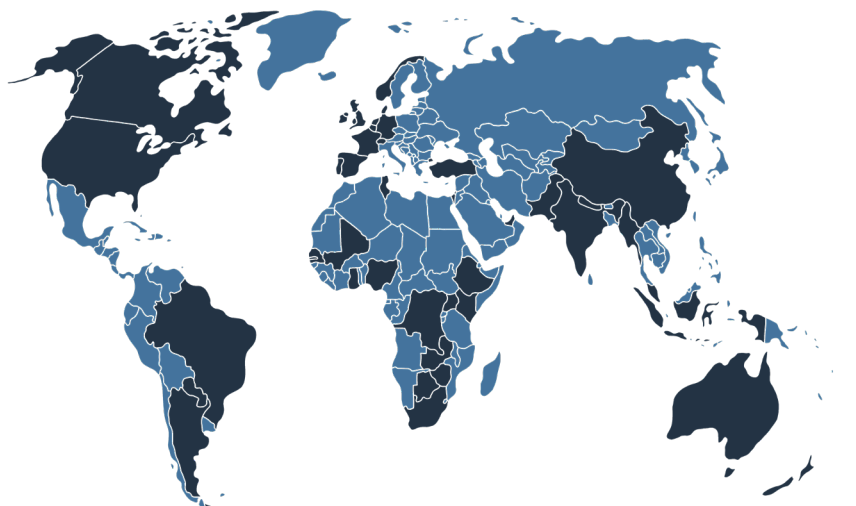


# PHA4GE

## Bridging the gap between Public Health & Bioinformatics

Alan Christoffels

South African National Bioinformatics Institute, University of the Western Cape



<https://www.pha4ge.org>

| <https://www.github.com/pha4ge>

|  @pha4ge

1

## **Vision**

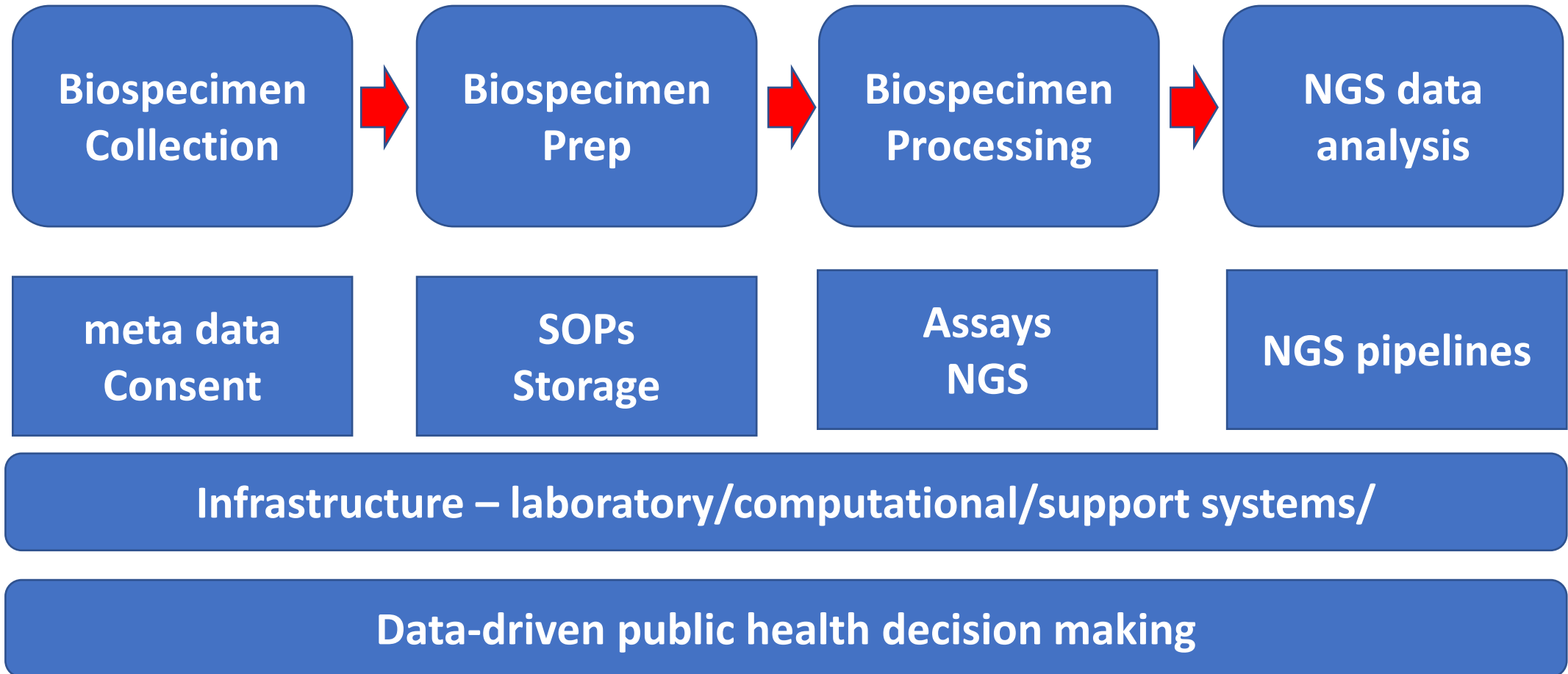
**a modular and scalable open-source platform** for public health bioinformatics. Architecture, APIs, ontologies, standards.

