# SARS-CoV-2 genomics contextual data curation and stewardship

**Emma Griffiths, PhD**

Chair, PHA4GE Data Structures Working Group

Hsiao Public Health Bioinformatics Lab

Faculty of Health Sciences, Simon Fraser University

Vancouver, Canada

Public Health Alliance for Genomic Epidemiology

# Outline

1. Challenges associated with genomics contextual data

2. Data standards: solutions for harmonization

3. PHA4GE SARS-CoV-2 specification
   - Resources for putting standards into practice

4. Data stewardship

5. Wrap up & links

Public Health Alliance for Genomic Epidemiology

# Contextual data is critical for interpreting the sequence data.

**Sequence data**



**Contextual data**

 Sample metadata

 Lab results

 Clinical/Epi data

 Methods

**Contextual data** (metadata) used for **surveillance** and **outbreak investigations**:

- **characterize** lineages and clusters
- identify variants with **clinical significance**
- correlate genomics trends with **outcomes, risk factors**
- **inform decision making** for public health responses and **monitor effects of interventions**

Public Health Alliance for Genomic Epidemiology

# Harmonizing fields of data is challenging.

*A field by any other name does NOT smell as sweet...*

SPECIMENSOURCE_1
Isolation
host_tissue_sampled
Source

The labs mean "**sample type**"

**Differences in labels, Same meaning**

*Computer doesn't recognize these as the same thing*

Source

The lab means "**submitting lab**"

**Same label, Different meaning**

*Computer doesn't recognize these as different*
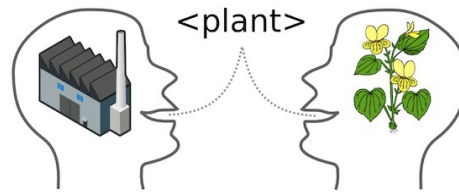
**...so, you can't just combine fields of data.**

Public Health Alliance for
Genomic Epidemiology

4

# Harmonizing information inside fields is challenging.

Free text = ☹

*Sever Acute Repiratry Sickness*

*SC2*

<plant>

Date:
2021-04-26
April 26, 2021
26-Apr-2021

MISSING
?

Errors

Short hand

Semantic ambiguity

Formats

Inconsistently collected

Public Health Alliance for Genomic Epidemiology

# Getting the right information to the right people is critical during health emergencies.

- Need to share data: **within** organization, with **trusted partners**, with **international** agencies/**public** repositories
- Data structure variability in local databases propagates to public repositories

### Private databases:

| Specimen Collected |
|---|
| ☐ Upper respiratory (e.g., Nasopharyngeal or oropharyngeal swab) |
| ☐ Lower respiratory (e.g., sputum, tracheal aspirate, BAL, pleural fluid) |

**6 - Specimen Type** (check all that apply)

**Specimen Collection Date:** yyyy / mm / dd     **(required)**

☐ NPS in UTM          **If possible:**

☐ Throat Swab in UTM          ☐ BAL

☐ Other (Specify):          ☐ Sputum

### Public databases:

| | |
|---|---|
| **isolate** | SARS-CoV-2/186197/human/2020/Malaysia |
| **collected by** | Universiti Malaya COVID Research group |
| **collection date** | 14-Mar-2020 |
| **geographic location** | Malaysia |
| **host** | Homo sapiens |
| **host disease** | COVID-19 |
| **isolation source** | Nasopharyngeal/throat swab |
| **latitude and longitude** | 3.1390 N 101.6869 E |

| | |
|---|---|
| **source name** | Lung sample from postmortem COVID-19 patient |
| **cell type** | Lung Biopsy |
| **strain** | NA |
| **subject status** | No treatment; >60 years old male COVID-19 deceased patient |

Different data structures make information less interoperable and more difficult to integrate.

That means you need to spend more time and resources to clean/re-structure information before you can use it.

Best practices for data **management/stewardship/structure** are critical parts of SARS-CoV-2 sequencing and analyses.

