



Monitoring Open Science Activities

Evgeny Bobrov

Berlin Institute of Health at Charité, QUEST Center

QUEST Seminar on Responsible Research

23.05.2023

BIH QUEST
Center for Responsible Research

BIH Berlin Institute
of Health
@Charité

Aus Forschung wird Gesundheit

Charité Dashboard on Responsible Research

Charité has committed itself to establish, promote and maintain a research environment which enhances the robustness of research and the reproducibility of results ([Rethinking Health - Charité 2030](#)).

This dashboard gives an overview of several metrics of open and responsible research at the Charité (including the Berlin Institute of Health). For more detailed information on the methods used to calculate those metrics, the dataset underlying the metrics, or resources to improve your own research practices, click one of the following buttons.

This dashboard is a pilot that is still under development. More metrics will be added in the future.

For more detailed open access metrics you can visit the [Charité Open Access Dashboard](#) developed by the Charité Medical Library.

[See methods](#)
[See resources](#)
[See dataset](#)

Open Science

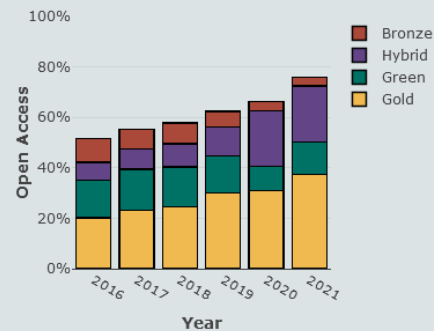
 Show absolute numbers

Double-click or select rectangular area inside any panel to zoom in

Open Access

72 %

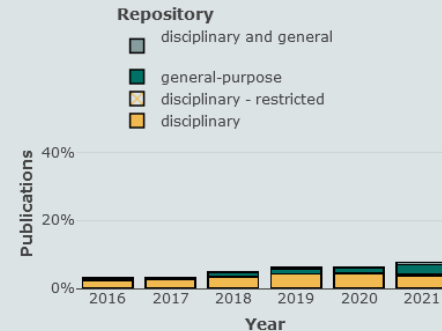
of publications were open access in 2021



Any Open Data

8 %

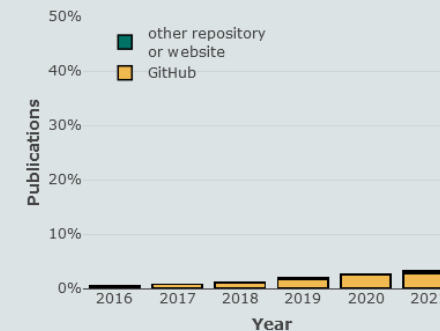
of publications mentioned sharing of data in 2021



Any Open Code

164

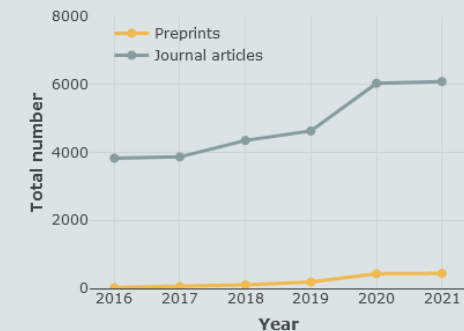
publications mentioned sharing of code in 2021