**Need for a community-driven open model, for public health bioinformatics**

**Critical incentive gap between academic software development and public health applications**

**2019 - 2024**

# PHA4GE embedded in a Public Health setting



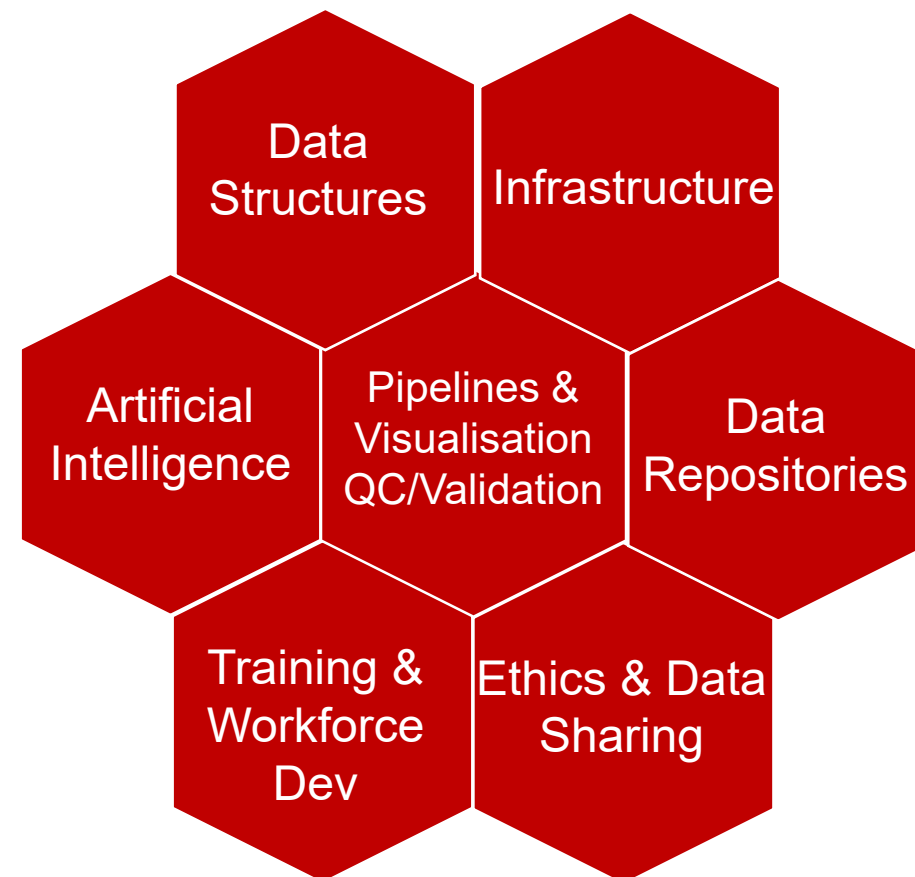
# Strategic Actions

Sustainability

Lead best practice  
implementation for public  
health laboratories

Human capital development

## Technical Working Groups:



<https://www.pha4ge.org>

<https://www.github.com/pha4ge>

 @pha4ge

4



# Data Structures Working Group



**Emma Griffiths**

Simon Fraser  
University



**Finlay Maguire**

Dalhousie University

<https://www.pha4ge.org>

| <https://www.github.com/pha4ge>

|  @pha4ge

5

# Contextual data is critical for interpreting SC2 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**

# Data structure variability in local databases propagates to public Repo

## Private databases:

Specimen Collected
<input type="checkbox"/> Upper respiratory (e.g., Nasopharyngeal or oropharyngeal swab)
<input type="checkbox"/> 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