Public Health Alliance for Genomic Epidemiology

# Data standards: Solutions for standardizing contextual data

**Data dictionaries:**
- **Fields, flat list of terms, formats**
- Data dictionaries fulfill particular purposes
- Organization or project specific project (ease of use)

e.g.     Pizza Type
            Cheese pizza
            Meat Lover's Pizza
         Veggie Pizza

**Minimum Information Checklists:**
- **Prescribed fields to describe something (in a particular context)**
- created by authoritative source
- core fields common between checklists (e.g. collected by, sample collection date, sample ID)
- specific packages provide fields for specific contexts (e.g. air, human gut)
- commonly used by international sequence repositories

e.g.     MIAP (Min Info About a Pizza)
         Pizza Type:
         Topping:
         Crust:
         Sauce:

Public Health Alliance for Genomic Epidemiology

8

# Ontology, A Way of Structuring Information
## Ontologies aim to represent truth. *Is this universal?*

- Controlled (standardized) vocabulary
- Hierarchy (granularity)
- Logic
- Definitions and unique IDs (disambiguation)

> e.g. Veggie Pizza (FoodOn:1234)
> A type of pizza that is topped with only vegetable and cheese toppings.

> e.g. Margarita – synonym of Margherita?
> IDs distinguish between pizza and cocktail

- Synonyms (facilitates mapping)
- Prepares data for more variety of analyses

Veggie pizza (FoodOn:1234)

hasRecipe
hasCookingTime
hasPreparationDate



Margherita pizza (FoodOn:1240)
synonym: Margarita pizza

Margarita cocktail (FoodOn:2376)
Synonym: Margarita

Public Health Alliance for Genomic Epidemiology

9

# Contextual data ecosystem: putting standards into practice

# The SARS-CoV-2 Contextual Data Standard

**SARS-CoV-2 Domain Content**

- Repository accession numbers and identifiers
- Sample collection and processing
- Host information
- Host exposure information
- Host reinfection information
- Host vaccination information
- Sequencing methods
- Bioinformatics and quality control metrics
- Lineage and variant information
- Pathogen diagnostic testing details
- Provenance and attribution

**Data Sources**

- Case report forms
- Public repository requirements
- Existing metadata standards
- Literature

**Mapping to Standards**

- MIxS  5.0
- MIGS Virus, Host-Associated
- Project/Sample Application Standard
- OBO Foundry Ontologies

Public Health Alliance for Genomic Epidemiology

# Putting standards into practice: Template and standard terminology



- **Standardized collection template** (colour-coded, yellow=required, purple=recommended, white=optional)
- **Pick lists**: standardized terms
- **Structured formats** e.g. for dates
- **JSON schema**

# Guidance documentation



- **Reference guide**: field labels, definitions, guidance, expected values

- **SOP**: how to curate contextual data

13

# PHA4GE standard quick FAQ

**Do I have to fill in the whole thing?**
*NO! Only use the parts you need. We've highlighted the most important bits.*

**Is this just for human/clinical samples?**
*NO! It's for ALL samples.*

**Do I have to share all my contextual data?**
*NO! It's all up to you!*

**What happens if your pick lists don't have the term I want?**
*1. Get in touch with us!*
*2. SOP shows you how to find a standardized term.*

Public Health Alliance for
Genomic Epidemiology

# Worked Examples

- **State X reference lab** in **Country Y** has a **border testing program**.
- A **36 year old female** from **Canada visiting Country Y** tests positive for **SARS-CoV-2**.
- An **NP swab** was collected on **March 16 2021**.
- The sample was sequenced by **State X reference lab** using an amplicon strategy according to the **ARTIC V4** protocol.
- The consensus genome was generated **iVar 1.3.2**.

specimen collector sample ID: ABCD1234
sample collected by: State X Reference Laboratory
sequence submitted by: State X Reference Laboratory
sample collection date: 2021-03-16
geo_loc name (country): Country Y
geo_loc name (state/province/territory): State X
purpose of sampling: Diagnostic testing
purpose of sampling details: Not Provided
purpose of sequencing: Travel-associated surveillance
purpose of sequencing details: border testing program
organism: Severe Acute Respiratory Syndrome Coronavirus 2

anatomical part: Nasopharynx (NP)
collection device: Swab
host (scientific name): Homo sapiens
host age: 36
host age unit: year
host gender: Female
host residence geo_loc name (country): Canada
travel history: individual travelled directly from Canada
amplicon pcr primer scheme: ARTIC V4
amplicon size: 400bp
consensus sequence software name: iVar
consensus sequence software version: 1.3.2