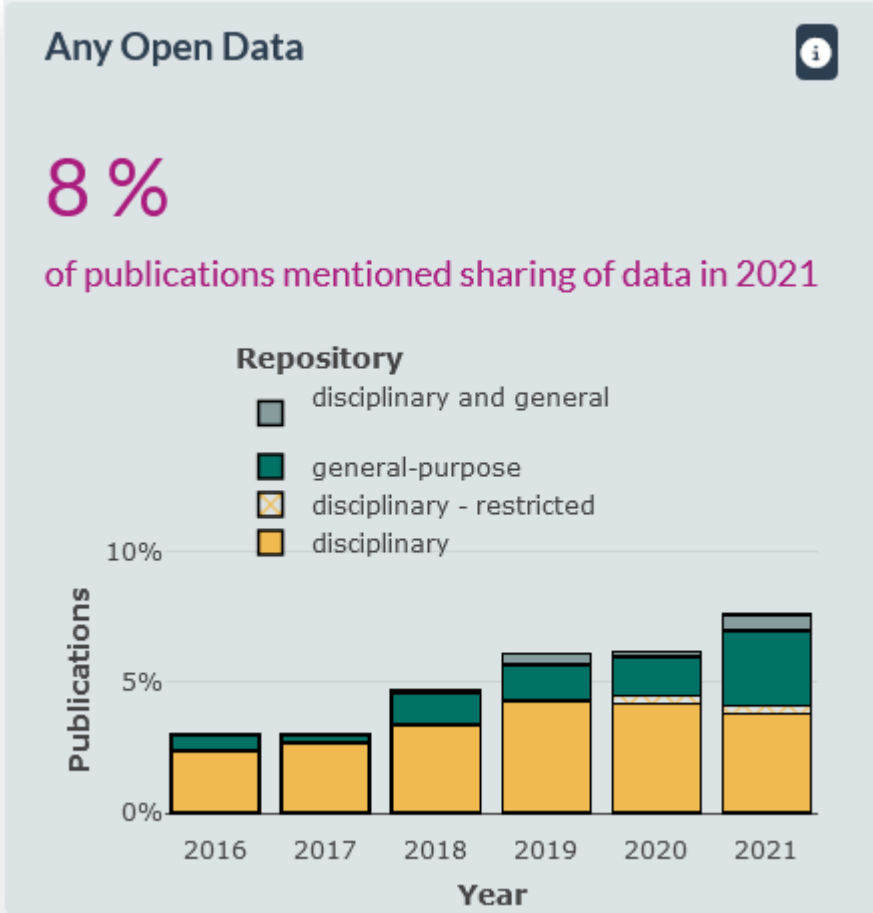
Preprints

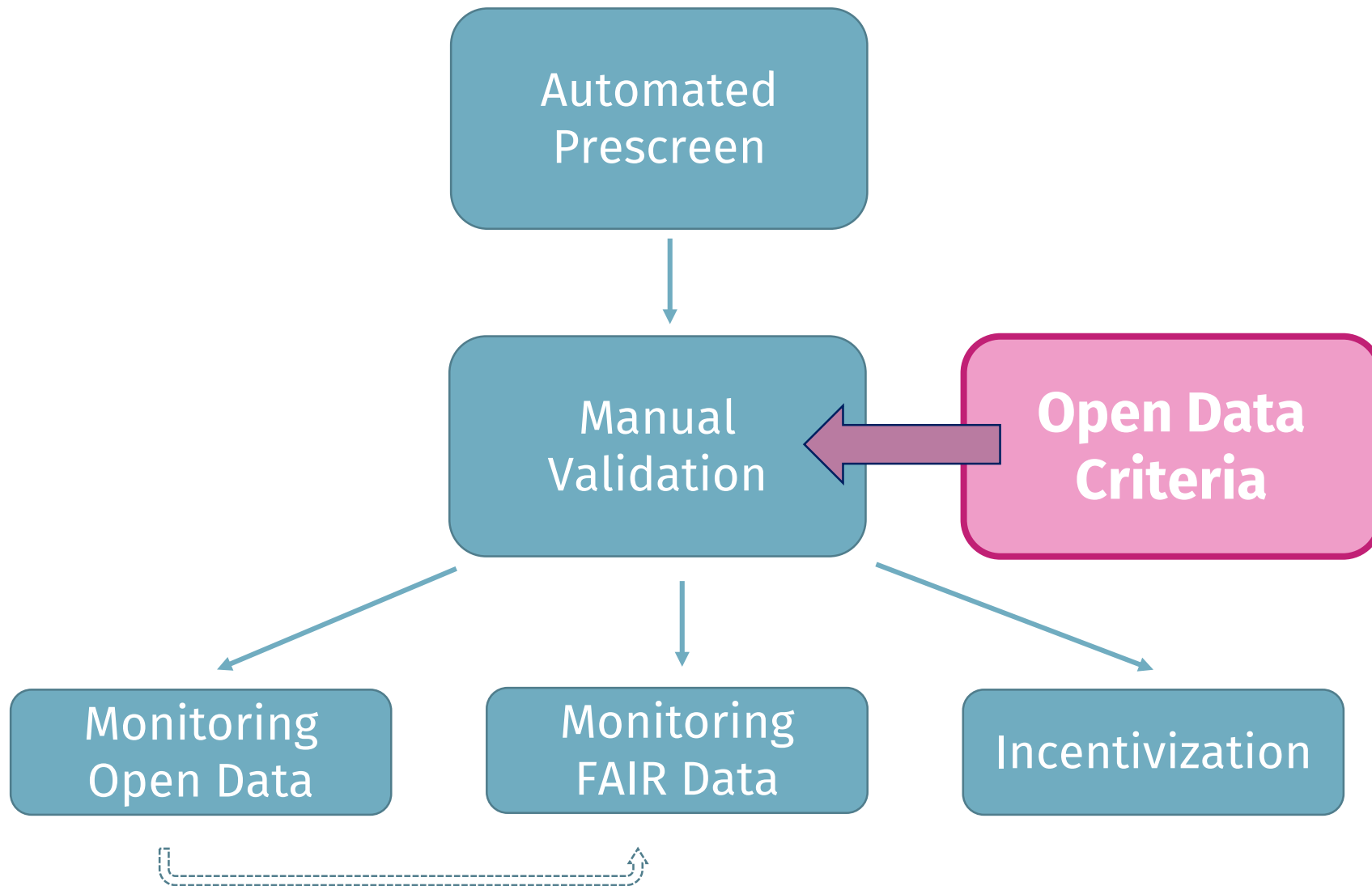
441

preprints published in 2021



Charité Dashboard on Responsible Research

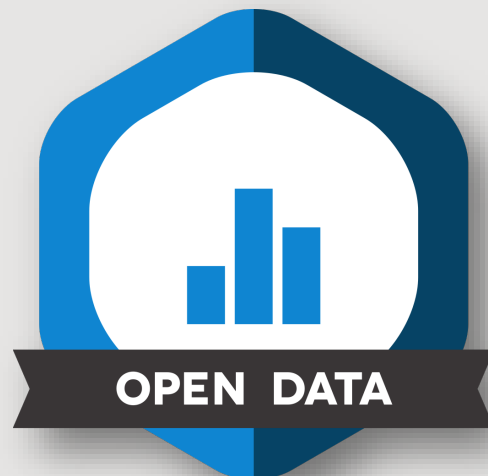






Open Data Definition

Open data is data that can be freely used, re-used and re-distributed by anyone - at most restricted by the obligation to name sources and “share-alike”.



