

Where Data Meets Meaning: Ontology Integration for Pathogen-Genomics Contextual Data

Rhiannon Cameron¹, Damion Dooley¹, Emma Griffiths¹, William Hsiao^{1,2}



Center for Infectious Disease Genomics and One Health
Simon Fraser University (SFU), Burnaby, BC, Canada

The Problem:

- The SARS-CoV-2 pandemic highlighted to the public health community the need to collect and compare viral genome contextual data.
- Non-standardized information systems across institutions result in datasets that are difficult to integrate and compare, as exemplified by the challenge of data harmonization within Canada's decentralized health system.
- These variations delay data integration and make consistent analysis difficult or impossible, ultimately hindering our ability to quantify public health threats and guide mitigation strategies.

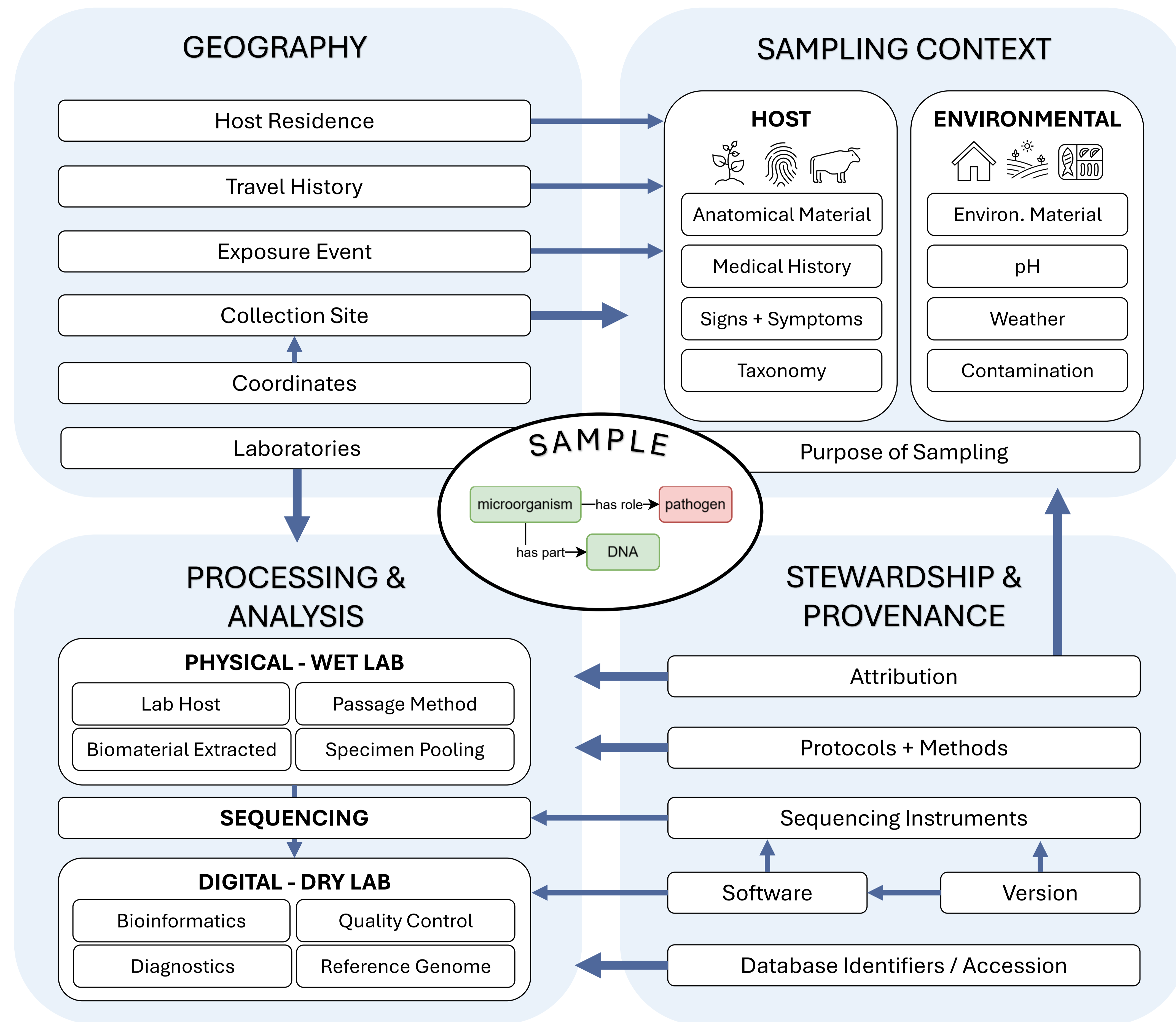
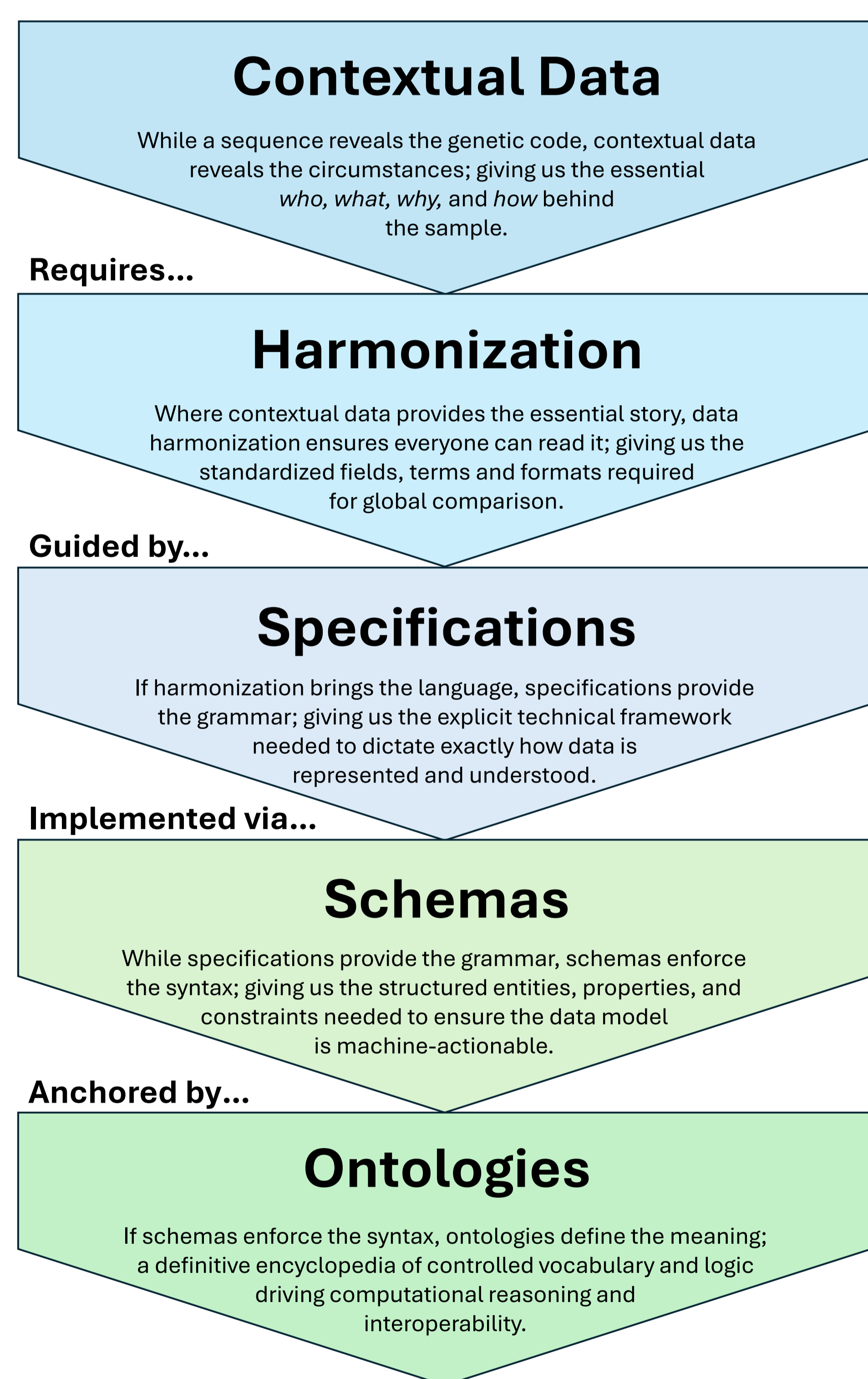


