# Overcoming challenges of SARS-CoV-2 genomics data harmonization and integration for public health surveillance, outbreak investigations and research using the PHA4GE SARS-CoV-2 contextual data specification

**Emma Griffiths, PhD**

Hsiao Public Health Bioinformatics Lab

Department of Pathology and Laboratory Medicine

University of British Columbia

Vancouver, Canada

PUBLIC HEALTH BIOINFORMATICS

Public Health Alliance for Genomic Epidemiology

# Sequencing SARS-CoV-2 genomes has helped track the spread of the virus worldwide.

**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

# **Challenge:** How data is encoded impacts how it can be integrated and used for analyses.

Specimen:

| 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

Patient Setting:

**7 - Patient Setting / Type**

☐ Assessment Centre    ☐ Family doctor/clinic    ☐ Outpatient/ER not admitted

Only if applicable, indicate the group:

☐ Healthcare worker      ☐ Institution / all group living settings

☐ Inpatient (hospitalized)      ☐ Confirmation (for use **ONLY** by a COVID testing lab). Enter your result (NEG/POS/or IND)

☐ Inpatient (ICU/CCU)

☐ First Nations / Inuit

☐ Unhoused / shelter      ☐ For clearance of disease

☐ ER - to be hospitalized      ☐ Other (Specify):

☐ Deceased / Autopsy

☐ Acute care facility      ☐ Long term care facility

☐ Group home (community living)

☐ Correctional facility      ☐ School or daycare

☐ Assisted living      ☐ Independent living

☐ Other residential facility type, specify: _____

☐ Shelter      ☐ Conference

Public Health Alliance for Genomic Epidemiology

3

# 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

| | |
|---|---|
| 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 |

Compare these public datasets:
- different **standards**
- different **granularity**
- **free text**
- different **formats**
- different **interpretation**
- **privacy** concerns
- **methods**

Public Health Alliance for Genomic Epidemiology

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.

Public Health Alliance for Genomic Epidemiology

Data Structures  |  Bioinformatic Pipelines and Visualizations  |  Infrastructure
Public Repositories  |  Reference, QC and Validation  |  Workforce Development
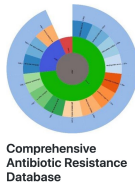Data Sharing and Ethics | Users and Applications

https://www.pha4ge.org    |    https://www.github.com/pha4ge    |    @pha4ge

# Data Structures Working Group

- 21 Members (9 countries, 4 continents)

**Goal: develop/promote data standards**
- COG-UK, SPHERES, CanCOGeN, Latin American Genomics Network
- identified need for fit-for-purpose contextual data standard for SARS-CoV-2 genomics

# The SARS-CoV-2 Contextual Data Specification Package

**PHA4GE SARS-CoV-2 Full Specification Content**
Repository accession numbers and identifiers
Sample collection and processing
Host information
Host exposure information
Sequencing methods
Bioinformatics and quality control metrics
Pathogen diagnostic testing details
Provenance and attribution



Shared with Trusted Partner

Full Metadata Set

Publicly Shared Metadata

Not Shared

Public Health Alliance for Genomic Epidemiology

# Template and terminology



- **Standardized collection template** (colour-coded)
- **Pick lists**: standardized terms
- **Reference guide**: field labels, definitions, guidance, expected values

https://github.com/pha4ge/SARS-CoV-2-Contextual-Data-Specification

Public Health Alliance for Genomic Epidemiology

9

# Supporting documentation



- **SOP**: how to use specification, find new terms, highlight practical/ethical/privacy issues
- **Field mapping to existing standards**: highlight alignment and gaps
- **JSON schema**: machine readable

https://github.com/pha4ge/SARS-CoV-2-Contextual-Data-Specification

Public Health Alliance for Genomic Epidemiology

# Protocols to mobilize harmonized data



- **7 public repository submission protocols** (GISAID, NCBI, EMBL-EBI) on **Protocols.io**