**Null values:**
Missing
Not Applicable
Not Collected
Not Provided
Restricted Access

Public Health Alliance for Genomic Epidemiology

# Worked Examples

- **State X reference lab** in **Country Y** in investigating an **outbreak** at **Hospital Z**.
- A previously **vaccinated** individual (**one dose AZ, Feb 1 2021**) tests positive for SARS-CoV-2, with a diagnostic PCR **CT value of 23 (E gene)**.
- The sample was sequenced by State X reference lab using an **Illumina MiSeq** (**Illumina prep kit**).
- The sequence data was **filtered, processed and dehosted** using the **SIGNAL pipeline**. The consensus genome was generated using **FreeBayes 1.2.3**.
- Using **Pangolin**, the genome was identified as a **Delta VOC**.

specimen collector sample ID: ABCD1234
sample collected by: State X Reference Laboratory
sequence submitted by: State X Reference Laboratory
geo_loc name (country): Country Y
geo_loc name (state/province/territory): State X
purpose of sampling: Diagnostic testing
purpose of sequencing: Cluster/Outbreak investigation
purpose of sequencing details: outbreak at Hospital Z
organism: Severe Acute Respiratory Syndrome Coronavirus 2
host (scientific name): Homo sapiens
exposure setting: Hospital
host role: Patient
host vaccination status: Partially vaccinated
vaccine name: AstraZeneca COVISHIELD COVID-19 vaccine (ChAdOx1-S)
number of vaccine doses received: 1
first dose vaccination date: 2021-02-01

sequencing instrument: Illumina MiSeq
library preparation kit: Illumina Prep Kit
raw sequence data processing method: https://github.com/phac-nml/covid-19-signal-nml
dehosting method: https://github.com/phac-nml/covid-19-signal-nml
consensus sequence software name: FreeBayes
consensus sequence software version: 1.2.3
lineage/clade name: B.1.617.2
lineage/clade analysis software name: Pangolin
lineage/clade analysis software version: 3.1.4
variant designation: Variant of Concern (VOC)
gene name 1: E gene (orf4)
diagnostic pcr Ct value 1: 23

# Protocols to mobilize harmonized data



- **7 public repository submission protocols (GISAID, NCBI, EMBL-EBI) on Protocols.io**
- **PHA4GE-adapted submission forms**
- **instructional videos**

**Different repositories have different fields, but PHA4GE helps standardize what goes into those fields.**

https://www.protocols.io/workspaces/pha4ge

Public Health Alliance for Genomic Epidemiology

# Data transformation tools: The DataHarmonizer

- Tool for data entry and validation
- Spreadsheet-style text editor application (PHA4GE template)
- Picklists, data structure, validation, export
- Guidance (reference guide), SOP

View all fields
View required fields
Move to desired field
Automate column fill

Validate
(check for errors
or missing info)

Learn your
way around
the system

Double click on
field labels for
guidance on how
to fill them

Save
Open existing file
Export to chosen format

| File ▾ | Settings ▾ | Validate | Help ▾ | Template | PHA4GE | Loaded file | |

| | Database Identifiers | Sample collection and processing | | | | |
|---|---|---|---|---|---|---|
| | specimen collector sample ID | sample collected by | sequence submitted by | sample collection date | geo_loc_name (country) | geo_loc_name (state/province/territory) | organism |
| 1 | | | | ▾ | ▾ | | |
| 2 | | | | ▾ | ▾ | | |

PUBLIC HEALTH BIOINFORMATICS

Public Health Alliance for Genomic Epidemiology

https://github.com/Public-Health-Bioinformatics/DataHarmonizer/wiki/PHA4GE-SARS-CoV-2-Template

