

Context is Everything: Integrating Genomics, Epidemiological and Clinical Data Using GenEpiO

**Context Matters** 

On behalf of the GenEpiO Development Team Will Hsiao & Damion Dooley (BC Public Health Lab), Emma Griffiths & Fiona Brinkman (SFU)

#### Genomic sequences don't mean much without contextual information.

#### Sequencing & Bioinformatics



- Sequencing, Assembly Pipeline Parameters
- QA/QC Metrics
- Tree Construction Details

#### **Isolate Source**



#### **Clinical and Epi Details**



- Food, Clinical, Environment
- Food category, Body Product
- Dates, Location

- Demographics
- Host disease, Symptoms
- Lab Test Results
- Exposures

# When requisitioning metadata, you need to anticipate the needs of downstream users.

FAIR Principles of Data Management:

- F Findable
- A Accessible
- I Interoperable
- R Reusable

Published in Nature, March 2016 (BioScience experts, ELIXIR)



## Descriptive – Organized - Standardized





•Free text, short hand, granularity, misspelling, paper format

# Lack of standardization results in semantic ambiguity.

When Words Can Mean Different Things.



## Ontology, A Way of Structuring Information



- Standardized, well-defined hierarchy of terms
- Interconnected with logical relationships e.g. Endive "is\_a" Lettuce type

- Ontologies help resolve issues of taxonomy, granularity and specificity
- Reduces time consuming manual processing and mining

#### **Ontology Acts Like A Mapping Tool.**



- Humans AND computers can read it
- Mapping allows interoperability AND customization
- Facilitates data sharing, reproducibility

# Ontologies are different from program interfaces.

#### Ontologies:



• How information is structured and linked

#### Interfaces:



• How information is presented to users

# Data standards and ontologies help sort data at the source.

Metadata Challenges:

- Collection
- Organization
- Storage/Archiving



#### Sorting at the source...

#### benefits from clarity



# Prospective Metadata Collection is Easier and More Efficient than Retrospective Retrieval.



# Standardization of Contextual Information

Facilitate Reporting and Quality Control

- Reproducibility
- Reproducibility
- Reproducibility
- Reproducibility



# To develop a useful ontology, engaging the end users is your TOP priority.





#### GenEpiO: Combining Different Epi, Lab, Genomics and Clinical Data Fields.

See draft version at https://github.com/GenEpiO/genepio/wiki



Environmental), BioSample

#### Current Ontology Development Focuses On 3 Key Areas



# FoodON: A Farm-to-Fork Food Ontology

See draft version at https://github.com/FoodOntology/foodon

#### Aim:

To provide food descriptors for food items, ingredients, production environments to facilitate outbreak investigations, risk assessments, source attribution etc



Farm to Fork Continuum

#### **Resources:**

- LanguaL (USDA)
- FoodEx2 (EFSA)
- Codex Alimentarius (WHO)
- USDA Nutrient Database
- Painter Classification (CDC)
- FoodO/FooDB
- AGROVOC
- Food Safety Information Network
- Compendium of Analytical Methods (HC)
- More...

Ontologies are commonly encoded using OWL (Web Ontology Language).

- Markup language for sharing ontologies on the web
- Machine and human-readable
- OWL statements written in RDF (XML syntax)
- Protégé (editor) *oprotégé*
- Ontology lookup services:





Explore GenEpiO with Proofsheet

http://tinyurl.com/uiproofsheet



# GenEpiO and FoodON are now part of the OBO Foundry library of ontologies.

Open Biomedical Ontologies - http://www.obofoundry.org/

- Prescribes best practices for ontology development
- Committed to common use, interoperability, collaborative development
- Common relations and syntax
- Accessible definitions, good documentation

144 ontologies accepted or under development → Describing genes and phylogenies to diseases and anatomy





Ontology for the description of lifescience and clinical investigations

## GenEpiO will be implemented in different interfaces.



0

ightarrow

ullet

 $\bullet$ 

BioSample Table # 157 Record The "Copy to New" button causes an existing filled-in record to be copied to a new one, for minor changer between records in large batch job submissions. Selection list fields have an enhanced as-you-type lookup function. Right-click in the field to access this. Fields with an asterisk (\*) are mandatory. Your Keyword PC / Mac Key Shortcut submission will fail if any mandatory fields are not "missing" Alt/Option + Right arrow completed. If information is unavailable for any "not collected" Alt/Option + Down arrow mandatory field (or to clarify the data collection status "not applicable" Alt/Option + Left arrow of any field) enter the appropriate keyword to the right. Light yellow input fields are included in the NCBI Biosample submission table. Other fields (in beige) are stored in this spreadsheet only - they are NOT included in the Biosample submission data. Alt. Id. CFIA, St-Hy\_10116 Clinical O Environmental LiDS0164 Listeria monocytogenes LIDS0164 2011 Alt. Facility id Alt. Isolate id Facility id Isolate id CFIA, St-Hy\_10116 iDS0164 Past id 1 NML 14-2968 Past id 2 CFIA, St-Hy\_10116 Canadian Food Inspection Agency See list Institution code Collection code Specimen Id Specimen CFIA Culture Id Culture CFIA Isolate testing

