

# Harmonizing the Implementation of PIDs across Repositories

Andrea Pörsch<sup>1</sup> // Kirsten Elger<sup>1</sup> // Emanuel Söding<sup>2</sup> // Dorothee Kottmeier<sup>3</sup> // Stanislav Malinovschii<sup>2</sup> // Yousef Razeghi<sup>4</sup> // Sören Lorenz<sup>2</sup> ● 0000-0003-4502-6223 ● 0000-0001-5140-8602 ● 0000-0002-4467-642X **D** 0000-0002-4263-4234 ⓑ 0000-0002-0007-630X
ⓑ 0000-0001-8577-6614 0009-0002-9792-6768

In our increasingly digital and interconnected world, the integration of Persistent Identifiers (PIDs) in metadata are essential for machine-readable and -understandable metadata as also described in the FAIR Guiding Principles for research data management. PIDs provide unique, permanent and machine-readable references to various types of digital objects, including publications, datasets, scientific software, individuals, organizations, samples that together represent the broad range of research outcomes.

# How can we encourage the use of identifiers?

Within the AK Metadata-PIDs working group (a joint initative between the HMC Hub Earth and Environment and the Helmholtz DataHub Earth and Environment), we discussed several PID systems and reached a consensus on recommending specific systems for different purposes: "ORCID" for identifying individuals, "ROR" for organizations, and the "PIDINST" PID for instruments. Our working group has focused on supporting the ongoing PID implementation in research infrastructures by conserving existing, well-established PID implementations (best practices) and promoting their integration in future systems. We further aim to provide support and guidance for new implementations.

😽 📥 ROR 🛛



### What are the differences between metadata schemas?

The two metadata schemas most used in the Earth and Environment sciences are DataCite and ISO 19115 (INSPIRE – EU Comm).

Both schemas are rich and comprehensive for their use cases. While ISO 19115 is dedicated to describe "classical" geodata (spatial data) and is extensively used by geological surveys and agencies, the DataCite schema was developed for DOI registration and is based in the research context. Especially the DataCite Schema has been significantly further developed to support data discovery, the use of PIDs and citations (with clear guidelines of how to include them in the metadata).

In contrast to this, ISO 19115 metadata allows for high granular description of individual data points (e.g., with the information about connected spatial and temporal information – which is not possible in the DataCite schema at the moment), but lacks, e.g. a clear option to add citations of publications or datasets to the metadata. Consequently, the ISO Metadata available via the assessed repositories are much less harmonized for these properties than the DataCite metadata files. Do we need to change this at all?



19115 example, this information is

provides a "Aggregate Information"

property. The ISO schema has the

the relation type from the DataCite

flexibility to link to external sources, here

schema. This is one example for adding

citations in the ISO metadata from GFZ

https://doi.org/10.5880/igets.po.l1.001).

Other repositories use different properties

Data Services (Neumeyer et al., 2017

### DataCite XML Schema

<relatedIdentifier relatedIdentifierType="DOI" relationType="IsReferencedBy">10.1007/s001900050078</relatedIdentifier> <relatedIdentifier relatedIdentifierType="DOI" relationType="IsReferencedBy">10.1007/978-3-642-78149-0\_49</relatedIdentifier> <relatedIdentifier relatedIdentifierType="DOI" relationType="IsDocumentedBy">10.2312/GFZ.b103-16087</relatedIdentifier> < <relatedIdentifier relatedIdentifierType="DOI" relationType="References">10.5880/igets.su.l1.001</relatedIdentifier> <relatedIdentifier relatedIdentifierType="DOI" relationType="IsCitedBy">10.1007/s40328-018-0212-5</relatedIdentifier> </relatedIdentifiers>

### ISO 19115 XML Schema

- <gmd:aggregationinfo></gmd:aggregationinfo>		
- <gmd:md aggregateinformation=""></gmd:md>		
- <gmd:aggregatedatasetidentifier></gmd:aggregatedatasetidentifier>		
- <gmd:rs identifier=""></gmd:rs>		
- <gmd:code></gmd:code>		
<pre><gco:characterstring>10.2312/GFZ.b103-16087</gco:characterstring></pre>	1g> ←	Identifie
	•	
- <gmd:codespace></gmd:codespace>	I de sette	C
<pre><gco:characterstring>DOI</gco:characterstring></pre>	Identi	ner type
	(L	DOI)
- <gmd:associationtype></gmd:associationtype>		
<pre><gmd:ds <="" associationtypecode="" codelist="http://datacite.org/schema/kernel/&lt;/pre&gt;&lt;/td&gt;&lt;td&gt;1-4" td=""><td></td></gmd:ds></pre>		
codeListValue="IsDocumentedBy">IsDocumentedBy <td>TypeCode&gt;</td> <td>6</td>	TypeCode>	6
	Relation	type with
	referen	ce to the
	DataCit	e schema

### Summary and Outlook

The metadata schemas ISO19115 and DataCite have been developed by different communities and for different purposes. Both have their strengths and weaknesses. It is possible to map a subset of metadata properties in a harmonized way (e.g. keywords, authors information), some metadata properties remain challenging (e.g. citations or PIDs in ISO metadata, dedicated spatial/temporal context in DataCite metadata). The full record, however, cannot be mapped. As a result, the content of the ISO metadata representations vary much more between the repositories than the DataCite metadata.

on type with

Do we need to align these schemas to be able to map the full metadata **between them?** For this, we need to further explore the use cases and especially the connections between the communities. If more geological surveys used DOIs, it would be reasonable to expand the ISO metadata schema for the inclusion of PIDs like ROR, ORCID and related references (citations). The DataCite Metadata Working group is already exploring option to be more specific and actively looks into properties already used by ISO 19115 (pers comm, K. Elger).

# HELMHOLTZ