Open Data Definition

Criterion: “Research data have been made **freely accessible** by researchers of the Charité” Main question: What does free access mean in the given context? Definition: Free access is understood as making the dataset available in such a fashion that any human can access it online and use it for any purpose without revealing one’s own identity, as long as this purpose does not explicitly infringe on legitimate rights of authors (e.g., their right to demand citation) or study subjects (e.g. not to be re-identified). Also see demarcations under 1d. Limitation: Criteria for cases where restricted access is considered justified are excluded, and are discussed further below. Demarcation: **Access requirements. Our definition includes** cases, where an agreement has to be accepted, as long as this agreement does not restrict reusers in the type of study or analysis [note that this is not yet included in the published version of the criteria in Kip et al. (2022, June 16)]. Reasoning: for machine-readability, any kind of agreement would be a considerable hurdle, and such obstacles should be avoided. However, query of datasets by humans is so far the norm, and for humans it is not a large obstacle to accept certain terms, although it is not “open” in the full sense. **Our definition excludes** cases where a registration is necessary [this is also not yet included in Kip et al. (2022, June 16)]. Reasoning: the threshold is too high for researchers who would like to look at the data to assess its reuse value for them, especially given that most datasets do not have sufficient metadata to assess their utility purely on these grounds. **Data availability upon request. Excludes** data available upon request. Reasoning: such data are not freely available, but rather only upon the discretion of the authors. The literature shows consistently that data upon request are difficult and sometimes impossible to come by (Vines et al., 2014; Tedersoo et al, 2021; Gabelica et al., 2022). Minimally, it is difficult and one has to reveal information about oneself. **Licenses. Does not require** any specific license. Reasoning: Many repositories do not provide standardized licenses at all. Thus, assessing the openness of licensing and use terms is time-consuming and very difficult to standardize. We do not encourage non-commercial (NC) licenses, similarly to many others (see e.g. Margoni & Tsiavos, 2018), and we are aware that it comes at a cost to reuse to apply them (Hagedorn et al., 2011; The Open Data Institute, 2015; Matthews, 2022). However, they do sometimes occur in biomedical research, and given the availability of data for research purposes, we prefer to incentivize them at this point. In addition, researchers who attached a restrictive license should not be disadvantaged compared to those who did not attach a license at all. No derivatives (ND) licenses would be inappropriate for data reuse, but we have never come across such licensing (Charité Metrics Dashboard - Data Reusability, no date), and as this so extremely rare for datasets, we decided not to address this case. **Data authorship. Requires** clarity about authorship of data by authors of the article. Thus, the definition **excludes** “data from data collections of consortia (“data pools”), if it is unclear whether the authors themselves have contributed to the pool.” Reasoning: It is impossible otherwise to distinguish between data sharing and data reuse. However, if it is explicitly stated in the article that the authors contributed to the data pool, this is considered sufficient. **Does not require** data authorship as listed in repositories specifically by the Charité-affiliated authors of the respective article. Reasoning: Author lists of datasets are often unavailable or list only one person, and it cannot be determined, whether the person listed as data depositor is actually the (only) creator or collector of the data. “The data can be **raw, primary, or secondary data** (e.g. from analyses of freely available datasets, meta-analyses, or health technology assessments); the data would thus allow the analytical replication (retracing of analysis steps) for at least a part of the study’s results; reporting of statistical values (means, standard deviations, p-values etc.) is not sufficient” Please note, that in our experience this is the most difficult of the criteria to check, and one cannot fully avoid heuristics which depend on the specific subfield and are influenced by assessor experience. However, the application of below definition and demarcations substantially constrains the assessor degrees of freedom. Main question: When does data allow analytical replication? Definition: Analytical replication is here understood as the re-tracing of the analysis (quantitative or qualitative) from shared data to the results presented in an article. This re-tracing can be based on raw data as they were collected, or data which were in any way cleaned, normalized or otherwise processed. Demarcation: “Source data” (tables with data points underlying figures). **Includes** so-called “source data”, which list individual measurements underlying figures, and are quite common in biomedical journals. Reasoning: these are considered open data, as they provide additional information to the article and its figures, and could be pooled with other data, even though these data might already be highly derived. **Statistical values. Excludes** statistical values. Reasoning: statistical values do not constitute additional data compared to what is normally reported in articles, and do not allow computational reproduction, as they already constitute the result. Reuse is possible in a meta-analytic way, but is greatly reduced compared to individual observations. **Other outputs than data. Excludes** analysis scripts, computer programs, and other methods, materials, and protocols, even if their development was the goal of the research project and/or their presentation was the focus of the publication. If data has been collected and shared for development or validation, these can, however, be included. Reasoning: the assessment in question focuses on openly available data. We acknowledge the importance, even the essentiality of

Automated Prescreen

ODDPub – Open Data Detection in Publications

- Open Source Tool in R, developed at QUEST
- Detection of keywords in articles → Open Data (and Open Code)
- Open Data detected as combination of terms regarding availability and deposit location, as well as partly identifiers
- Charité: ca. 5500 publications/year → ca. 770 Hits
- → Have to be screened manually

Manual validation

Screening of open data statements in Numbat

- Open source tool in PHP, originally developed for systematic reviews
- Adaptation for open data screening
- → Detailed protocol of the complete workflow in protocols.io



The screenshot shows the Numbat tool interface on protocols.io. At the top, there is a dark blue navigation bar with a search input field containing the text "Search", a magnifying glass icon, and links for "Features" and "Plans". Below the navigation bar, the main content area features a white background. On the left, there is a clipboard icon with a checklist. To the right of the icon, the title "Semi-automated extraction of information on open datasets mentioned in articles" is displayed in bold black text, followed by a downward arrow. Below the title, the authors "Anastasiia Iarkaeva¹, Evgeny Bobrov¹, Jan Taubitz¹, Benjamin Gregory Carlisle¹, Nico Riedel¹" are listed, with a superscripted "1" indicating their affiliation. Below the authors, the affiliation "1Berlin Institute of Health at Charité (BIH), QUEST Center for Responsible Research" is shown. To the left of the main content, there is a vertical sidebar with several buttons: "May 12, 2022", "Bookmark" (with a star icon), "Run" (with a play icon), and "Copy / Fork" (with a fork icon). To the right of the sidebar, there are three buttons: "3 Works for me" (with a number 3), "Share" (with a share icon), and a URL "dx.doi.org/10.17504/protocols.io.q26g74p39gwz/v1". Below the share button, there is a user profile icon and the name "evgeny.bobrov" with a lightning bolt icon.