**IRIDA Isolate NCBI Biosample Submission Form** 

Creates BioSample-Compliant Genome Submission Forms.  $\bullet$ 

# GenEpiO Will Help Integrate Genomics and Epidemiological Data.

	Genetic D	Genetic Distance									
Use computers to identify common exposures, symptoms etc among genomics clusters	Phylogenetic Tree	Cluster	Line List ID	Patient Name	Prov. Health No.	Age	Sex	Location	Sample ID	Collection Date	Culture Result
		A	1	John Smith	4513253244	26	м	Vancouver	F14231	14/03/21	Salmonella sp.
		A	2	Sally Smith	4519567458	24	F	Vancouver	F14235	14/03/21	Salmonella sp.
	₋┥_┌╴	в	3	Tom Jones	4517543216	35	м	Vancouver	M6542	14/03/24	Salmonella sp.
	ொட	в	4	Helen Jones	9856321124	35	F	Vancouver	S1245	14/03/22	Salmonella sp.
		с	5	Jennifer Lee	4516853122	29	F	Vancouver	S5642	14/03/22	Salmonella sp.
		с	6	Michael	9456534561	45	м	Victoria	T68954	14/03/25	Salmonella

Example: Automating Case Definition generation Correlate Genomics Salmonella Cluster A cases between 01 Mar 2015- 15 Mar 2015 with High-Risk Food Types Spinach → Leafy Greens and Geographical Location of Vancouver

#### Line List Visualizations of Selectable Data Based on GenEpiO Fields.





22

## Summary of the Advantages Genomic Epidemiology Ontology offers Public Health.

- 1. Eliminates semantic ambiguity
- 2. Term-mapping allows customization
- 3. Faster data integration

4. Standardized quality control and result reporting trigger actionable events in same way

Investigation power!

5. Reproducibility (accreditation, validation)





Genomic Epidemiology Ontology is like instrumentation for your contextual information... it needs maintenance and continual improvement



To achieve consensus and uptake  $\rightarrow$  International GenEpiO and FoodON Consortia (>60 members from 15 countries)

> Join us! E-mail: IRIDA-mail@sfu.ca

# Metadata management is crucial, but tricky.

AFTER

#### BEFORE



**GenEpiO** helps sort it out.

## Context is everything in foodborne outbreak investigations.



#### **GenEpiO** helps fit the data together.

Contact the GenEpiO Dev Team at: ontology-group-irida@googlegroups.com



#### **Project Leaders**

Fiona Brinkman – SFU Will Hsiao – PHMRL Gary Van Domselaar - NML

#### Simon Fraser University (SFU)

**Emma Griffiths Geoff Winsor** Julie Shay Matthew Laird **Bhav Dhillon** 

**McMaster University** Andrew McArthur Daim Sardar

**European Food Safety Agency** Leibana Criado Ernesto Vernazza Francesco **Rizzi Valentina** 

National Microbiology Laboratory (NML) Franklin Bristow Aaron Petkau **Thomas Matthews** Josh Adam Adam Olsen Tara Lynch Shaun Tyler Philip Mabon Philip Au Celine Nadon Matthew Stuart-Edwards Morag Graham **Chrystal Berry** Lorelee Tschetter Eduardo Toboada Peter Kruczkiewicz Chad Laing Vic Gannon Matthew Whiteside **Ross Duncan** Steven Mutschall

**University of Lisbon** João Carriço

**European Bioinformatics Institute** Melanie Courtot **Helen Parkinson** 

JBC

g Genome **Genome**Canada BritishColumbia

SFU **PHSA Laboratories** 



**BC Public Laboratory and BC Centre for Disease Control** (BCCDC) Damion Dooley Judy Isaac-Renton Patrick Tang Natalie Prystajecky Jennifer Gardy Linda Hoang Kim MacDonald

Yin Chang Eleni Galanis Marsha Taylor

**University of Maryland** Lynn Schriml

Jennifer Law

**Canadian Food Inspection Agency** (CFIA) Adam Koziol **Burton Blais Catherine Carrillo** 

**Dalhousie University** Rob Beiko Alex Keddy

#### www.irida.ca









Agence de la santé publique du Canada