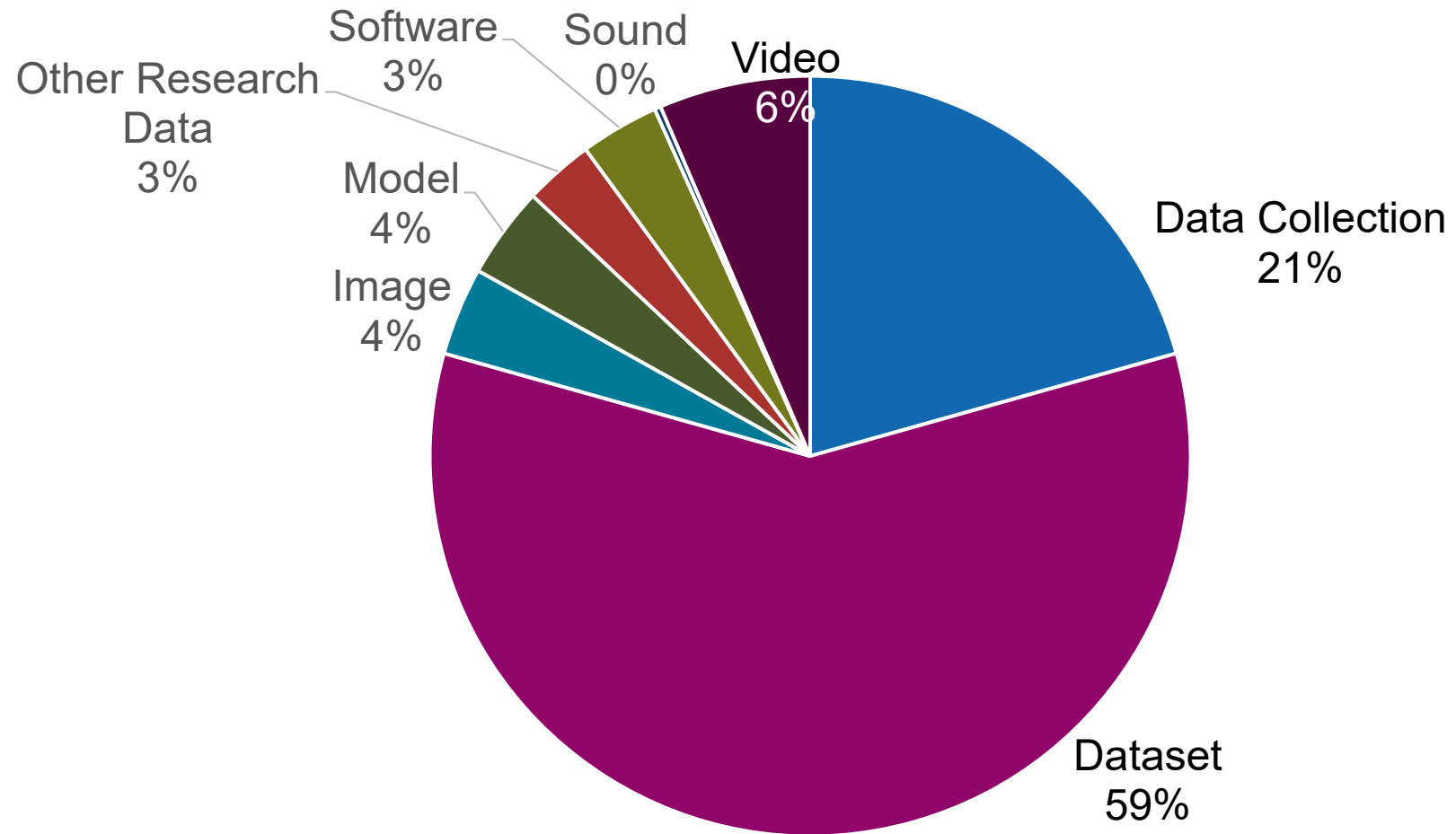
2. How ETH researchers use the repository

