

# CURATION AT GFBIO DATA CENTERS

Tanja Weibulat<sup>1,2</sup>, Maren Gleisberg<sup>3</sup>, Birgit Klasen<sup>4</sup>, Ivaylo Kostadinov<sup>2</sup>, Jimena Linares<sup>2</sup>, Anke Penzlin<sup>5</sup>, Judith Weber<sup>6</sup>

<sup>1</sup>SNSB IT Center, <sup>2</sup>GFBio e.V., <sup>3</sup>BGBM, <sup>4</sup>ZFMK, <sup>5</sup>SGN, <sup>6</sup>Uni Bremen-MARUM/PANGAEA

## GFBIO DATA CENTERS

The ten GFBio Data Centers (fig. 1) act as infrastructure partners in the GFBio broker network. Data submitted to GFBio are transmitted to and curated by data curators at the GFBio Data Centers (fig. 2). The submitted data are enriched with metadata, dependent on the data type, data format and data content. This is done in close collaboration with the data producer.



Figure 1: List of the GFBio Data Centers (<https://www.gfbio.org/data-centers>)

Descriptors of Curation Levels				
	Level 1	Level 2	Level 3	
Exchange with data producer regarding metadata				
Metadata are curated by data curator				
Metadata are assigned to GFBio consensus elements of biodiversity community agreed standards for data exchange				
Stable Identifiers (e.g. DOI, ENA-accession numbers) are assigned to published datasets				
Exchange with data producer regarding research data content and quality				
Research data are curated by data curator				
Research data are assigned to GFBio consensus elements of biodiversity community agreed standards for data exchange				
Research data are semantically enriched (e.g. by linking to ontologies or identifier services)				
Long-term collaboration between data producer and data curator regarding dynamic datasets				
Remote curation by the data producer				
Continuous versioning of dynamic datasets				

Figure 2: Curation descriptors of GFBio Data Centers arranged according to three curation levels (See also <https://kb.gfbio.org/pages/viewpage.action?pageId=39256097>)

## DATA CURATION DESCRIPTORS

Data curation is done in various ways and at all steps of the data life cycle. The efforts may improve data and metadata quality but are no direct indication for the quality of the original research data. Descriptors concerning metadata, research data and long-term data curation offered by GFBio Data Centers are listed in figure 2.

## GFBIO CURATION LEVELS

The curation levels set up by Task Group 2 – Curation Levels & Criteria in GFBio are intended to provide guidance for data producers. They give them a first impression of the kind and amount of curation efforts needed for handling different types of data and metadata structured according various standard schemas and formats (fig. 3).

**Level 1 curation** addresses metadata curation only. It is applied for example to data packages with metadata structured according community agreed standards for metadata exchange, e.g. EML (Ecological Metadata Language). The research data themselves are structured according proprietary schemes. After curation at the GFBio Data Centers data will be published once with PID assignment.

**Level 2 curation** comprises metadata curation as for level 1. In addition the curation efforts include aspects of research data curation, e.g. conceptual schema and ontology assignment. After curation at the GFBio Data Centers data will be published once with PID assignment.

**Level 3 curation** is meant in addition to level 1 and 2 for data packages from long-term research projects. The data are curated for a longer period and published dynamically with ‘snapshots’ (versions).

Not all GFBio associated Data Centers offer Level 3 data curation. See the Data Center profiles for more details.

## COMMUNITY STANDARDS

As part of the GFBio project, approximately 40 community agreed and domain-specific standards for metadata and data exchange were identified and documented. The profiles of the ten GFBio Data Centers address 20 of them (fig. 3, see table in the GFBio Public Wiki). Community standards as well as domain-specific vocabularies are determining the results of the curation work and of the reusability of the delivered data packages.

Data exchange standards, protocols and formats relevant for the collection data domain within the GFBio network

The National History Collections and Culture Collection BGBM, DSMZ, MfN, SGN, SNSB, ZFMK with their evolving GFBio Collection Data Centers/ Data Archives are partners of several national and international initiatives and projects developing and using data exchange protocols and standards. During the first months of the GFBio project the partners collected and evaluated relevant technical documentation of existing domain-specific data exchange formats, interfaces and protocols with relevance for the integration and harmonization between collection management systems and archive infrastructure as a whole (see table below). This table deals with collection standards, but it was decided to include EML and GML, as well, as they are used and closely connected to the work within GFBio Collection Data Centers/ Data Archives.

This documentation is part of the process to identify existing and to install and integrate new data exchange mechanisms and protocols appropriate for the GFBio network. It provides useful information for software developers and data scientists to set up and run GFBio agreed standard exchange software solutions at the GFBio Collection Data Centers.