☐ Throat Swab in UTM

☐ Other (Specify):

If possible:

☐ BAL

☐ 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	<a href="#">Malaysia</a>
host	Homo sapiens
host disease	COVID-19
isolation source	Nasopharyngeal/throat swab
latitude and longitude	<a href="#">3.1390 N 101.6869 E</a>

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

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

# 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**

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



## Putting standards into practice: Template and standard terminology

**Home**   **Insert**   **Draw**   **Page Layout**   **Formulas**   **Data**   **Review**   **View**

Paste   Cut   Copy   Format   Arial   10   Bold   Italic   Underline   Merge & Center   General   Conditional Formatting   Format as Table   Cell Styles   Insert   Delete   Format   AutoSum   Fill   Sort & Filter   Find & Select

U6   fx

	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	<b>Sample collection and processing</b>													
	<b>sequence submitted by</b>	<b>sequence submitter contact email</b>	<b>sequence submitter contact address</b>	<b>sample collection date</b>	<b>sample received date</b>	<b>geo_loc name (country)</b>	<b>geo_loc name (state/province/region)</b>	<b>organism</b>	<b>isolate</b>	<b>purpose of sampling</b>	<b>anatomical material</b>	<b>anatomical part</b>	<b>body product</b>	<b>environmental material</b>
2														
3														
4														
5														
6														
7														
8														
9														
10														
11														
12														
13														
14														
15														
16														
17														
18														
19														
20														
21														
22														
23														
24														
25														
26														
27														
28														
29														
30														
31														
32														
33														
34														
35														
36														
37														
38														
39														
40														

Template   Reference Guide   Vocabulary   +

- **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

PHA4GE SARS-CoV-2 Contextual Data Template_demo			
Home Insert Draw Page Layout Formulas Data Review View			
E3			
A	B	C	D
Database Identifiers	Definition	Guidance	Examples
specimen collector sample ID	The user-defined name for the sample.	Every Sample ID from a single submitter must be unique.	prov_rona_99
bioproject umbrella accession	The INSDC umbrella accession number of the BioProject	Required if submission is linked to an umbrella	PRJNA623807
bioproject accession	The INSDC accession number of the BioProject(s) to	Required if submission is linked to a BioProject.	PRJNA12345
biosample accession	The identifier assigned to a BioSample in INSDC arch	Store the accession returned from the BioSample	SAMN14180202
SRA accession	The Sequence Read Archive (SRA), European Nucleo	Store the accession assigned to the submitted "run".	SRR11177792
GenBank/ENA/DBJ accession	The GenBank/ENA/DBJ identifier assigned to the se	Store the accession returned from a GenBank/ENA/DBJ	MN908947.3
GISAID accession	The GISAID accession number assigned to the seque	Store the accession returned from the GISAID	EPI_ISL_123456
GISAID virus name	The user-defined GISAID virus name assigned to the	GISAID virus names should be in the format "hCoV-	hCoV-19/Canada/prov_rona_99/2020
host specimen voucher	Identifier for the physical specimen.	Include a URI (Uniform Resource Identifier) in the form of	URI example:
Sample collection and processing	Definition	Guidance	Examples
sample collected by	The name of the agency that collected the original sar	The name of the agency should be written out in full, (with	Public Health Agency of Canada
sample collector contact email	The email address of the contact responsible for follow	The email address can represent a specific individual or	johnnyblogs@lab.ca
sample collector contact address	The mailing address of the agency submitting the sam	The mailing address should be in the format: Street	655 Lab St, Vancouver, British Columbia,
sequence submitted by	The name of the agency that generated the sequence.	The name of the agency should be written out in full, (with	Centers for Disease Control and Prevention
sequence submitter contact email	The email address of the contact responsible for follow	The email address can represent a specific individual or	Resplab@lab.ca
sequence submitter contact address	The mailing address of the agency submitting the seq	The mailing address should be in the format: Street	123 Sunnybrook St, Toronto, Ontario, M4P
sample collection date	The date on which the sample was collected.	Record the collection date accurately in the template.	2020-03-19
sample received date	The date on which the sample was received.	The date the sample was received by a lab that was not	2020-03-20
geo_loc name (country)	The country of origin of the sample.	Provide the country name from the pick list in the	South Africa
geo_loc name (state/province/territory)	The state/province/territory of origin of the sample.	Provide the state/province/territory name from the GAZ	Western Cape
geo_loc name (county/region)	The county/region of origin of the sample.	Provide the county/region name from the GAZ geography	Derbyshire
geo_loc name (city)	The city of origin of the sample.	Provide the city name from the GAZ geography ontology.	Vancouver
geo_loc latitude	The latitude coordinates of the geographical location o	Provide latitude coordinates if available. Do not use the	38.98 N
geo_loc longitude	The longitude coordinates of the geographical location	Provide longitude coordinates if available. Do not use the	77.11 W
organism	Taxonomic name of the organism.	Select "Severe acute respiratory syndrome coronavirus	Severe acute respiratory syndrome
isolate	Identifier of the specific isolate.	This identifier should be an unique, indexed, alpha-	SARS-CoV-2/human/USA/CA-CDPH-
culture collection	The name of the source collection and unique culture	Format: "<institution-code>:[<collection-	/culture_collection="ATCC:26370"
purpose of sampling	The reason that the sample was collected.	Select a value from the pick list in the template.	Diagnostic testing
purpose of sampling details	Further details pertaining to the reason the sample wa	Provide a free text description of the sampling strategy or	Screening of bat specimens in museum