Number of published datasets per year

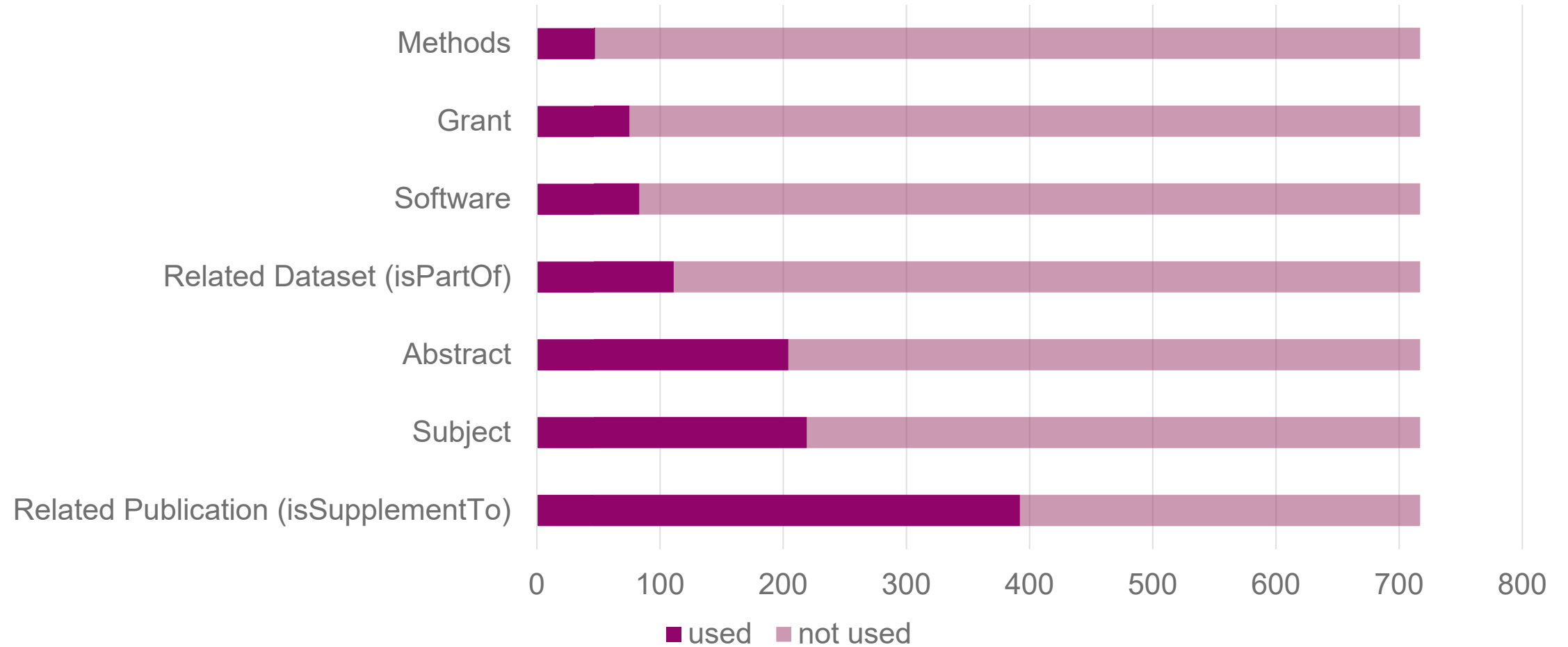


Total:
778 Datasets

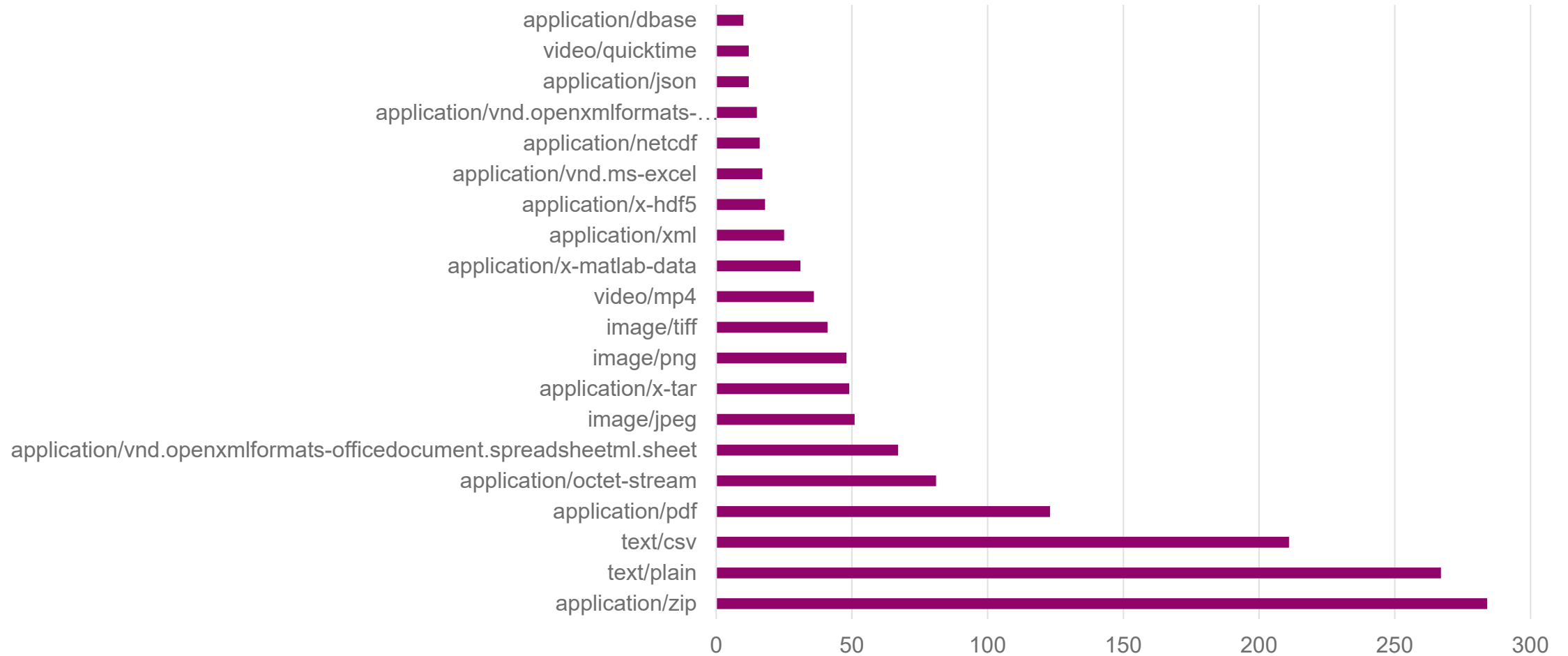