Figure 1: Sample-Centric Specification. An abstracted contextual map tracking the transformation of a physical sample from its environmental or clinical origin into a digital record.

Schema vs Ontology:

Feature	Schema	Ontology
Primary Role	Defines the structure	Defines the meaning
Validation	"Does the data fit the shape?"	"Does the data make logical sense?"
Key Components	Classes, slots, enumerations, values, and data properties	Classes, instances, URIs, semantic triplets, logic heritage
Implementation	LinkML ¹	OBO Foundry ² GENEPIO ³
Tooling Example	DataHarmonizer Templates ⁴	Knowledge Graph Database

¹ Linked Modelling Language (linkml.io/linkml)
² Open Biological and Biomedical Ontology Foundry (<http://www.obofoundry.org/>)
³ Genomic Epidemiology Ontology (<https://genepio.org/>)
⁴ The DataHarmonizer (<https://github.com/cidgoh/DataHarmonizer>)

Value of a LinkML Schema:

Ontologies can express schema-like elements but are designed for open-world semantic reasoning (deriving truths) rather than full schema definition (checking for completeness or validity). By contrast, LinkML is a schema language for defining structured data, validation rules, and interoperable models. Offloading schema concerns to LinkML provides several advantages:

- Structured data modeling;** including non-relational style patterns used by traditional databases and applications.
- Schema-driven generation;** source for generation artifacts (documentation, code, even ontology patterns).
- Multi-target conversions;** can be converted into multiple representations (e.g., JSON, SQL, OWL, RDF, etc.).
- Explicit validation;** offers sophisticated validation rules including patterns, ranges, and dynamic enumerations.
- Post-composition;** enables composite assertions that go beyond ontology axioms.

Logic of an Ontology:

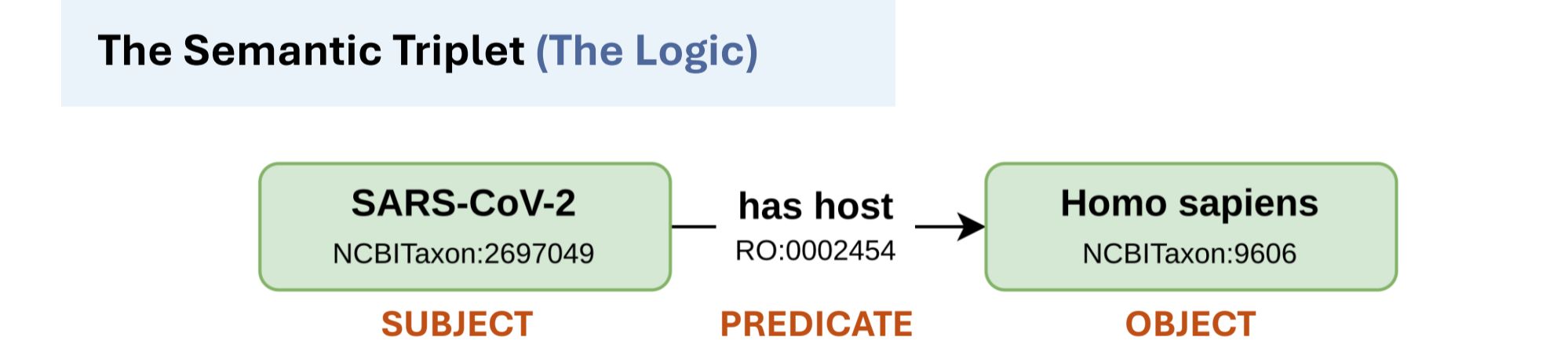


Figure: RDF triple. A machine-readable statement that can unambiguously reasoned over and queried. Subject/Object = Classes, Predicate = Relational Axiom.

Ontology Inheritance (The Hierarchy)

