

## **Title: Standardizing sample information to support efficient tracking, data integration, and reuse in Environmental Systems Science**

Joan Damerow<sup>1</sup>, Deb Agarwal<sup>1</sup>, Kristin Boye<sup>2</sup>, Eoin Brodie<sup>1</sup>, Shreyas Cholia<sup>1</sup>, Hesham Elbashandy<sup>1</sup>, Kim Ely<sup>3</sup>, Amy Goldman<sup>4</sup>, Valerie Hendrix<sup>1</sup>, Christopher Jones<sup>5,6</sup>, Matthew Jones<sup>5,6</sup>, Zarine Kakalia<sup>1</sup>, Kenneth Kemner<sup>7</sup>, Annie Kersting<sup>8</sup>, Kate Maher<sup>9</sup>, Nancy Merino<sup>8</sup>, Fianna O'Brien<sup>1</sup>, Zach Perzan<sup>9</sup>, Emily Robles<sup>1</sup>, Cory Snavely<sup>10</sup>, Patrick Sorensen<sup>1</sup>, James Stegen<sup>4</sup>, Pamela Weisenhorn<sup>7</sup>, Karen Whitenack<sup>1</sup>, Mavrik Zavarin<sup>8</sup>, and Charuleka Varadharajan<sup>1</sup>

<sup>1</sup>Lawrence Berkeley National Laboratory, Berkeley, CA; <sup>2</sup>SLAC National Accelerator Laboratory, Menlo Park, CA;

<sup>3</sup>Brookhaven National Laboratory, Upton, NY;

<sup>4</sup>Pacific Northwest National Laboratory, Richland, WA;

<sup>5</sup>National Center for Ecological Analysis and Synthesis (NCEAS), Santa Barbara, CA

<sup>6</sup>DataONE, Santa Barbara, CA

<sup>7</sup>Argonne National Laboratory, Lemont, IL

<sup>8</sup>Lawrence Livermore National Laboratory, Livermore, CA

<sup>9</sup>Department of Earth Systems Science, Stanford University, Palo Alto, CA

<sup>10</sup>National Energy Research Scientific Computing Center (NERSC), Berkeley, CA;

**Contact:** ([JoanDamerow@lbl.gov](mailto:JoanDamerow@lbl.gov))

**Project Lead Principle Investigator (PI):** Deb Agarwal, Charuleka Varadharajan

**BER Program:** CESD Data Management

**Project:** Environmental Systems Science Data Infrastructure for a Virtual Ecosystem (ESS-DIVE)

**Project Website:** <https://ess-dive.lbl.gov/>

**Project Abstract:** Physical samples are foundational entities for research in earth and environmental sciences; they are not only the basis of individual studies but could also be integrated with other data to inform new and broader-scale questions. Data contributors to the Department of Energy's (DOE) Environmental Systems Science Data Infrastructure for a Virtual Ecosystem (ESS-DIVE) repository often work in large, interdisciplinary teams and send samples to multiple facilities for analyses. This community needs an efficient system for persistent sample identification to enable efficient sample tracking, data integration, and reuse.

We therefore conducted a community pilot test on the use of persistent identifiers for physical samples—specifically, International Geo Sample Numbers (IGSNs). Eight projects with a variety of sample types registered samples for IGSNs, standardized sample collection metadata, published sample metadata in the System for Earth Sample Registration (SESAR) sample catalog and have or will publish related datasets in ESS-DIVE. We compare existing sample-related standards and evaluate the experience of users to develop practical recommendations for sample identification and documentation. We gathered information for the pilot test through discussions with project teams and documented several components, such as the efficiency of the process (i.e. use of templates, labeling, registering samples, and updating metadata) and any apparent problems. We resolved uncertainties in the allocation of related sample identifiers, use of metadata elements, and added standard terms as needed. Throughout the pilot test, we also gathered feedback on desired use cases, which include: improvements in data management, advanced search capabilities, ability to link identifiers, promote interoperability of biological and geological samples, and ability to integrate and reuse sample data.

The pilot test has informed community-driven standards and tools for sample identifiers, tracking, and metadata in the ESS-DIVE repository. Our overall goal is to provide practical recommendations for efficient sample data management while also preserving and maximizing the potential value of samples into the future.