

Dave Vieglais¹, Stephen M. Richard⁷, Quan Gan², Yuxuan Zhou², Hong Cui², Danny Mandel², Neil Davies³, John Deck³, Eric C Kansa⁴, Sarah Whitcher Kansa⁴, John Kunze⁴, Christopher Meyer⁵, Thomas Orrell⁵, Sarah Ramdeen⁶, Rebecca Snyder⁶, Ramona L. Walls⁶, Kerstin Lehner⁶
 1 University of Kansas; 2 University of Arizona; 3 University of California, Berkeley; 4 The Alexandria Archive Institute; 5 National Museum of Natural History, Smithsonian Institution; 6 Columbia University; 7 independent Contractor

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.

Controlled Vocabularies (CVs)

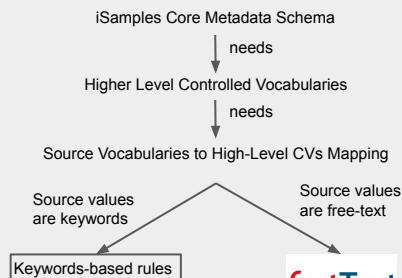
Material Type: What kind of material is the specimen?

Specimen Type: What kind of thing is the specimen?

Sampled Feature: What was the sample collected originally to represent?

The vocabularies provide **consistent semantics for high-level integration** of existing vocabularies used in the source collections

Source Terminology to iSamples CVs Mapping

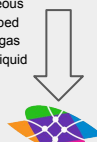


GeOME



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

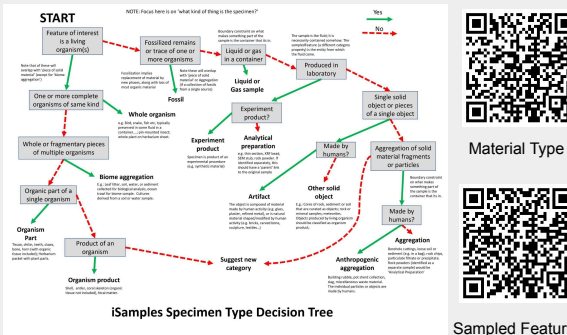
specimen type = liquid-aqueous in SESAR records was mapped to specimen type = liquid or gas sample and material type = liquid water.



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

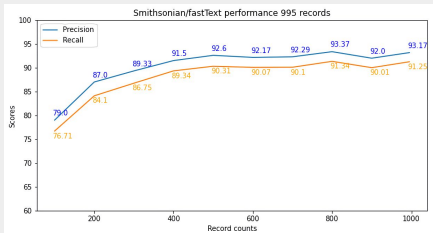


Material Type



Sampled Feature

Training Sizes Impact fastText Performance



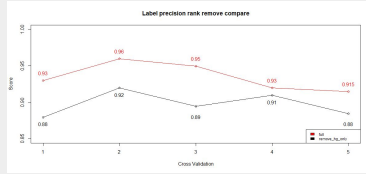
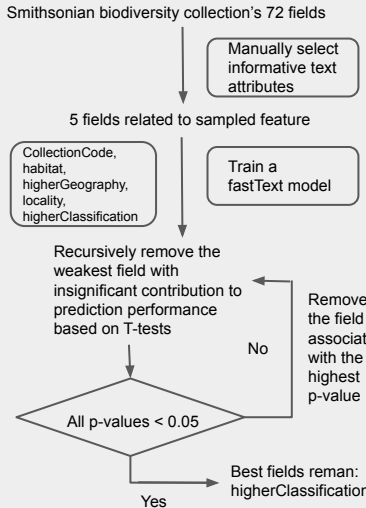
We trained models with the 100 records to 995 records and found the more training records are used to train an ML model, the higher precision and recall performances.

Acknowledgement

Funded by the US National Science Foundation (CSSI)



Feature Selection for Sampled Feature Prediction



The result showed only one attribute mainly influenced performance