



# Natural Language Processing and Machine Learning: Applying Advances in BioMedicine to Earth Science

Ruth Duerr, Skatje Myers, Martha Palmer, Chris Jenkins, Anne Thessen, James Martin



An NSF-Funded project at INSTAAR and CLEAR at the University of Colorado Boulder, the Smithsonian Institute, and Ronin Institute for Independent Scholarship

# What is ClearEarth ?

ClearEarth is a collaborative project aimed at bringing semantic technologies from the biomedical field into the earth-surface earth, ice and life sciences.

# Why Biomedicine ?

Biomedicine is very advanced in building large, deep semantic resources using Natural Language Processing (NLP) and Machine Learning (ML) techniques.

# Goals

## Overall

Enable use of advanced industry and research software within the geo-, bio- and cryospheric sciences for operations such as query and reasoning.

# Year 1

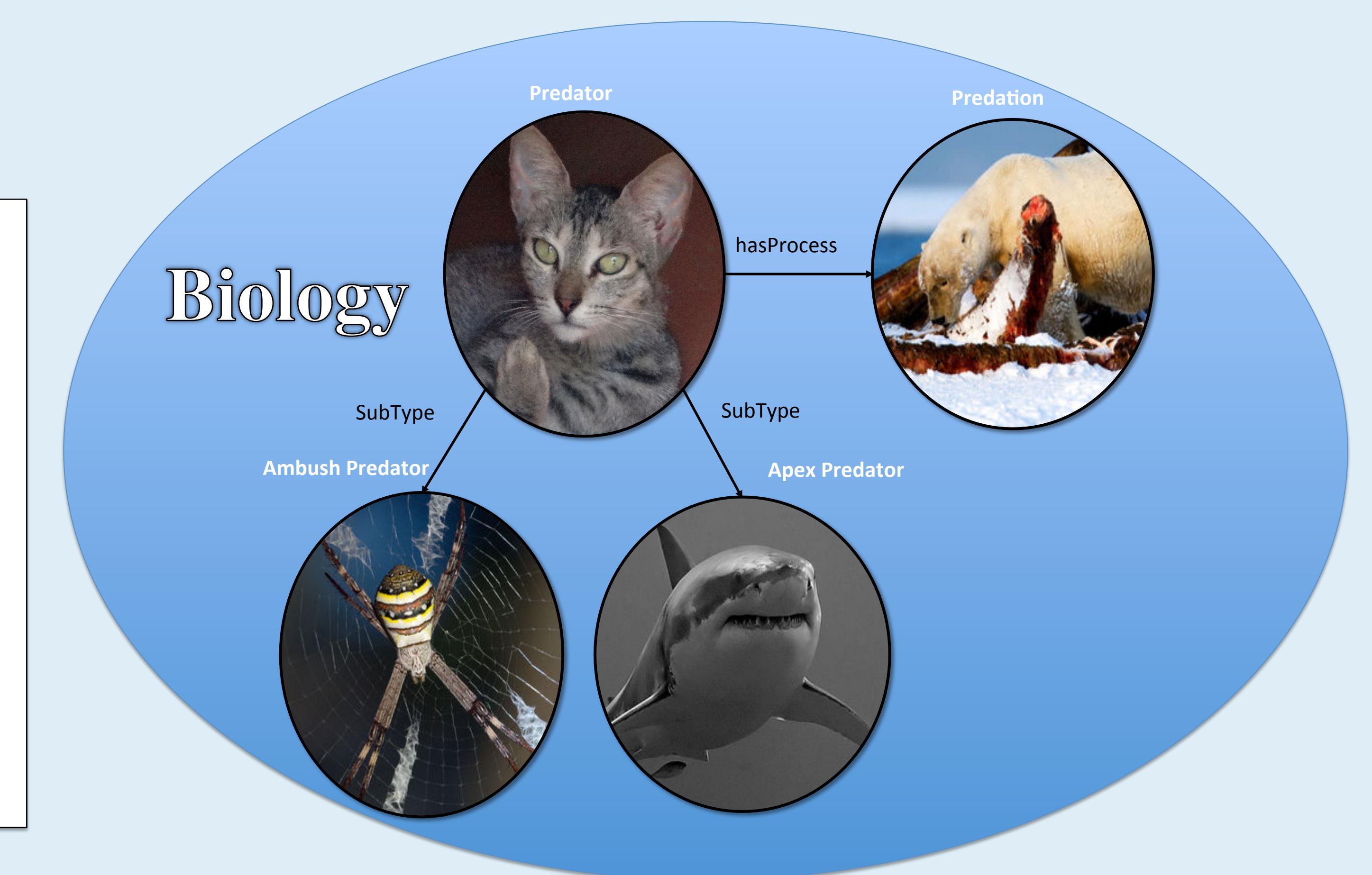
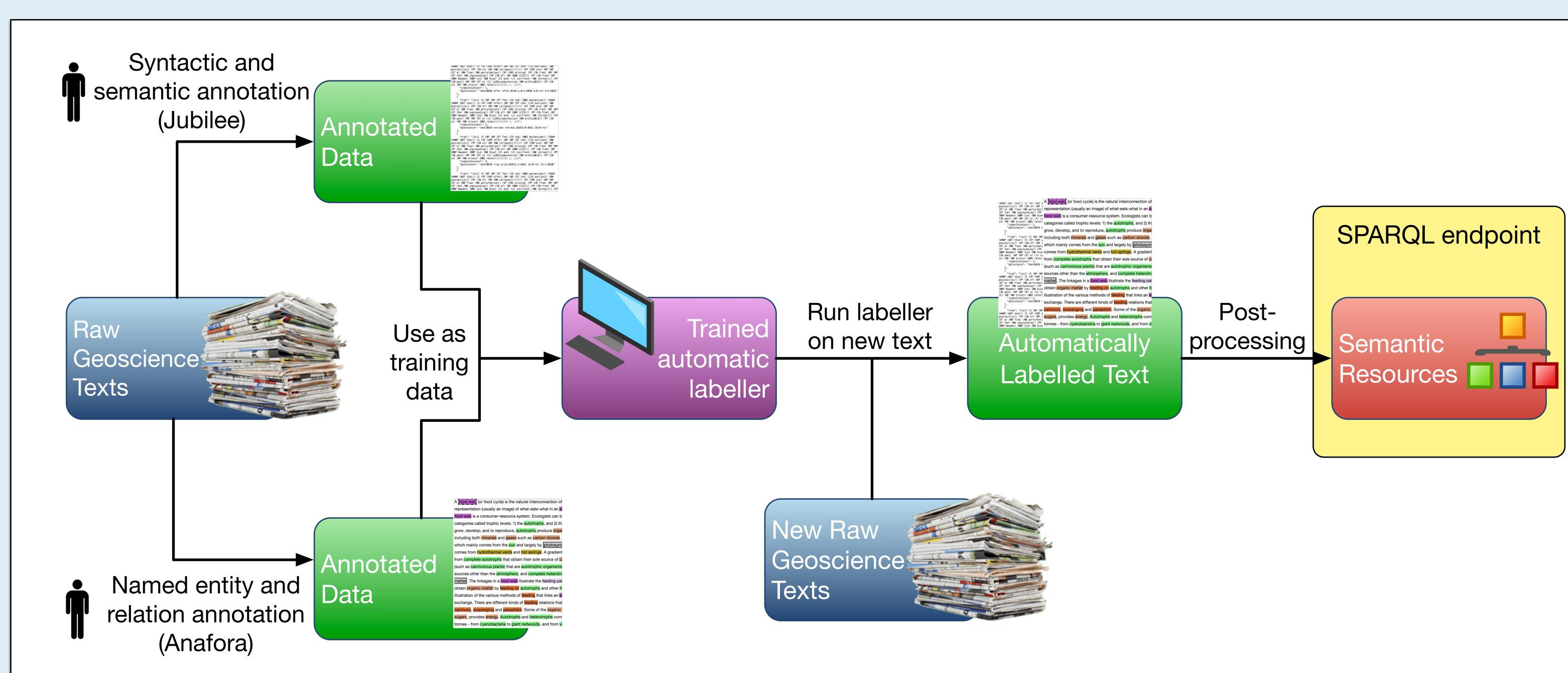
- Develop efficient and accurate methods of annotating earth-ice-biology texts with syntactic and semantic content.
  - Investigate pathways for successful adoption of the NLP and ML methods from the biomedical semantic technologies.

# Early Lessons

- Obtaining high quality, internally consistent annotation which is useful to machine learning algorithms is non-trivial
  - Thought structures and concepts used by domain scientists are different and messier than those required by machines
  - Interdisciplinary collaboration between highly disparate groups is challenging: terminologies and priorities
  - Genres of text are distinguished: for example, professional publication, popular article, newscast, twitter, field notes. We find all these in the geo-bio-cryo fields, just as much in biomedicine.