Dataset types



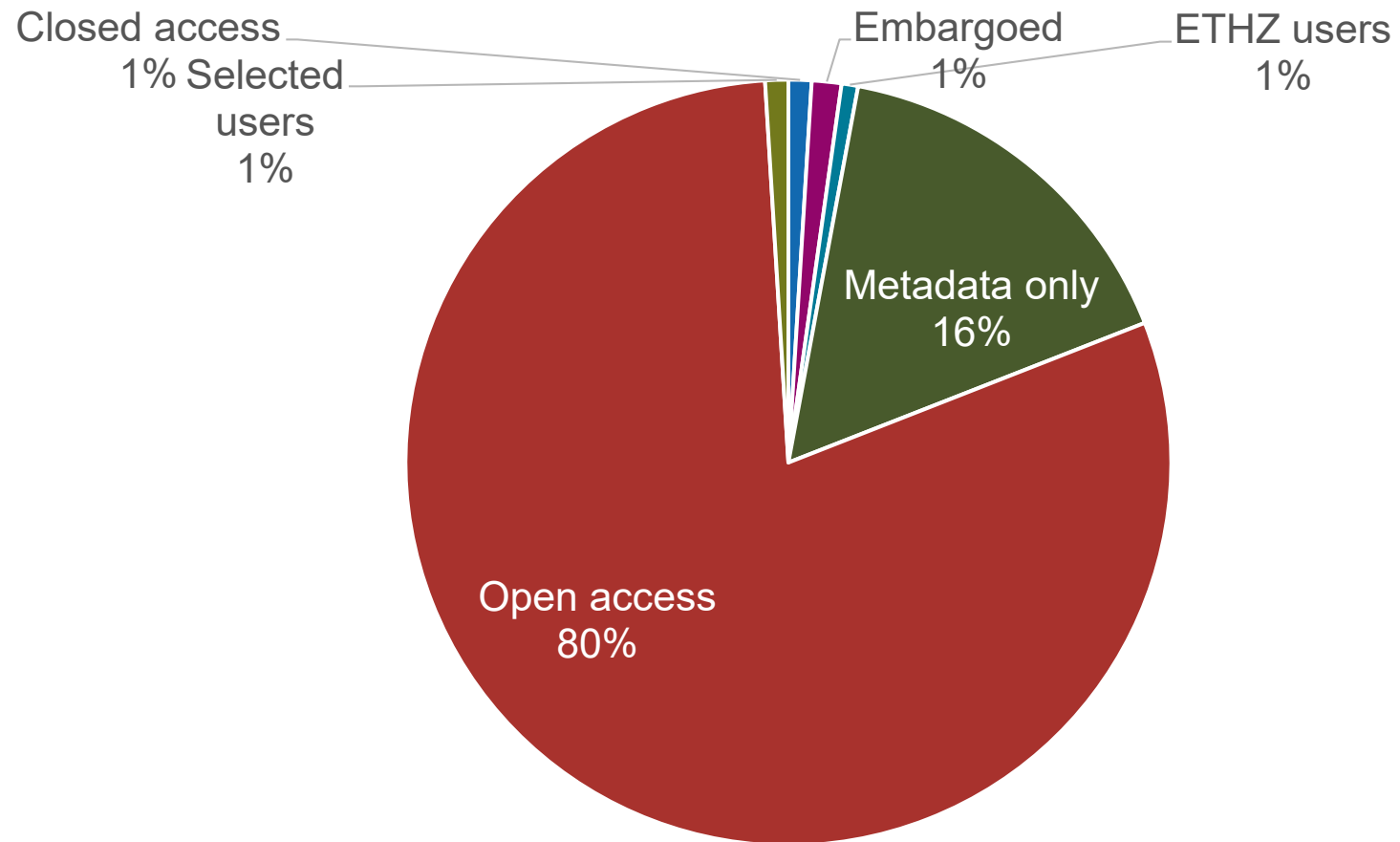
Optional metadata



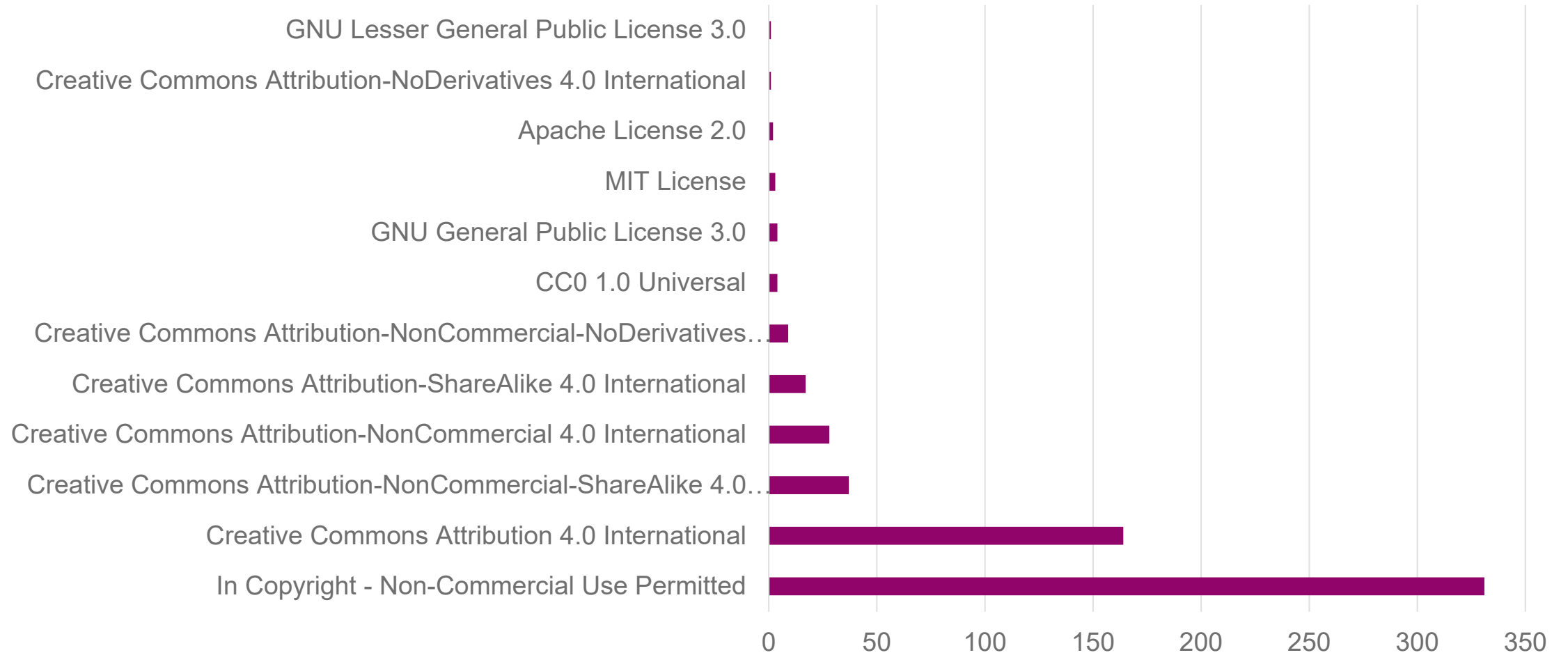
File types (types used in > 10 items)



