

# Internet of Samples: Creating and Mapping Controlled Vocabularies for Specimen Type, Material Type, and Sampled Feature

Dave Viegla<sup>1</sup>, Quan Gan<sup>2</sup>, Yuxuan Zhou<sup>2</sup>, Stephen Richard<sup>3</sup>, Hong Cui<sup>2</sup>, Neil Davies<sup>4</sup>, John Deck<sup>4</sup>, Eric Kansa<sup>5</sup>, Sarah Kansa<sup>5</sup>, John Kunze<sup>6</sup>, Danny Mandel<sup>2</sup>, Chris Meyer<sup>7</sup>, Thomas Orrell<sup>8</sup>, Sarah Ramdeen<sup>9</sup>, Rebecca Snyder<sup>8</sup>, Ramona Walls<sup>2</sup>, and Kerstin Lehnert<sup>9</sup>

<sup>1</sup>University of Kansas

<sup>2</sup>University of Arizona

<sup>3</sup>U.S. Geoscience Information Network

<sup>4</sup>University of California Berkeley

<sup>5</sup>Open Context, The Alexandria Archive Institute

<sup>6</sup>California Digital Library

<sup>7</sup>Smithsonian National Museum of Natural History

<sup>8</sup>Smithsonian Institution

<sup>9</sup>Columbia University

November 24, 2022

## Abstract

Material samples are vital across multiple scientific disciplines with samples collected for one project often proving valuable for additional studies. The Internet of Samples (iSamples) project aims to integrate large, diverse, cross-discipline sample repositories and enable access and discovery of material samples as FAIR data (Findable, Accessible, Interoperable, and Reusable). Here we report our recent progress in controlled vocabulary development and mapping. In addition to a core metadata schema to integrate SESAR, GEOME, Open Context, and Smithsonian natural history collections, three small but important controlled vocabularies (CVs) describing specimen type, material type, and sampled feature were created. The new CVs provide consistent semantics for high-level integration of existing vocabularies used in the source collections. Two methods were used to map source record properties to terms in the new CVs: Keyword-based heuristic rules were manually created where existing terminologies were similar to the new CVs, such as in records from SESAR, GEOME, and Open Context and some aspects of Smithsonian Darwin Core records. For example specimen type = *liquid* > *aqueous* in SESAR records mapped to specimen type = *liquid or gas sample* and material type = *liquid water*. A machine learning approach was applied to Smithsonian Darwin Core records to infer sampled feature terms from record text describing habitat, locality, higher geography, and higher classification fields. Applying fastText with a 600-billion-token corpus in the general domain, we provided the machine a level of “understanding” of English words. With 200 and 995-record training sets, 87%, 94% precision and 85%, 92% recall were obtained respectively, yielding performance sufficient for production use. Applying these approaches, more than 3x10<sup>6</sup> records of the four large collections have been mapped successfully to a common core data model facilitating cross-domain discovery and retrieval of the sample records.

## Introduction

Material samples play vital roles in multiple scientific disciplines. A sample initially collected for one project may prove valuable for many more studies. The Internet of Samples (iSamples) project aims to integrate large, diverse, cross-discipline sample repositories and enable access and discovery of material samples as FAIR data (Findable, Accessible, Interoperable, and Reusable). In this poster we report our recent progress in controlled vocabulary development and mapping.

## Repository Descriptions

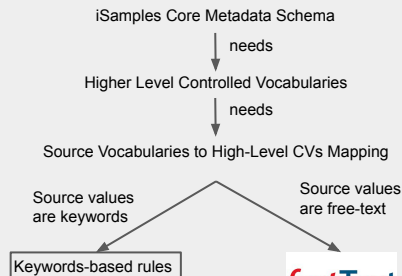
**SESAR** is a community platform that helps make Earth Sciences samples more discoverable, accessible, reusable and connects samples with the knowledge ecosystem derived from them.

**GEOME** is a web-based database that captures the who, what, where, and when of biological samples and associated genetic sequences.

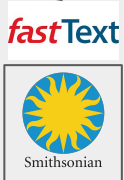
**Open Context** holds archaeology samples and goes beyond the archive by richly integrating the totality of your analyses, maps, media, and journals together so they can support your interpretations.

**Smithsonian** Institution is the world's largest museum, education, and research complex and holds natural history of biodiversity.

## Source Terminology to iSamples CVs Mapping



GeOME



Machine learning prediction based on habitat, locality, higher geography, and higher classification fields in biodiversity collections

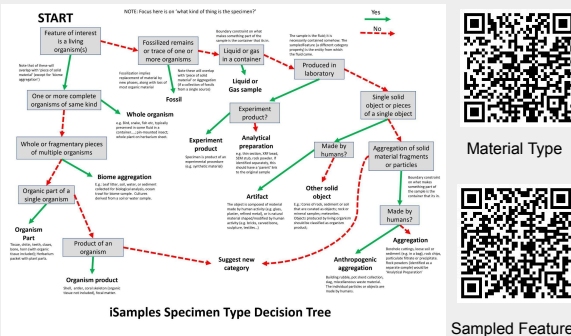
specimen type = liquid-aqueous in SESAR records was mapped to specimen type = liquid or gas sample and material type = liquid water.



iSamples  
The internet of samples

Applying these approaches, more than 3M records of the four large collections have been mapped successfully to a common core data model facilitating cross-domain discovery and retrieval of the sample records.

## Vocabulary Decision Tree Graphs

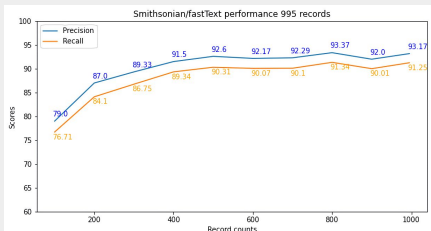


Material Type



Sampled Feature

## Training Sizes Impact fastText Performance



We trained models with the 100 records to 995 records and found the more training records are used to train an ML model, the higher precision and recall performances.

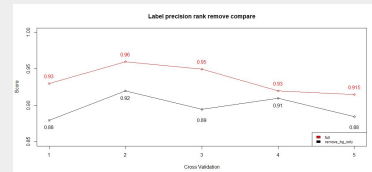
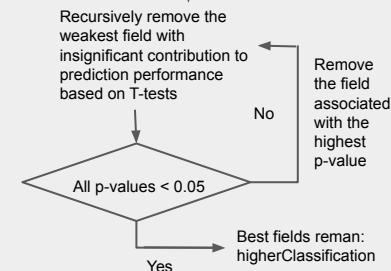
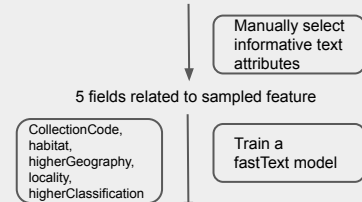
## Acknowledgement

Funded by the US National Science Foundation (CSSI)



## Feature Selection for Sampled Feature Prediction

Smithsonian biodiversity collection's 72 fields



The result showed only one attribute mainly influenced performance