Table: Data exchange standards, protocols and formats relevant for the collection data domain within the GFBio network (edit | edit source)

Notes: Cells highlighted in grey indicate standards and protocols for which the GFBio Collection Data Centers/ Archives have expertise or which they use directly or indirectly.

Standard/ Protocol - Acronym	Full name/ Version	Short description	Documentation/ Schema (URL)	Category	Data domain	Status	GFBio Collection Data Centers - Expertise	Related References
ABCD 2.06	Access to Biological Collections Databases v2.06 (2007-06-05)	Standard for the access to and exchange of data about specimens and observations	<a href="#">ISC File v2. Schema v2.06</a> (see also certain sub-versions of ABCD 2.06)	Data Exchange Standard	Collection Archives	accepted (TDWG-6)	BGBM, DSMZ, MfN, SGN, SNSB, ZFMK	[1][2][3]
ABCD 2.1	Access to Biological Collections Databases v2.1 (2014-05)	Enhanced version developed for GGBN and BIRN, but it will not be used by GGBN	<a href="#">ISC File v2. Schema v2.1</a>	Data Exchange Standard	Collection Archives	published	BGBM, MfN, SNSB, SNSB, ZFMK	
ABCD 3.0	Access to Biological Collections Databases v3 (2019-01-01)	New version now described as an IRI, Schema and as Ontology. This allows the access of the standard through semantic queries, encourages element reuse and serves as basis for future software and services in the area of semantic web.	<a href="#">ISC File v3. Schema v3.0</a> (see also v3.0)	Data Exchange Standard	Collection Archives	published	BGBM, MfN	[4]
ABCDONA	DNA Extension for ABCD v2.06 (2009-05-27)	Standardized DNA Schema extension for ABCD to facilitate storage and exchange of data related to DNA collection units. It offers a rudimentary set of DNA-specific data (Sequences).	<a href="#">ISC File v2. Schema v2.06</a>	Data Exchange Standard	Collection Archives	draft (TDWG-6)	BGBM, DSMZ, SGN, SNSB (DSM)	[5]
ABCDDEFG	Access to Biological Collection Databases Extended for Biodiversity	Standard developed for use with paleontological, mineralogical and geological digitalized collection data	<a href="#">ISC File v2. Schema v2.06</a>	Data Exchange Standard	Collection Archives	proposed (TDWG-6)	BGBM, DSMZ, MfN, SGN, SNSB, SNSB	[6][7]
ABCDHISPID	Herbarium Information Standards and Protocols for Interchange of Data (HISPID)	HISPID is a file format serving as extension for ABCD v2.06. It was developed by Australian herbaria to enable the interchange of plant specimen data.	<a href="#">Documentation v1. google code ZIP v1</a>	Data Exchange Standard	Collection Archives	published	(DSM) expertise: HISPID is used in Australia	[8]
AC	Auditorium Core Multimedia Resources Metadata Standard v1.0	A set of vocabularies designed to represent metadata for biodiversity multimedia resources and collections.	<a href="#">Documentation v1. AC on github</a>	Metadata Standard	Collection Archives	accepted (TDWG-6)	BGBM, SNSB, SNSB, ZFMK (all BIRN/Partners)	[9]
BioCASE UMBIT <sup>®</sup>	Universal Access Protocol v1.31 (2012-11-14)	The protocol used in the BioCASE universal network for communication between the central software and the wrapper software sitting on top of the providers databases.	<a href="#">Documentation v1. UMBIT</a>	Data Exchange Protocol	Collection Archives	published	BGBM, DSMZ, MfN, SGN, SNSB, SNSB, ZFMK	[10]
CDM-light	CDM-light v2.04 (2019-12-14)	A light version of the Common Data Model (CDM) developed in the context of the EDC platform for system interoperability.	<a href="#">Documentation v2.04</a>	Data Exchange Standard	Collection Archives	published	BGBM	[11]
DC	Dublin Core: Metadata Element Set <sup>®</sup> v1.1 (2003-09-15)	Dublin Core is a Metadata Standard that was originally developed for libraries but its elements have been reused in many other formats as well, e.g. DDC, Dublin Core	<a href="#">Documentation v1.1 DC v1.1</a>	Metadata Standard	Literature, General	ISO Standard v1.1	BGBM, DSMZ, MfN, SGN, SNSB	[12]

Figure 3: Data exchange standards, protocols and formats relevant for the biodiversity and collection data domain within the GFBio network (<https://gfbio.biowikifarm.net>).

## THE FUTURE

The GFBio Data Centers are partners of the NFDI4Biodiversity consortium and will continue their work on data curation and publication according to their profiles and portfolios. NFDI4Biodiversity is widening the scope and network of data repositories which adds a variety of curation approaches, metadata standards and data schemas. Ongoing documentation of work will be essential.