Availability



Licences



3. Quality assurance and compliance checks

What does ETH Library do in terms of quality assurance?

Metadata

- Check submitter metadata for consistency with repository rules and spelling errors
- Check whether related data and publications are linked correctly
- Add formal metadata (access rights, format, size, publisher,...)

What does ETH Library do in terms of quality assurance?

Files

- Virus check
 - Check readability (open files with common viewer / tool, random sample for large collections)
 - Detect file formats with DROID
 - Check whether file formats are compatible with the chosen retention period
 - Add new formats and support level to file formats registry
 - Check whether file names, folder names and structure are comprehensible
- Quality assurance does not involve repository staff editing / manipulating files
- Repository staff contacts submitters with recommendations, submitters decide whether they apply the recommended changes and resubmit their files

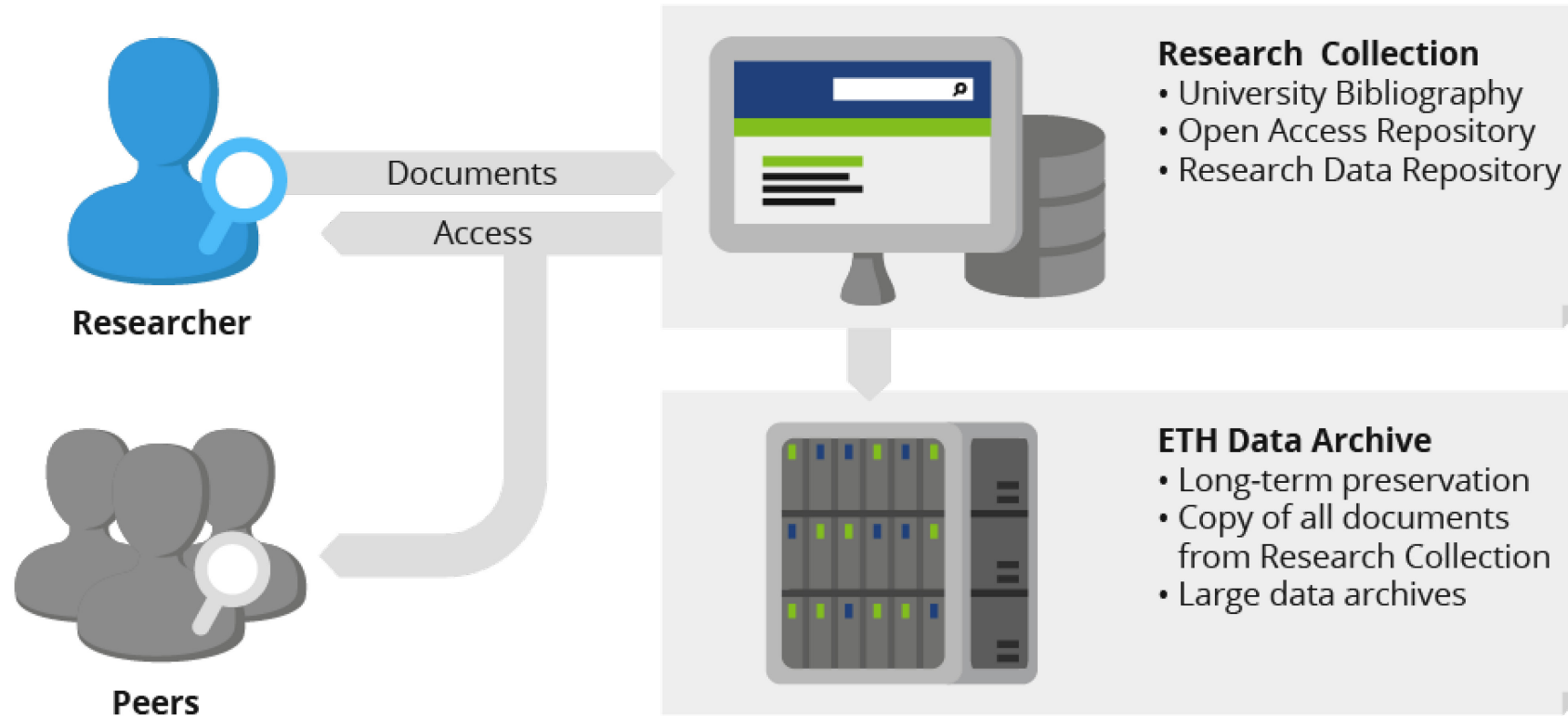
Compliance checks

- In principle, according to Research Collection's [Terms of Use](#), end users are responsible for compliance with the following laws and policies and confirm non-violation of these norms when submitting their data:
 - Swiss Copyright Act
 - Guidelines for the Financial Exploitation of Research Results at ETH Zurich
 - ETH Zurich Compliance Guide
- However, in terms of risk management and as a service to ETH researchers, ETH Library does inform users if repository staff detects violations of these norms during their quality assurance tasks.

Compliance management

Topic	Risk	Mitigation
Copyright	<ul style="list-style-type: none">• Third-party-copyrighted material included in data collection• Conflicting licence statements for overall data collection vs. individual files	<ul style="list-style-type: none">• Inform submitter of potential copyright violations or licence incompatibilities
Exploitation of research results	<ul style="list-style-type: none">• Software not registered and licenced according to ETH guidelines	<ul style="list-style-type: none">• Inform submitter about ETH Software Licencing Policy
Protection of personal data	<ul style="list-style-type: none">• Non-anonymised personal data included in dataset without explicit permission by human subjects	<ul style="list-style-type: none">• Advice users to consult with data protection experts on a case-by-case basis before upload

How to reconcile publishing and preservation requirements



