ESS-DIVE January 2020 Webinar: Community use of persistent sample identifiers and standard metadata to support sample management, tracking, and search

> Joan Damerow, Deb Agarwal, Kristin Boye, Eoin Brodie, Madison Burrus, Shreyas Cholia, Hesham Elbashandy, Ricardo Eloy Alves, Kim Ely, Amy E Goldman, Val Hendrix, Zarine Kakalia, Ken M Kemner, Annie B Kersting, Katharine Maher, Nancy Shiao-Lynn Merino, Fianna O'Brien, Zach Perzan, Emily Robles, Cory Snavely, Patrick Sorensen, James Stegen, Pamela Weisenhorn, Karen Whitenack, Mavrik Zavarin and Charuleka Varadharajan

#### ESS-DIVE Deep Insight for Earth Science Data



Office of Science BERKELE





### **ESS-DIVE: Toward a Fusion Database**



Future: Fusion Database Advanced data search and integration across datasets

se nd ets

Interoperable, standardized metadata needed for Fusion DB`

## Address Community Need - Samples

Challenge: Sample identification and tracking

Need practical, standardized sample identification, description, and tracking system to **integrate and publish data** from field, lab, DOE user facilities

Support Data Linking/Interoperability and Reuse: Advanced data search and integration

**Solution:** International Geo Sample Numbers (IGSNs) for ESS samples

- SESAR provides standardized core sample metadata, templates
- Linking to other samples, online metadata profiles, datasets, publications



# Recent and Upcoming Sample-Related Activities

May 2019 - Start IGSN pilot test, gathered first round feedback through Nov 2019

Oct 2019 - National Microbiome Data Collaborative (NMDC) Ontology workshop

Dec 2019 - AGU: IGSN Community Meeting

Jan 2020 - Earth Science Information Partners (ESIP): Defining the Bull's Eye of Sample Metadata

Jan 2020 - USGS: National Digital Catalog Sample Collections Metadata workshop

March 2020 - Research Data Alliance (RDA) <u>15th Plenary Meeting</u>: Identifying physical samples and linking them to the global research data ecosystem



#### Internet of Samples: Toward an Interdisciplinary Global Infrastructure for Material Samples

**iSamples** 53 Neil Davies Local Global Collection John Deck Transcription Authorities Index Index Sarah Kansa IGSN Kerstin Lehnert John Kunze iSamples ARK iSamples Central In-A-Box Sarah Ramdeen Spreadsheet Dave Vieglais DOI Ramona Walls ... Raymond Yee ന Metadata Selected iSamples Profiles Proposal: Oct 2019 Profiles App NMDC and ESS-DIVE submitted letters of Individual Global Organization collaboration

## Sample PID and Metadata Pilot Test

**Diverse, Interdisciplinary Projects**  $\rightarrow$ Biogeochemical responses to contamination, warming, disturbance

SLAC SBR SFA - Kristin Boye, Zach Perzan LLNL SBR SFA - Mavrik Zavarin, Nancy Merino PNNL SBR SFA - James Stegen, Amy Goldman ANL SBR SFA - Pamela Weisenhorn LBNL Watershed Function SFA and TES SFA FICUS -Eoin Brodie, Patrick Sorenson, Ricardo Eloy Alves BNL TES SFA NGEE Tropics - Kim Ely









# Basic Sample Metadata for Pilot Test

#### **Required Fields**

- PID, and Parent PID where relevant
- Sample Name (unique)
- Chief Scientist, Collector
- Sample type: Object type (controlled), Material (controlled)
- Collection date (ISO)
- Latitude, Longitude (WGS 84), if relevant
- Project (controlled)
- Collection method
- Bio: Scientific Name
- Processing details (required for certain analyses

#### **Recommended Fields**



- NEW: Site ID, Collection Event ID
- Coordinate uncertainty or Instrument used
- Location description
- Sample description
- Sample type: Classification (controlled)
- Primary Physiographic feature (ENVO)
- Preservation method
- Purpose
- Collection time
- Depth (min, max), Depth scale (required if depth entered)
- Elevation
- Physiographic Feature (local environment)
- NEW: Biome (broader ecosystem)

Crosswalk: <a href="https://docs.google.com/spreadsheets/d/1pV1zmqRFbAhA9kxcbOnz2RDWO8siXjnlfL2VhYJHDiM/edit#gid=0">https://docs.google.com/spreadsheets/d/1pV1zmqRFbAhA9kxcbOnz2RDWO8siXjnlfL2VhYJHDiM/edit#gid=0</a>