## PHA4GE – SARS-CoV-2 Contextual Data Template User Guide and SOP 2.0

introduced to capture different kinds of anatomical and environmental samples, as well as collection devices and methods. These fields include "anatomical material", "anatomical part", "body product", "environmental material", "environmental site", "collection device", and "collection method". **Populate only the fields that pertain to your sample.** Provide the most granular information allowable according to your organization's data sharing policies.

**e.g. nasal swab** should be recorded:

host (scientific name)	host (common name)	host disease	anatomical part	collection device
Homo sapiens	Human	COVID-19	Nasopharynx	Swab

**e.g. saliva** should be recorded:

host (scientific name)	host (common name)	host disease	anatomical material
Homo sapiens	Human	COVID-19	Saliva

**e.g. human feces** should be recorded:

host (scientific name)	host (common name)	host disease	body product
Homo sapiens	Human	COVID-19	Feces

**e.g. sewage from treatment plant** should be recorded:

environmental site	environmental material
Sewage Plant	Sewage

**e.g. swab of a hospital bed rail** should be recorded:

environmental site	environmental material	collection device
Hospital	Bed Rail	Swab

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

# Protocols to mobilize harmonized data

The screenshot shows the PHA4GE workspace on Protocols.io. The header includes the workspace name 'PHA4GE' and its description 'The Public Health Alliance for Genomic Epidemiology'. Below this, the interests are listed: 'Public Health, Pathogen Genomics, Bioinformatics, Open Data, Open Source, Interoperability, Reproducibility, Standards, Metadata'. The navigation bar shows 'Publications' with 7 items, 'Members' with 4, 'Discussions' with 1, 'Resources', and 'News'. The main content area displays two publications from the ENA (European Nucleotide Archive) dated Jul 09, 2020. Both publications are titled 'SARS-CoV2 EBI assembly submission protocol' and 'SOP for populating EBI submission templates (ENA)'. They list authors: Nabil-Fareed Alikhan<sup>1</sup>, Emma Griffiths<sup>2</sup>, Ruth Timme<sup>3</sup>, and Duncan MacCannell<sup>4</sup>. The affiliations are: <sup>1</sup>Quadram Institute Bioscience, <sup>2</sup>University of British Columbia, <sup>3</sup>US Food and Drug Administration, and <sup>4</sup>Centers for Disease... Both publications are associated with the 'Coronavirus Method Development Community' and 'PHA4GE' tags. The contact person is Nabil-Fareed Alikhan. The first publication has 49 views and the second has 28 views.