18

# Data transformation tools: The DataHarmonizer



- Enter data once, export in different submission formats i.e. GISAID and NCBI (BioSample, SRA, GenBank)

- Complementary to **multiSub** (interchange between submission formats)
- https://github.com/maximilianh/multiSub

https://github.com/Public-Health-Bioinformatics/DataHarmonizer/wiki/PHA4GE-SARS-CoV-2-Template

# Data stewardship: oversight and practices to ensure data is **accessible, usable, safe, trusted.**

**Privacy protection (sharing):**

- Public trust essential, loss of trust has consequences (protection, transparency)
- De-identified data (no names/addresses)
- Be careful of 1) geographical granularity, 2) small case numbers in defined geo_loc/time, 3) combinations of fields
- Age/gender often shareable, sharing of clinical/epi information may be restricted at individual level (e.g. hospitalization, exposures)
- Track identifiers (chain of custody), but personal health IDs may be considered PHII
- Consult privacy officer (jurisdictional policies)
- *See curation SOP for details…*

*Contextual data can provide critical information in investigations of life or death situations. Please consider when evaluating privacy concerns (see GA4GH's Human rights argument for prioritizing sharing).*
https://www.ga4gh.org/genomic-data-toolkit/regulatory-ethics-toolkit/framework-for-responsible-sharing-of-genomic-and-health-related-data/

Public Health Alliance for Genomic Epidemiology

# Data stewardship: oversight and practices to ensure data is **accessible, usable, safe, trusted**

**FAIR stewardship best practices (Wilkinson et al (2016), Nature):**
- Findable, Accessible, Interoperable, Reusable
- Standards help preserve integrity and meaning of data (now and in future) for YOU and others
- Machine-actionable → indexed, searchable
- Formal, broadly applicable language (data standards)

**Security & Quality:**
- Provenance, methods (rich details) → attribution, auditability, reproducibility (track methods), accountability
- Contextual data may require storage with higher security than seq data
- Errors corrected, update as required

Public Health Alliance for Genomic Epidemiology

# Summary: What can the PHA4GE spec do for you?

✔ 1. Data is more interpretable by <span style="color:red">humans and computers</span>

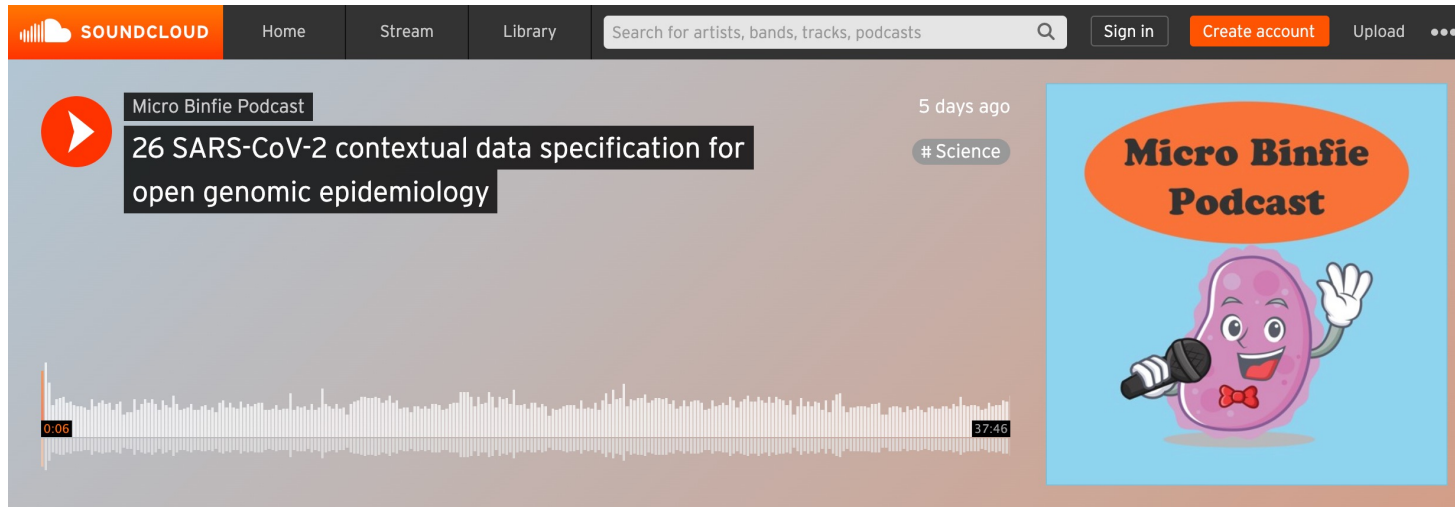✔ 2. <span style="color:red">Future-proof</span> contextual data