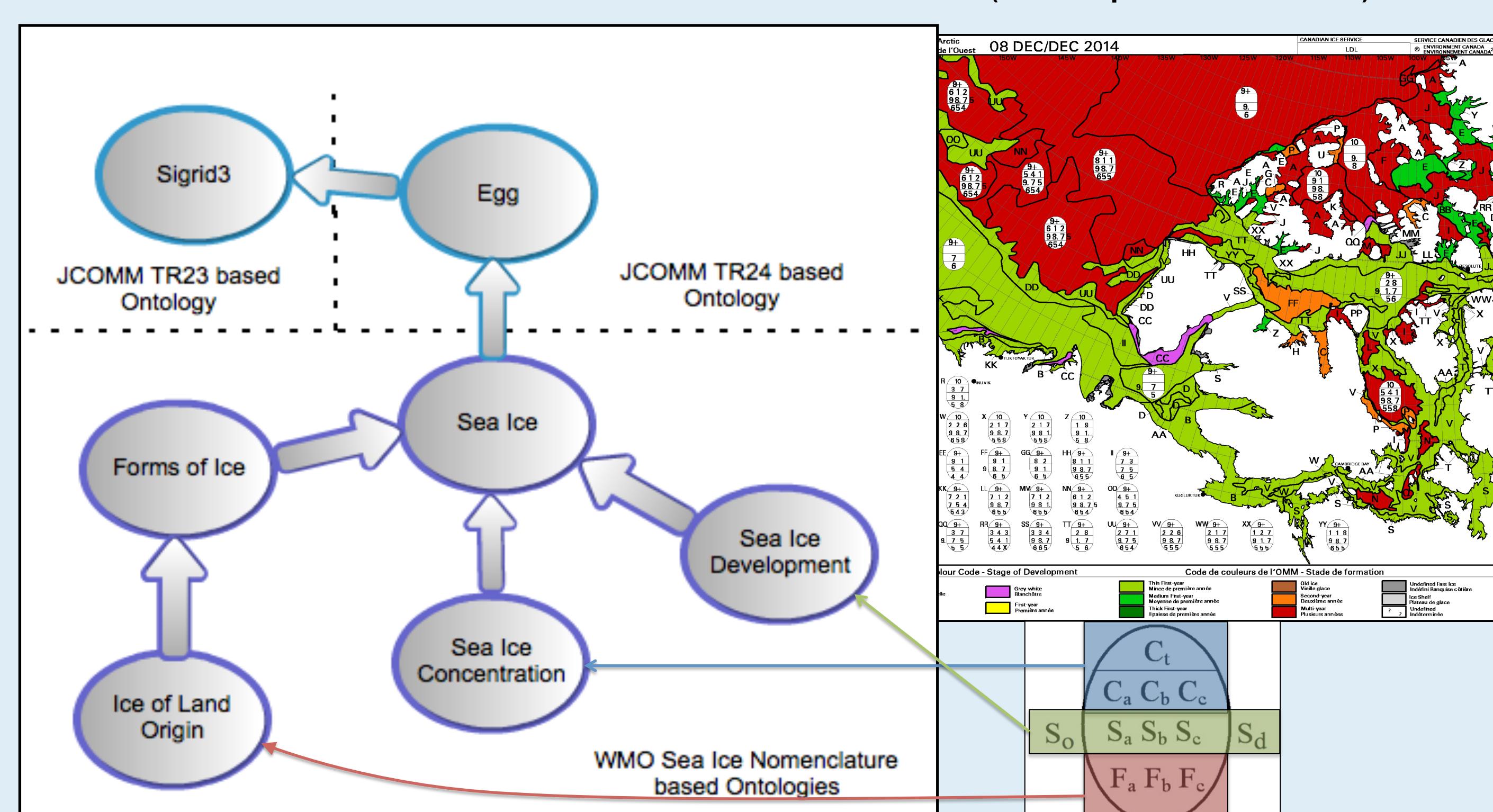
# Methods

The most successful current approach for automatically extracting semantic content from text is by using supervised machine learning.



# Validation

Verify that these approaches can re-create hand built ontologies developed where both methods use the same source documents. Using ontologies developed through the Semantic Sea Ice Interoperability Initiative funded Award ACI 0956010 INTEROP (as depicted below).