- **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

# Pipelines & Visualisation Working Group



**Victoria Dyster**  
Mitra Bio



**Gültekin Ünal**  
Ankara University



**Emily Smith**  
Theiagen Genomics

## Co-chairs

<https://www.pha4ge.org>

| <https://www.github.com/pha4ge>

|  @pha4ge<sub>12</sub>

22-23 July 2024, Washington DC, USA

# Pipelines & Visualisation Working Group

## Ten Best Practices for Public Health Bioinformatics Pipelines

### PHA4GE Bioinformatics Pipelines & Visualization Working Group

Libuit KG, Guthrie J, Ambrosio F, Kapsak C, Unal Gultekin, Holmes J, Wright S, Nguinkal J, Doughty E, Southgate J, O'Cathail C, Carleton H, Kingwara L, Khan W, Baker K, Diallo A, Connor T, Kanwar S, Maturure P, James S, Cuesta I, Dyster V, Gaskin A, Williams C, Smith E, Rokney A, Petkau A, Varona S, Gnimpieba E, Rey S, Macori G, & Mboowa G

*Updates and modifications to this documented are captured in the repository [changelog](#).*

Public accessible  
Repositories

Semantic versioning of  
stable releases

Workflow management  
systems

Packaged software

<https://pha4ge.org>

|

<https://www.github.com/pha4ge>

|

 @pha4ge



# Public Health Implementation

## Success Stories



Public Health Alliance for  
Genomic Epidemiology

### Nigeria



20 public health scientists complete training on SARS-CoV-2 genome analysis, data sharing and metadata standards

### Malawi



Delivery of the MinION sequencer and reagents for SARS-CoV-2 sample collection and analysis

### Kenya



SARS-CoV-2 sample collection and analysis using updated equipment (sequencers and high performance computing infrastructure)  
Established a platform for genomic training

### Malaysia



Inter-Continental collaboration between Malaysia and researchers in Argentina and Tokyo to implement the PHA4GE hAMRionization metadata specification

### Ethics & Data Sharing Workshop - Mauritius



A 4-day Workshop on ethical data sharing practices and inclusive, equitable health research for teams conducting research into ethics and data governance challenges in under-resourced or challenging research environments

