



The GHGA Metadata Model: A Framework for National Data Sharing

Karoline Mauer^{1,2,7}, Anandhi Iyappan^{3,7}, Deepak Unni^{3,7}, Florian Kraus^{4,7}, Simon Parker^{5,7}, Galina Tremper^{5,7}, Bilge Sürün^{4,7}, Paul Menges^{5,7}, Koray Kirli^{5,7}, Joachim L. Schultze^{1,2,6,7}, Thomas Ulas^{1,2,6,7}, Sven Nahnsen^{4,7}

- ^I Systems Medicine, German Center for Neurodegenerative Diseases (DZNE) e.V., Venusberg-Campus 1/99, 53127 Bonn, Germany,

About GHGA

The German Human Genome-Phenome Archive (GHGA) was established as infrastructure to store and share human omics data. GHGA is funded as part of participants from 21 institutions, GHGA is organized in six data hubs across omics data producers and high-performance computing centers committed to pr

GHGA's Mission and Interaction

GHGA will receive data from major sequencing centers and other institution international initiatives such as genomDE and the European Genome-Phenome and comparability of data for large-scale analyses. Data access will be regulated up by GHGA providing a secure space for data storage and access. GHGA will establish a cloud based analysis platform.



www.ghga.de

National EGA functionality Streamlined data deposition

Central infrastructure for distributed data access committees

Unified ethico-legal framework

Harmonization of (meta)data





newsletter:





² German Center for Neurodegenerative Diseases (DZNE), PRECISE Platform for Genomics at DZNE, and University of Bonn, Bonn, Germany, ³ European Molecular Biology Laboratory (EMBL), Heidelberg, Germany, ⁴ Quantitative Biology Center (QBiC), University of Tübingen, Tübingen, Germany, ⁴ Germany, ⁴ Quantitative Biology Center (QBiC), University of Tübingen, Tübingen, Germany, ⁴ German ⁵ German Center for Cancer Research (DKFZ), Heidelberg, Germany, ⁶ Life and Medical Sciences (LIMES) Institute, University of Bonn, Germany, ⁷ German Human Genome-Phenome Archive (GHGA), Heidelberg, Germany

		The
a consortium in 2020 with the aim t of the National Research Data Infrastruc	to build a national federated ture (NFDI) initiative. With 46	
Germany combining leading institutions rovide scalable infrastructure.	; in genomic medicine, major	
ions in Germany, while exchanging (ne Archive (EGA). Technical harmonization de via Data Access Committees with a solution build on existing infrastructures for high andardized data halysis ata visualization atistics and aggregation	neta)data with national and on of data will ensure quality lid ethico-legal framework set ph-performance computing to Cloud-based analytics platform (PaaS) for large-scale data portals	Technical Metadata
tegration of multiple nics modalities and onnecting omics data to henotype data	Community-specific data portals	
ies and free text possible	 FAIR Data Principles Global Alliance for Genomic Health Research Data Alliance Genomics Standards Consortium Genomic Data Commons 	
by humans and machines ccess Policy , submission of restricted data possible	 Dublin Core Minimum Information Standards 	Metad The go and the
ed and highlighted in the submission o link between objects	 Data Use Ontology Human Phenotype Ontology Experimental Factor Ontology 	represe GHGA The GF
nformation about the file generation indards evel projects etadata schema	- SNOMED CT - ICD-11 - NCTt - HANCESTRO - other controlled vocabularies	commu core th from E Further minimu
¹ Wilkinson et al., The FAIR Guidi and Stewardship. <i>Sci Data</i> 3 , 1 doi: www.doi.org/10.1038/sda	ng Principles for Scientific Data Management 60018 (2016) ta.2016.18	Model The GF seman
@GHGA_DE www.ghga-de.	.github.io/ghga-metadata-schema/	metada







data in GHGA

bal of GHGA is to support all types of human omics data; most prominently, high throughput sequencing data. Depending on the type of study e type of experiment, there can be different metadata properties that are relevant and need to be captured. Based on existing international ards, the Metadata Workstream at the GHGA Consortium has developed a metadata schema to provide a systematic and standardized way of enting metadata by adopting and adhering to the FAIR (Findable, Accessible, Interoperable, and Reusable) data principles¹.

Metadata Schema

HGA Metadata Schema is built to support domain specific requirements for representing information about data generated by various research unities starting with the Cancer and Rare Disease communities. The schema is built incrementally to ensure that (1) the schema has a basic nat is robust and (2) extensions to the schema can be added as new use cases arise. The schema itself is inspired by the metadata schema EGA. This is primarily because a lot of the data in the initial phase of GHGA will be similar to the kind of information that is submitted to EGA. rmore, GHGA's commitment to be a national node in the Federated EGA network requires the schema to be easily translated to EGA with um loss of information.

ling Framework

HGA Metadata Schema is implemented using LinkML, a modeling language and a framework that can be used to build a schema along with ntics. The schema exists as a Yet another Model Language (YAML), a human and machine readable file format that is used to define the lata schema. Using the LinkML framework and the model definition as a YAML file, we generate technology-specific artifacts which are used throughout the GHGA software stack.