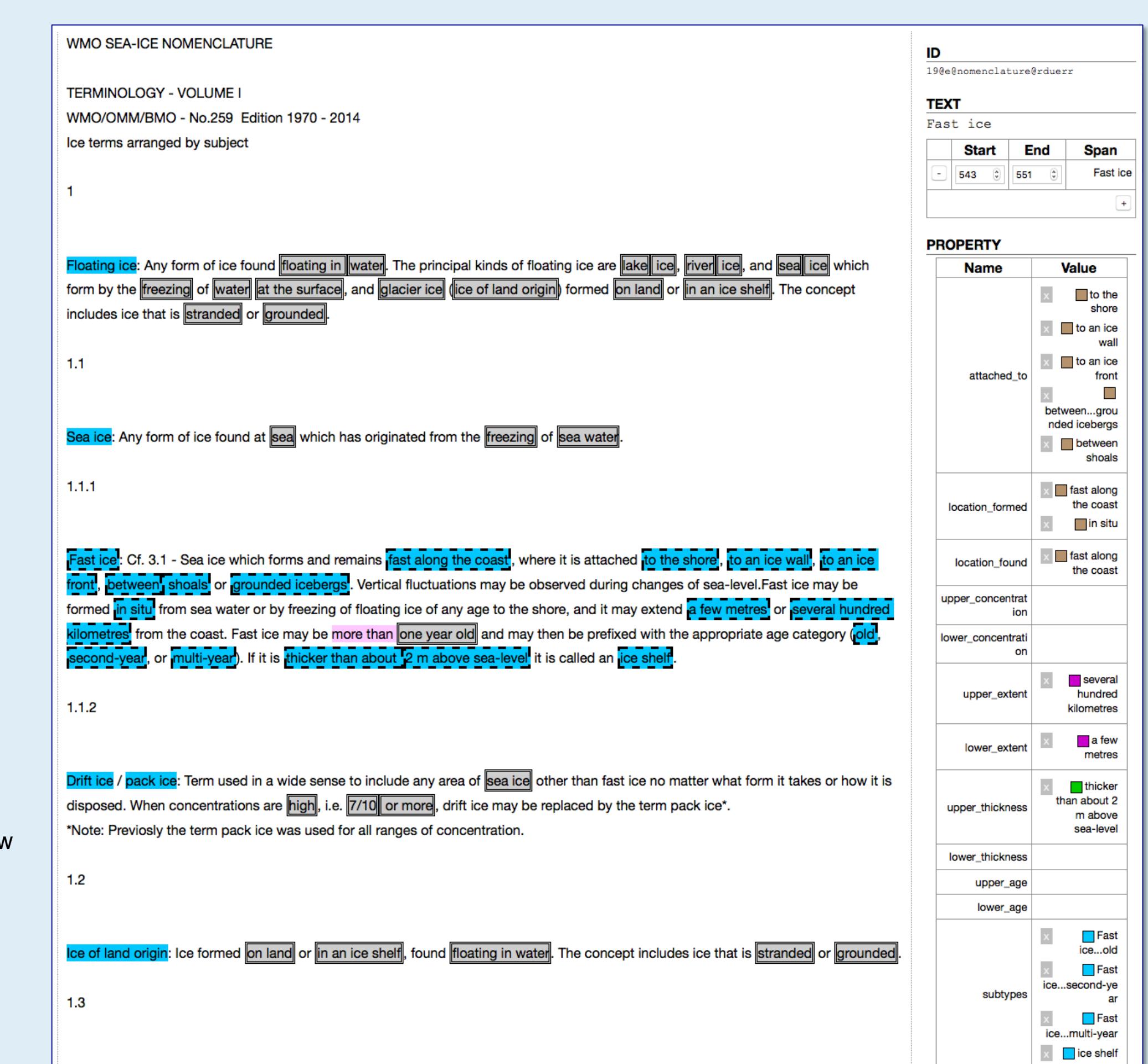


# Cryospheric Science



Beyond the WMO Sea Ice Nomenclature, other cryospheric glossaries are being gathered as input corpora for this project. For example, glossaries of permafrost<sup>1</sup>, snow terms<sup>2</sup>, glacier<sup>3</sup> and glacier mass balance<sup>4</sup> terminology, and general cryospheric glossary<sup>5</sup>.

**We welcome pointers to other glossaries!!!**



The extreme importance of the annotation (here, Anaphora) stage must be emphasized. The Machine Learning (ML) software ported from BioMedicine requires sufficiently large (large !) annotated corpus. Noisy entries work against it, so annotation must be careful. But once that threshold is reached, virtually all texts in the domain are treatable by the ML-NLP methods.