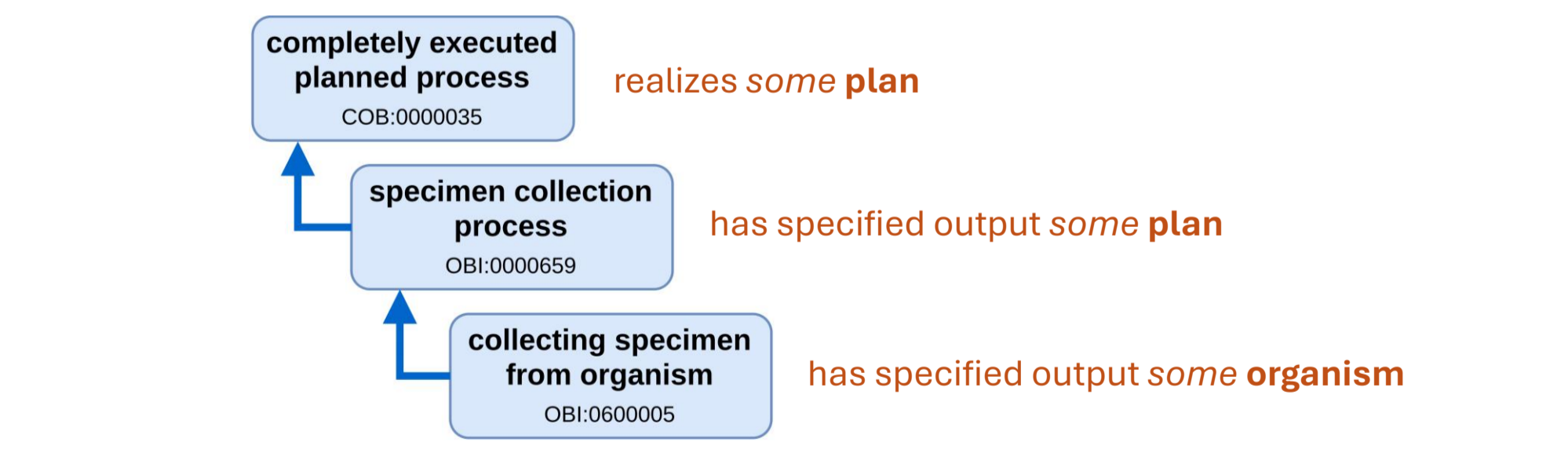


Figure: Inheritance. The "is a" hierarchy allows subclasses to automatically inherit parent axioms. For example, "collecting specimen from organism" inherits its parent's axioms, eliminating the need to re-code logic across classes.

Why You Should Care:

- Offload manual curation to computers, allowing automated validation to handle the heavy lifting of data formatting.
- Leverage existing data structures created and vetted by experts, rather than building models from scratch.
- Streamline integration across bioinformatics ecosystems via interchange mappings to databases, tools, and alternate models.
- Facilitate cross-jurisdictional comparisons by ensuring disparate datasets share the semantic meaning.
- Build scalable, community-adaptable data infrastructure that can evolve with dynamic stakeholder needs and emerging threats.

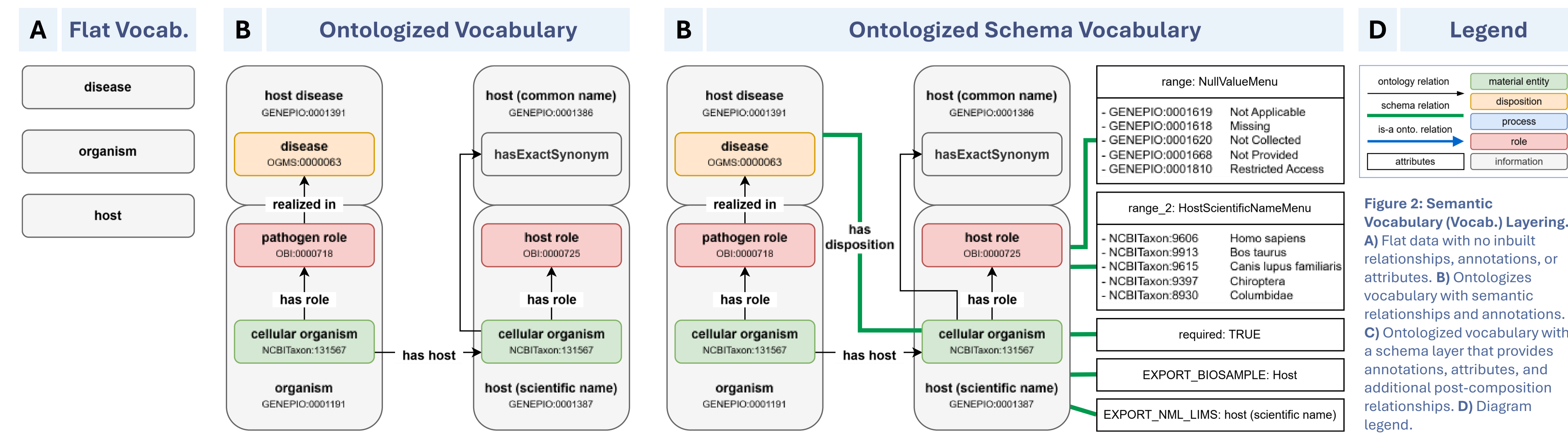


Figure 2: Semantic Vocabulary (Vocab.) Layering. A) Flat data with no inbuilt relationships, annotations, or attributes. B) Ontologizes vocabulary with semantic relationships and annotations. C) Ontologized vocabulary with a schema layer that provides annotations, attributes, and additional post-composition relationships. D) Diagram legend.

ACKNOWLEDGEMENTS: The authors respectfully acknowledges the unceded traditional territories of the Coast Salish peoples, including the səliwətaf (Tseil-Waututh), kwikwəḷəm (Kwkwetlem), Skwxwú7mesh Úxwumixw (Squamish) and x'məəθk'əyəm (Musqueam) Nations, on which SFU Burnaby is located.

The authors also acknowledge our provincial and national collaborators who continue to provide valuable feedback, along with other project collaborators and alumni: **Charlie Barclay, Anoocha Sehar, Madeline Iseminger** and the Public Health Alliance for Genomic Epidemiology (PHA4GE).