How to reconcile publishing and preservation requirements

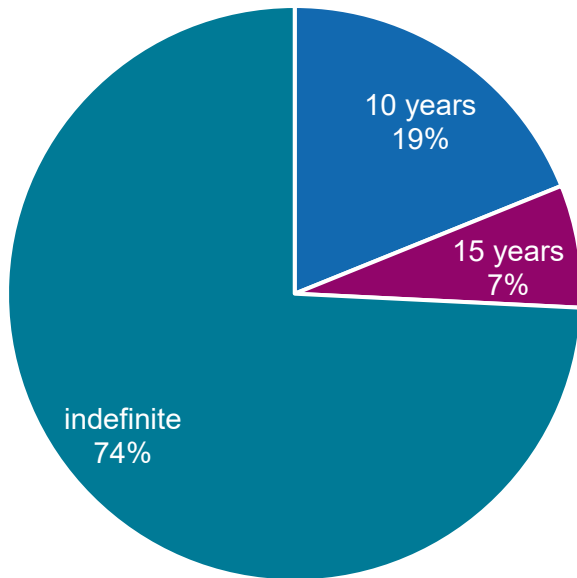
Requirements for preservation

- Submit data in formats suitable for long term preservation
- Either researchers or repository staff must invest resources in file format conversion

User priorities for data publishing

- Publish data for immediate reuse by peers and/or peer reviewers
- Comply with data management regulations of funders and institution
- Invest minimal time in data preparation
- «Better safe than sorry»: users choose indefinite retention period even if file formats are not suitable for long term preservation

How to reconcile publishing and preservation requirements



➤ **New approach** as of 2020:

- ✓ Chosen retention period is no longer used as indicator for how long we must «keep the files readable» but only for how long we must «keep the files» (e.g. bitstream preservation)
- ✓ Submitters need to activate a checkbox if they are actually interested in keeping their data readable over the long term
- ✓ If activated our team provides recommendations on how to convert the submitted files so they become suitable for long term preservation

Defining the scope of your quality assurance / data curation activities

	Level 1	Level 2	Level 3
INGEST	<ul style="list-style-type: none"> Authentication Documentation Chain of custody Metadata Deposit Agreement File Validation 		
APPRAISE/ACCEPT	<ul style="list-style-type: none"> (licenses) Rights management (DUAs) 	<ul style="list-style-type: none"> (file review) Risk management (remediation) 	
CURATE	<ul style="list-style-type: none"> Arrangement & Description Indexing Persistent identifier Transcoding File inventory or manifest 	<ul style="list-style-type: none"> Contextualize Curation log File Renaming Restructure Quality Assurance File Format Transformations 	<ul style="list-style-type: none"> Code Review Conversion (analog) Interoperability Software registry Data cleaning De-identification Peer review
ACCESS	<ul style="list-style-type: none"> Full-text indexing (system-automated) Discovery services Data citation File download Metadata brokerage Contact information Embargo 	<ul style="list-style-type: none"> Restricted Access (mediated requests) 	<ul style="list-style-type: none"> Data visualization
PRESERVE	<ul style="list-style-type: none"> File Audit Migration Cease data curation Secure storage Succession planning Tech monitoring/refresh Versioning 	<ul style="list-style-type: none"> Repository certification 	<ul style="list-style-type: none"> Emulation