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

# Uptake: early adopters

# Want to know more?

## 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 Si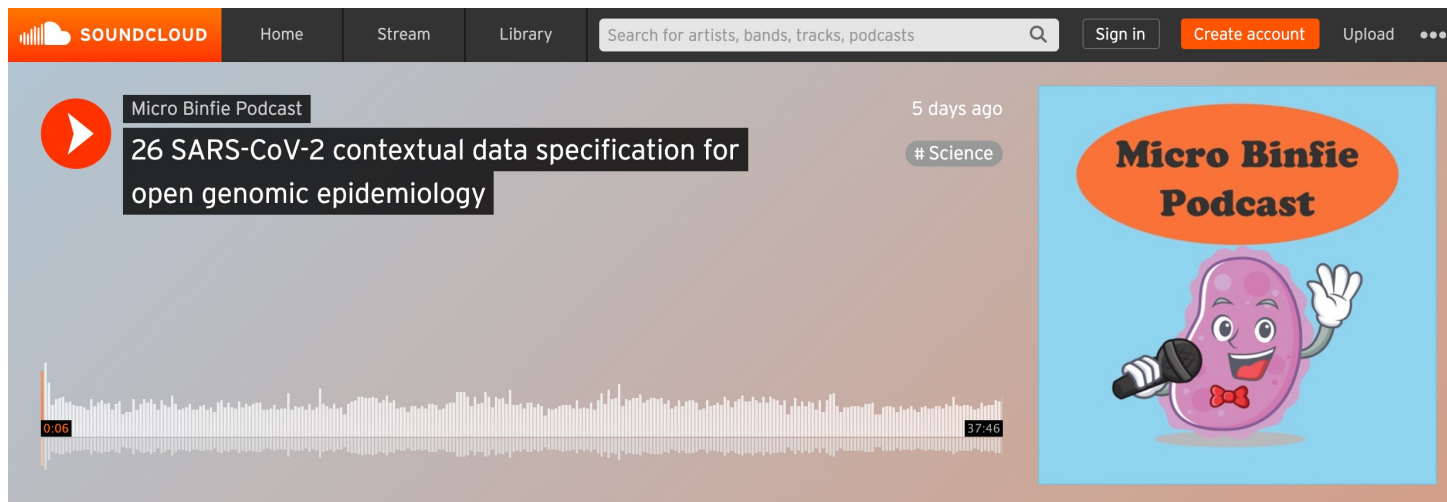lva, Emma Hodcroft, William W.L. Hsiao, Lee S. Katz, Samuel M. Nicholls, Paul E. Oluniyi, Idowu B. Olawoye, Amogelang R. Raphenya, Ana Tereza R. Vasconcelos, Adam A. Witney, Duncan R. MacCannell

Get PDF

Cite

Stay tuned for workshops coming soon!

Public Health Alliance for Genomic Epidemiology

# Summary

- contextual data is critical for **interpreting** genomics data/analyses
- specification structures metadata **consistently across labs**
- share with **public** repos, **trusted partners**, use for more **efficient private analyses**
- **future-proof** metadata

- <span style="color:red">success depends on community uptake</span>

<span style="color:red">Ongoing challenges: trust, equitable data sharing, better understanding of risk, co-ordination, sharing mechanisms</span>
<span style="color:red">→ **public health genomics ecosystem**</span>

Public Health Alliance for
Genomic Epidemiology

Thank you to
the Data Structures WG & friends,
PHA4GE consortium,
The Bill and Melinda Gates Foundation,
and to you!

Public Health Alliance for
Genomic Epidemiology

**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

**Steering Committee & Secretariat**
Jamie Southgate
Alecia Naidu
Rangarirai Matima
Nawaal Nacerodien
Peter Van Heusden
Danny Park

This work would not be possible without the contributions and dedication of these wonderful people.

Find us:
https://www.pha4ge.org
https://www.github.com/pha4ge
@pha4ge

Public Health Alliance for Genomic Epidemiology