# Sample Metadata Research - Existing standards and templates

General: Dublin Core, DataCite, ISO Observations and Materials

Geoscience: IGSN Descriptive Metadata Schema, SESAR IGSN template, USGS, EPA

**Biology:** Darwin Core, Minimum Information about any sequence (MIxS) and Environmental Packages

Other Relevant Ontologies or Shared Vocabularies: e.g. ENVO and GAZ

**Crosswalk:** <u>https://docs.google.com/spreadsheets/d/1pV1zmqRFbAhA9kxcbOnz2RDWO8siXjnlfL2VhYJHDiM/edit#gid=0</u>

#### 2,476 Registered IGSNs







1,525: Gas,

Porewater

# User Feedback: Sample Registration

- More automation for large campaigns
  - Field metadata collection apps
  - Scripts to standardize more quickly
  - Convert to xml or json and use API
  - Inherit relevant metadata from parent
- Describe collection Landing pages for collections

	IGSN: Sample Name: Other Name(s): Sample Type: Parent IGSN:	IEBWE000P BWE201406031C0010 Core Section IEBWE0001
Relation To Parent Depth in Core (min): Depth in Core (max):	0 Centin 10 Centi	SESAR Sample Landing Page
Related Samples		
Parents:	IEBWEOO	01 BWE201406031C
Siblings:	<ul> <li>IEBWI</li> </ul>	E000Q BWE201406031C1020 E000R BWE201406031C2030 E000S BWE201406031C3040 E000T BWE201406031C4050 E000U BWE201406031C5060 E000V BWE201406031C6070 E000W BWE201406031C7080 E000X BWE201406031C8090
Children:	No Child	ren



# GFZ

#### Helmholtz Centre POTSDAM

#### HELMHOLTZ CENTRE POTSDAM **GFZ GERMAN RESEARCH CENTRE** FOR GEOSCIENCES

General Identifiers		
Program:	ICDP	<b>逆(適該該該</b> 首
Expedition:	ICDP 5054	法
Type:	Core	
Name:	5054_1_A_658_Z	
IGSN:	ICDP5054ECYD101	
Parent IGSN:	ICDP5054EHW1001	
Release Date:	2017-3-1	24
Sampling Location		
Latitude:	63.40163	
Longitude:	13.202917	
Coordinate System:	WGS84	
Elevation:	-1735.452	
Final Depth:	-1741.562	
Location Type:	N/A	
Location Name:	Åre, Jämtlands län, Sweden	
Location Description:	COSC-1 is located in the vicinity of the mine	ne abandoned Fröå
Country:	Sweden	
Province:	Jämtlands län	
County:	N/A	
City:	Åre	
Geology		
Material:	Rock	
Rock Classification:	N/A	
From Corrected Depth:	2257.962	
To Corrected Depth:	2264.072	
Depth Reference:	meter below ground level	
Geological Age:	mid-paleozoic	
Geological Unit:	N/A	
Methods		
MSCL	yes	
XRF	yes	
Lithological Description	yes	
Core Overview	yes	
Core Section Scan	yes	
Core Catcher Scan	no	
Drilling		
Drilling Method:	Coring>RockCorer	

wireline diamond coring, HQ and NQ bit size Lund University, Engineering Geology

Swedish Research Council (Vetenskapsrådet)

Larsson Drilling Consulting AB

2400.1m

Sample Family
▼ ⊕ 5054_1_A
5054_1_A_1_Z
▶ 🗂 5054_1_A_2_Z
5054_1_A_3_Z
▶ 🗂 5054_1_A_4_Z
▶ 🗂 5054_1_A_5_Z
▶ 🛅 5054_1_A_6_Z



 $\Phi$ =Hole,  $\Box$ =Core,  $\blacksquare$ =Core-Section,  $\blacksquare$ =Core-Sample

The Sample Family shows a sub-sampling graph. Select entries to navigate samples. Core-Samples are issued to scientists on request. The naming convention for a Core-Sample is:

Expedition\_Site\_Hole\_Core\_Section,from-to(cm). Hole, Core, and Core-Section are following the same schema respectively.