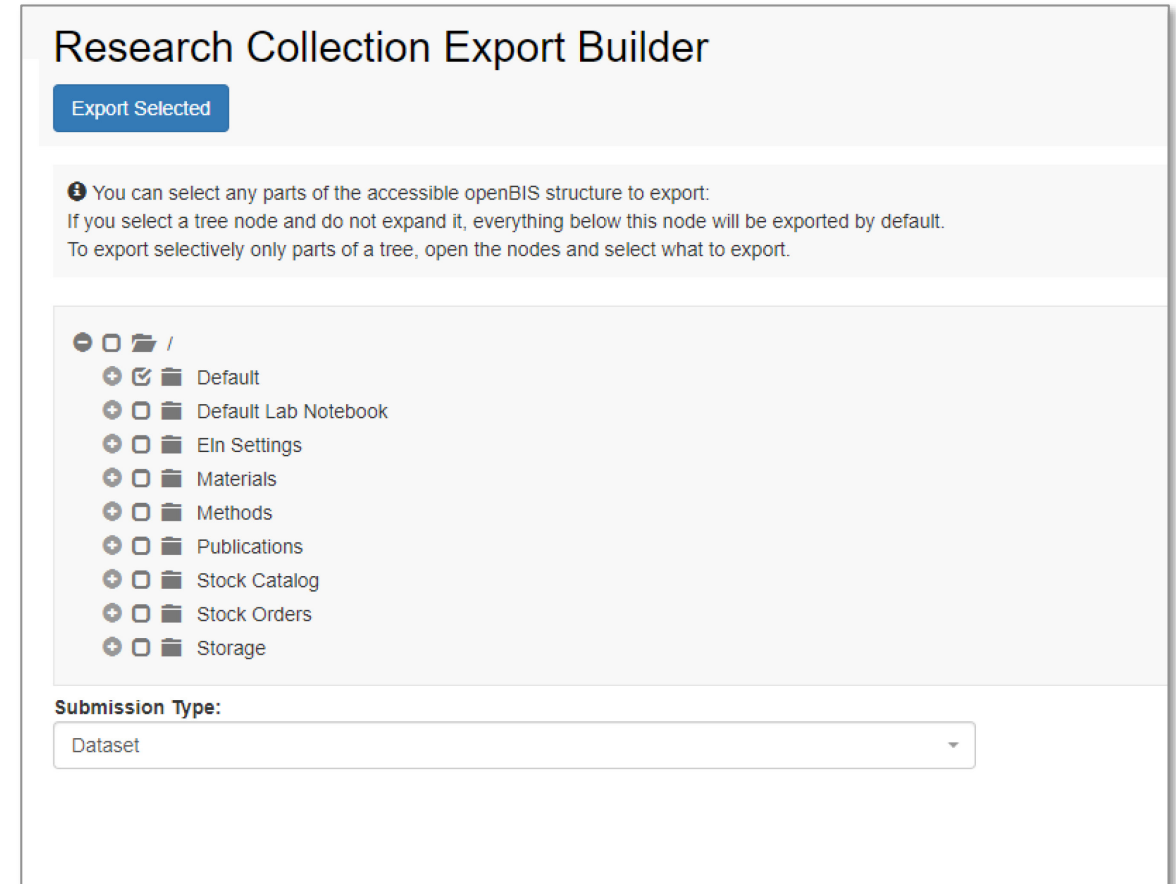
Lafferty-Hess, Sophia, Julie Rudder, Moira Downey et al. (2020). Conceptualizing Data Curation Activities Within Two Academic Libraries. In: *Journal of Librarianship and Scholarly Communication* 8(1). eP2347. <https://doi.org/10.7710/2162-3309.2347>