✔ 3. Harmonize and integrate data across labs/databases (<span style="color:red">interoperability</span>)

✔ 4. <span style="color:red">Tools</span>

Public Health Alliance for Genomic Epidemiology

# Learn more…



https://soundcloud.com/microbinfie/26-sars-cov-2-metadata#t=0:00

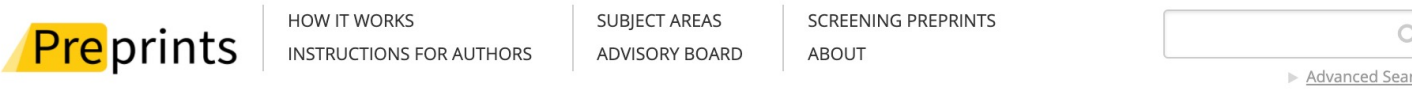## Listen to episode 26 Micro Binfie podcast



preprints.org > doi: 10.20944/preprints202008.0220.v1

Preprint    Article    Version 1    This version is not peer-reviewed

### The PHA4GE SARS-CoV-2 Contextual Data Specification for Open Genomic Epidemiology

Emma J. Griffiths * , Ruth E. Timme , Andrew J. Page , Nabil-Fareed Alikhan , Dan Fornika , Finlay Maguire , Catarina Inês Mendes , Simon H. Tausch , Allison Black , Thomas R. Connor , Gregory H. Tyson , David M. Aanensen , Brian Alcock , Josefina Campos , Alan Christoffels , Anders Gonçalves da Silva , Emma Hodcroft , William W.L. Hsiao , Lee S. Katz , Samuel M. Nicholls , Paul E. Oluniyi , Idowu B. Olawoye , Amogelang R. Raphenya , Ana Tereza R. Vasconcelos , Adam A. Witney , Duncan R. MacCannell

Version 1 : Received: 7 August 2020 / Approved: 9 August 2020 / Online: 9 August 2020 (15:53:58 CEST)

https://www.preprints.org/manuscript/202008.0220/v1

## Read our preprint
Update coming out soon!

Public Health Alliance for Genomic Epidemiology

# Special thanks to...

**Data Structures Team**
Ana Ribeiro de Vasconcelos
Josefina Campos
Idowu Olawoye
Paul Oluniyi
Adam Witney
Andrew Page
David Aanensen
Ines Mendes
Emma Hodcroft
Simon Tausch
Allison Black
Ruth Timme
Greg Tyson
Mike Feldgarden
Lee Katz
Brian Alcock
Amos Raphenya
Finlay Maguire
Dan Fornika
Duncan MacCannell

**Specification Contributors & Partners**
Nabil-Fareed Alikhan
Alan Christoffels
Will Hsiao
Sam Nicholls
Tom Connor
Anders Gonçalves da Silva
Dominique Anderson
Danny Park
CDC TOAST Team

**Steering Committee & Secretariat**
Jamie Southgate
Alecia Naidu
Rangarirai Matima
Nawaal Nacerodien
Peter van Heusden
Kevin Libuit
Nicola Mulder
Nicki Tiffin

Find the spec package:
https://github.com/pha4ge/SARS-CoV-2-Contextual-Data-Specification

Find us:
https://www.pha4ge.org
@pha4ge
datastructures@pha4ge.org

Public Health Alliance for Genomic Epidemiology