BUA Open Science Dashboards

- Project funded by Berlin University Alliance (BUA), Objective 3
- **Systematic Monitoring of open science practices** → Application of (semi-)automated tools + presentation of open science indicators in dashboards
- with Open Access Office Berlin, 10/2021-9/2023

FAIR Data Dashboard

Charité Metrics Dashboard

Start page

FAIR data

Berlin Science Survey

Methods/Resources/Data ▾

About

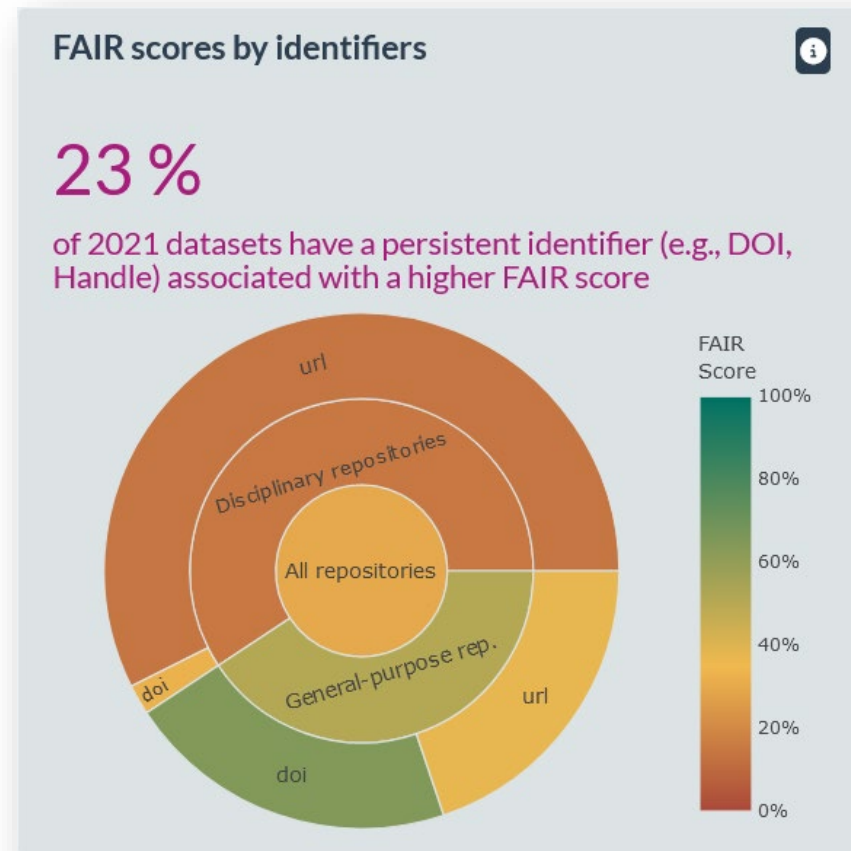
Data Reusability (FAIR data)

Good data management is a key factor in generating, reproducing, and reusing scientific knowledge. The [FAIR principles](#) provide guidance to increase the findability, accessibility, interoperability, and reusability of research data objects. As a result of the increase in volume, complexity, and creation speed of data the FAIR principles emphasize the machine-actionability of data reuse

The FAIR data metrics in this dashboard indicate how well research data objects shared by Charité researchers and the repositories used to deposit them conform with the FAIR principles.