4. New developments

Integration with openBIS

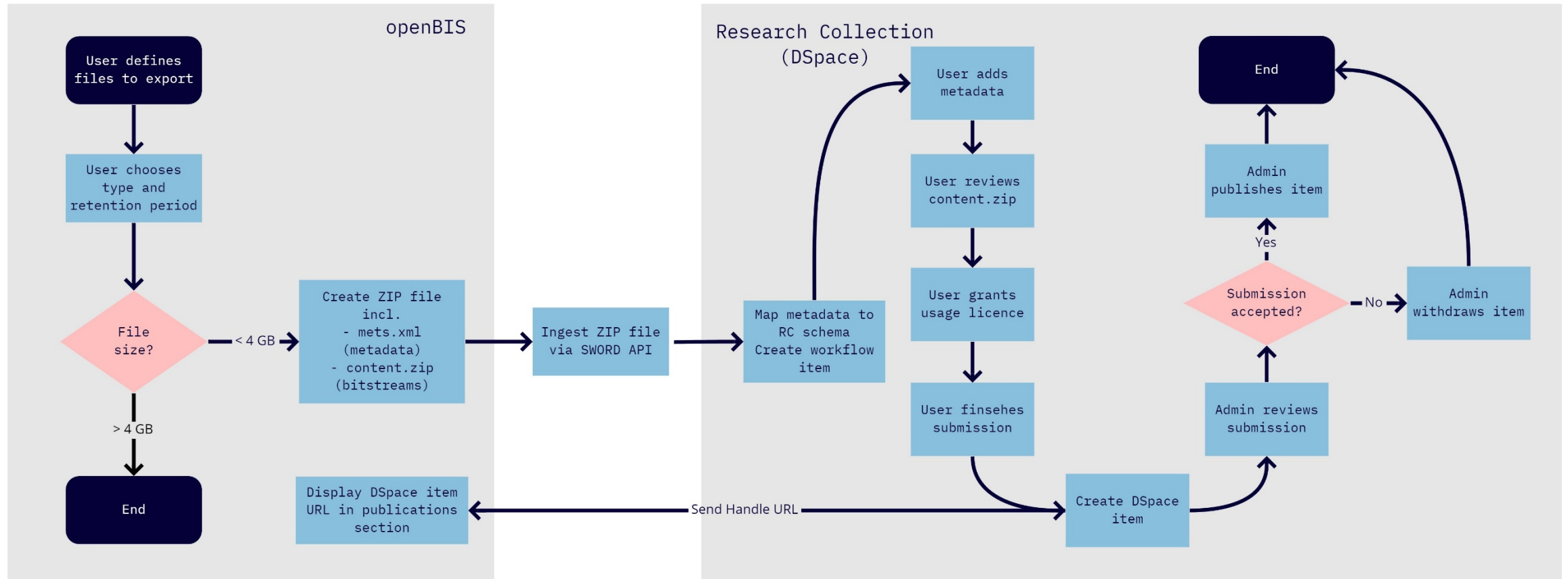
Motivation

- Researchers at ETH can use **openBIS** to manage their daily research and data.
- Researchers use the **Research Collection** to publish their data related to a publication.
- Facilitating the transfer of selected data + metadata from openBIS to the Research Collection benefits researchers who use openBIS for their active data management and provides for an **integrated solution** supporting data management, data publishing, and data preservation.



The screenshot shows the 'Research Collection Export Builder' interface. At the top, there is a blue button labeled 'Export Selected'. Below this, a light gray box contains an information icon and text: 'You can select any parts of the accessible openBIS structure to export. If you select a tree node and do not expand it, everything below this node will be exported by default. To export selectively only parts of a tree, open the nodes and select what to export.' Below the text is a tree view of the openBIS structure. The root node is a folder icon with a minus sign. Underneath, there are several sub-nodes, each with a plus sign and a folder icon: 'Default' (with a checkmark), 'Default Lab Notebook', 'EIn Settings', 'Materials', 'Methods', 'Publications', 'Stock Catalog', 'Stock Orders', and 'Storage'. At the bottom of the interface, there is a 'Submission Type:' label and a dropdown menu currently set to 'Dataset'.

Integration with openBIS



Solution for publishing large datasets

Motivation

- Users want to deposit much larger files than what we can currently accommodate (ca. 10 GB per file)

Solution

- Integrating two existing services at ETH Zurich: Research Collection (DSpace) and polybox (ownCloud)

Solution for publishing large datasets



Research Collection:

- Metadata record / landing page incl. DOI
- Link to download page on «libdrive»



libdrive:

- File upload
- File download

Solution for publishing large datasets

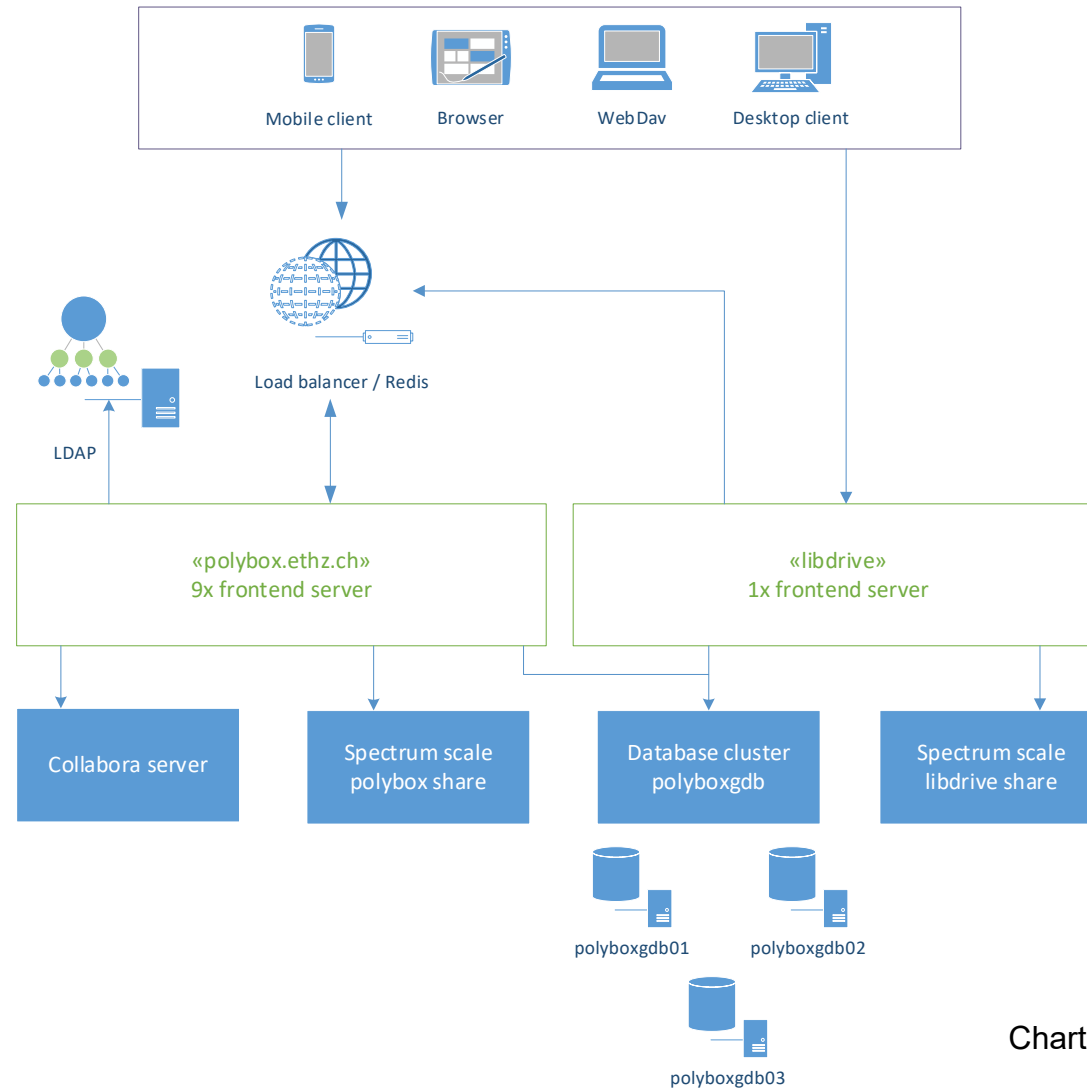


Chart by Gianluca Caratsch, ETH IT Services

Solution for publishing large datasets

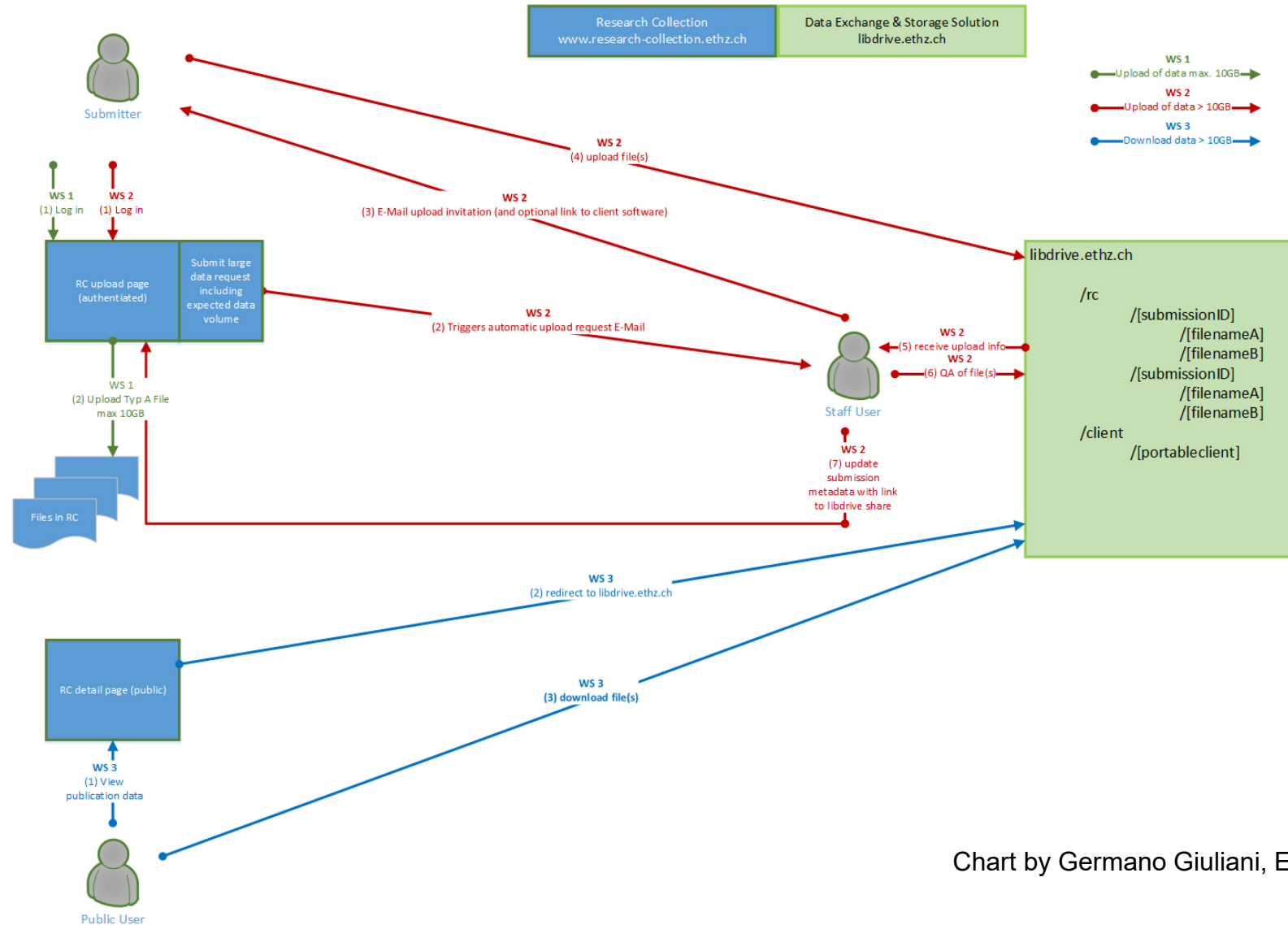


Chart by Germano Giuliani, ETH Library

Solution for publishing large datasets

File size	Upload via	Download via
< 10GB	Research Collection submission form	Browser
10 – 20GB	ownCloud web client	Browser
20 – ca. 200GB	ownCloud client or WebDAV	Browser or client (tbd)
200GB - ca. 1TB	Offline transfer via USB device	Offline transfer via USB device (request access via email form)

What's next? Plans for the future



Certification

Applying for the Core Trust Seal



Geo-referencing

Improved geo-location
referencing and search



Google Dataset search

Implement schema.org for
inclusion of datasets

Barbara Hirschmann
Head E-Publishing
barbara.hirschmann@library.ethz.ch

ETH Library
Rämistrasse 101
8049 Zurich, Switzerland

www.library.ethz.ch