#### Location Map



Drilling Start/End: 2013-09-05 / 2014-08-26 \* Latitude: 63.40163 \* Longitude: 13.20292 \* Åre, Jämtlands län, Sweden

ist, B. S. G., Berthet, T., ... Tsang, C.-F. (2015). COSC-1 - drilling of a subduction-related allochthon in the Palaeozoic Caledonide orogen of Scandinavia. Scientific

Conze, Ronald; Gee, David G.; Klonowska, Iwona; Pas-

#### Publications & Datasets

Lorenz, H., Rosberg, J.-E., Juhlin, C., Bjelm, L., Almqv-Drilling, 19, 1-11. doi:10.5194/sd-19-1-2015

Lorenz, Henning; Rosberg, Jan-Erik; Juhlin, Christopher; Bjelm, Leif; Almqvist, Bjarne; Berthet, Théo; cal, Christophe; Pedersen, Karsten; Roberts, Nick;

Operator:

Funding Agency:

Total Length:



# Proposed Option: Collections, Events, Sites



#### Metadata in Existing, Published Data Packages

- Searches for soil, water, leaf samples
- Pangaea, ORNL DAAC, ESS-DIVE, DataONE
- 19 data packages
- Consider separate requirements for collection/samples

Collection and sample metadata from published data packages



#### **User Feedback: IGSNs for Subsamples**





# Proposed Option: ID Extensions for Subsamples

				Add extensions for subsamples and record parent-child relationships in metadata
	Soil Core	Collection of water samples at one	→ Collection of highly related samples	IEBWE0001
		place and time		
	Core Section	Varying depths	$\rightarrow$ Primary sample	IEBWE00011M IEBWE00012M
¢	Subsamples sent for analyses	Each individual container sent for analyses	→ Subsamples/ Containers	IEBWE00012MA IEBWE00012MB IEBWE00012MC

### **ARK Qualifiers for Subsamples**





# User Feedback: IGSN Collection Metadata



Some metadata not previously recorded supports integration and reuse

- Sample relationships parent IGSN,
- Sample Type

Additional classification needed: plant (biology), water

SESAR IGSN Liquid>aqueous informal classifications



# User Feedback: IGSN Metadata



Some metadata not previously recorded supports integration and reuse

- Sample relationships parent IGSN,
- Sample Type

Additional classification needed: plant (biology), water

#### Key Need: Linking Geo and BioSamples

# General Metadata provided for liquid>aqueous samples





• Need links between IDs

## IGSN and BioSample Interoperability



- Need links between IDs
- IGSN metadata → BioSample MIxS/MIMS, Env Packages
- Use of standard ontologies
- Classifications for water samples
- Add biome, sample processing

eroperability	
IGSN:	BioSample: SAMN04007148
IGSN	Source_material_id = IGSN Link
Material	env_medium* = ENVO; isolation_source
populate from material	organism* ( <u>e.g. soil metagenome</u> )
Physiographic feature	Env_local_scale* = ENVO
	Env_broad_scale* = ENVO
Country; Locality	Geo_loc_name* = GAZ
	samp mat process

# ESS-DIVE Plans for Sample PID registration and management



**Near term** - Continue to use SESAR to register sample PIDs, with options:

- Register samples and update metadata through SESAR portal or API
- Customize SESAR IGSN template with some additional fields and picklists
- Offline generation of IDs that you can register later, ID extensions
- Document recommended practices and metadata template

#### **Address Efficiency**

- Sample collection app to obtain standardized info in the field
- Separate collection, site, event metadata, scripts transfer to samples

Continue discussions with OSTI about administering DOE sample PIDs

## ESS-DIVE Plans for Fusion Database and Sample Tracking



Prototyping for Fusion database

Need pilot testers to publish sample related datasets in ESS-DIVE

Advanced search and integration - sample metadata parsing/indexing

Sample tracking - connect relevant collection metadata and data using PIDs

- Work with DOE user facilities

## Plans to Address Geo and BioSample Interoperability



Scripts to generate MIxS template for registered IGSNs

- Link IGSN and BioSample IDs: IGSN URL in BioSample metadata, vice versa

Paper on pilot test and geo/bio sample interoperability, recommended practices

Collaborate with broader community: NMDC, GSC, IGSN, DOE user facilities, USGS, sample groups at ESIP and RDA

### Sample Metadata Crosswalk



https://docs.google.com/spreadsheets/d/1pV1zmqRFbAhA9kxcbOnz2RDWO8siXjnlfL2VhYJHDiM/edit#gid=0

#### Meetings/form to get detailed feedback on sample metadata requirements

Elements that should be required

Elements that are more appropriate or feasible at a collection-level