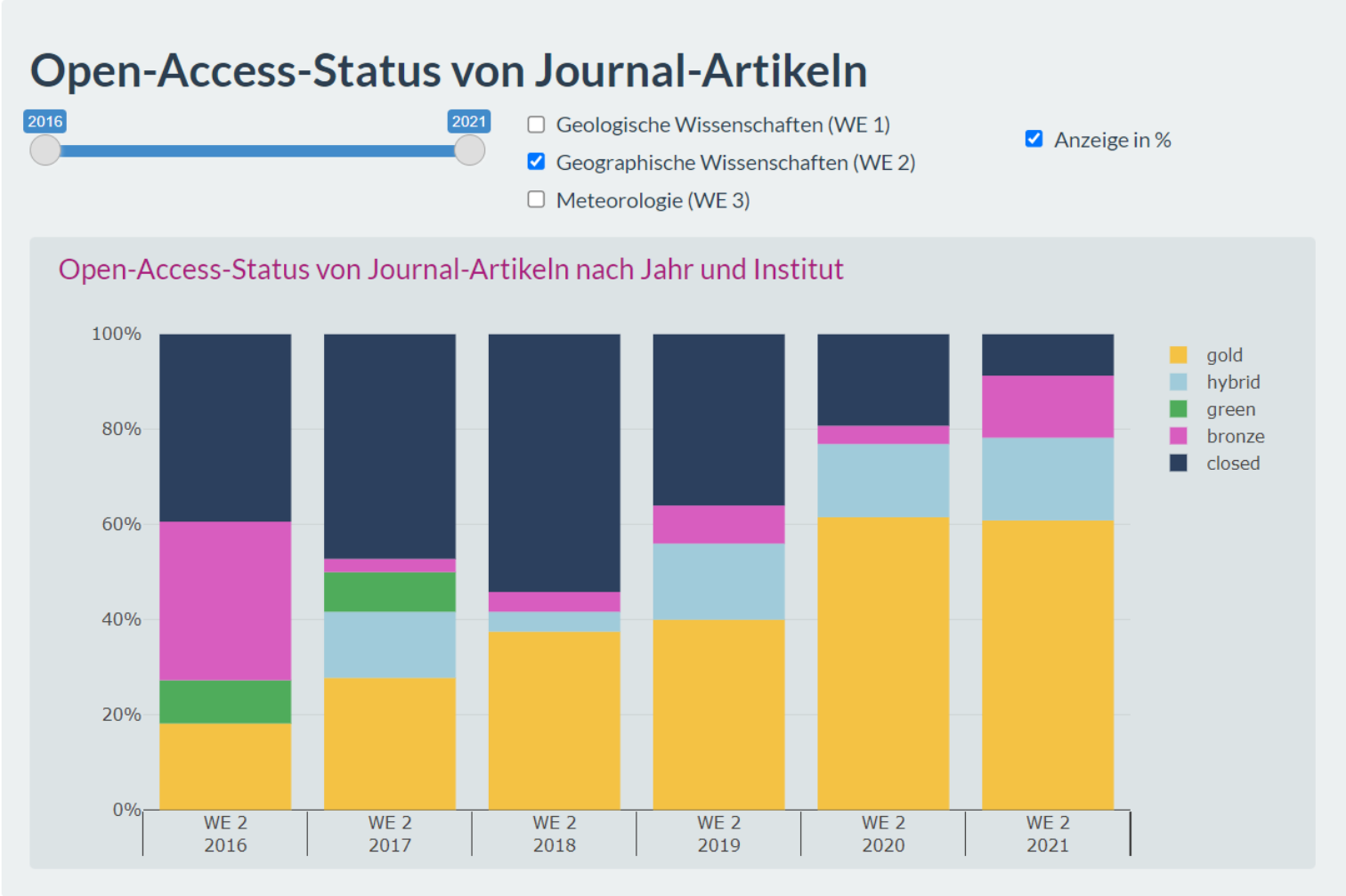
It is important that the FAIR metrics are not to be understood as evaluations, but rather as assistance. This is true at the repository level and even more so at the publication level. Individual researchers have limited influence on the FAIRness of research data objects, which is primarily determined by the data repositories.

FAIR Data Dashboard



<https://quest-dashboard.charite.de/#tabFAIR>

Open Science Dashboard in other fields



From Dashboards to „Magnifiers“



Thanks to collaborators and team members!



Anastasiia
Iarkaeva
*Project Data
Scientist*



Alumnus:
Jan Taubitz
*Project
Data Scientist*



Vladislav
Nachev
*BIH QUEST
Data Scientist*



Alumnus:
Nico Riedel
*BIH QUEST
Data Scientist*

- *Project Team Open Science Dashboards @Open Access Office Berlin*
- *Miriam Kip (Open Data criteria)*
- *Benjamin Carlisle (Numbat)*
- *+ many others for validation of open data extraction*

Thank you!

quest.bihealth.org/

quest-dashboard.charite.de

BIH QUEST
Center for Responsible Research

BIH Berlin Institute
of Health
@Charité

Aus Forschung wird Gesundheit