# Thank you



BC Centre for Disease Control



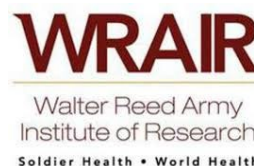
FRED HUTCH™



EMBL-EBI



BILL & MELINDA  
GATES foundation



<https://www.pha4ge.org>

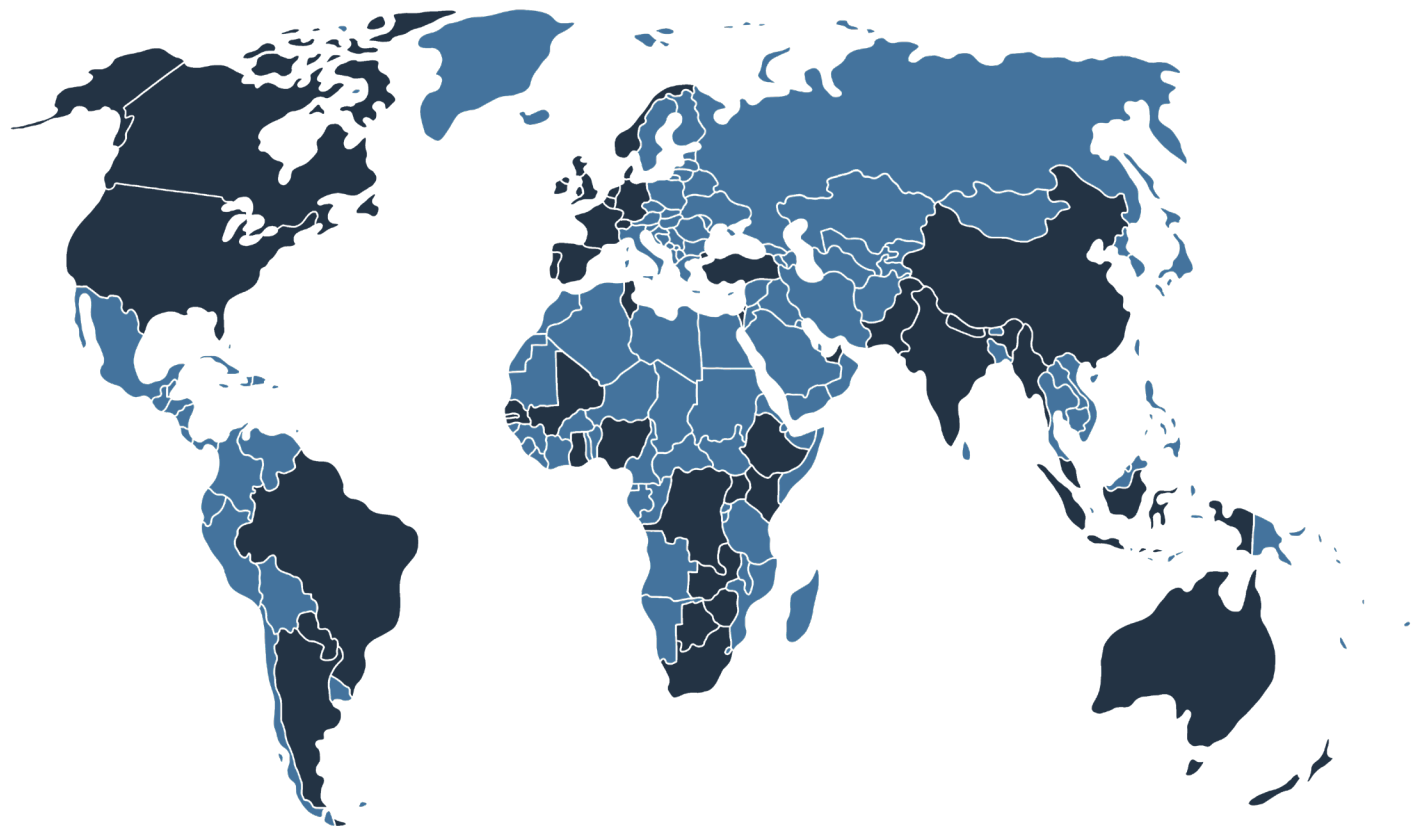
<https://www.github.com/pha4ge>

 @pha4ge<sub>15</sub>

Extra slides for the audience after the participants



# Sustainability



**237 members:**  
63 – Africa  
82 – North America  
61 – Europe  
17 – Asia  
8 – Oceania  
5 – South America

<https://www.pha4ge.org>

| <https://www.github.com/pha4ge>

|  @pha4ge<sub>17</sub>



# Leading Best Practices

## In Progress

### Data structures WG

- Data Object Model Paper
- SARS-COV-2 Primer standardization paper
- AMR hAMRonization tool paper

### Ethics and Data Sharing WG

- Systematic review paper in progress

### Pipelines & Viz WG

- Influenza: Guidance Document in progress covering seasonal and veterinary
- HIV, m.TB, Wastewater Sequencing Guidance Documents in progress.
- Identifying SC2-Recombinants, SC-2 Omicron, MPXV Bioinformatics Solutions
- SC-2 Quality Control Guidance Document

### Infrastructure WG

- Framework for Compute infrastructure for Pathogen Genomics Laboratories in progress



# Leading Best Practices

Published since 2022

- Future-Proofing & maximizing the utility of metadata:...
- <https://academic.oup.com/gigascience/article/doi/10.1093/gigascience/giac003/6529104>
- Framework for the promotion of ethical benefit sharing...
- <https://gh.bmj.com/content/7/2/e008096>
- **Omicron variant guidance document**
- <https://github.com/pha4ge/pipeline-resources/blob/main/docs/omicron-resources.md>
- **Identifying SC2 Recombinants Guidance Doc**
- <https://github.com/pha4ge/pipeline-resources/blob/main/docs/sc2-recombinants.md>
- **SC2 Quality Control Document**
- <https://github.com/pha4ge/pipeline-resources/blob/main/docs/qc-solutions.md>
- **Proposed Standards for Public Health Bioinformatics Software**
- <https://github.com/pha4ge/pipeline-resources/blob/main/docs/pipeline-standards.md>



# Subawards

**Focusing on AMR, SARS-CoV-2 and Ethics & Data Sharing themes**

Botswana

Ethiopia

Kenya

Malawi

Nigeria

South Africa

Sudan

The Gambia

Uganda

Zimbabwe

Cambodia

Congo

DR Congo

Fiji

Malaysia

Nepal

Pakistan

Philippines

Zambia



## **Ethics and Data sharing Working Group**

Nicki Tiffin (UWC, SA)

Framework for sharing data during outbreaks  
<https://ethics-forum.pha4ge.org/>

# ELSI subgrants

## Democratic Republic of Congo & Gabon: (Ethics Training)

Setup, run dedicated online ethics training resources

## Kenya: (Ethics Implementation)

Engage with ethics review boards - review of genomic research applications

## Zimbabwe: (Scientific Citizenship)

Videos made by local researchers targeting school learners using local language

## Nigeria: (Community engagement)

Vaccine Hesitancy in Nigeria