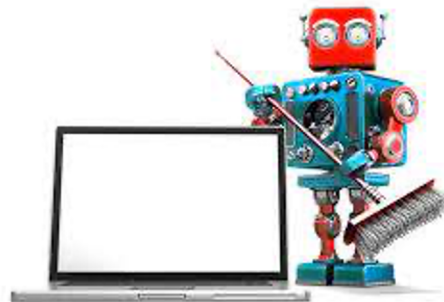


# Overcoming challenges of SARS-CoV-2 genomics data sharing for public health surveillance, outbreak investigations and research using the PHA4GE SARS-CoV-2 contextual data specification

Emma Griffiths, Ruth Timme, Finlay Maguire, Ines Mendes, Lee Katz, Damion Dooley, Rhiannon Cameron, Dominique Anderson, Anders Gonçalves da Silva, William Hsiao, Duncan MacCannell

# Housekeeping

1. Session is being recorded
2. Please keep mics muted until Q&A
3. Please put questions in the chat
4. Please keep cameras off if internet unstable/not presenting
5. Keep phone/apps on silent
6. Slides will be made available after workshop
7. If you'd like to tweet #FAIRConvergence



# Who Are We?



Public Health Alliance for  
Genomic Epidemiology

# Workshop Overview

## 1. Public health microbial genomics

- Importance for COVID-19 response
- Challenges in data harmonization/integration
- Overview of PHA4GE SARS-CoV-2 specification package
- How PHA4GE specification makes genomics contextual data FAIR

## 2. Demo of spec: putting standards into practice

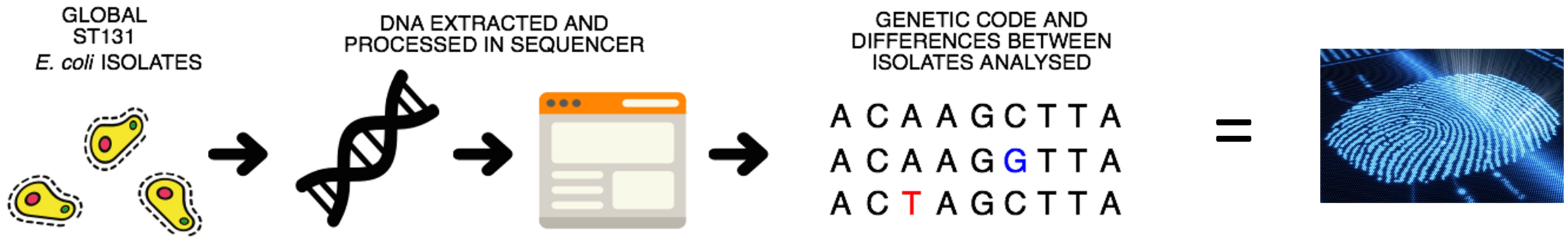
- from chaos to control
- improving the quality of open data

## 3. Implementations of specification

- DataHarmonizer (Canada)
- AusTrakka (Australia)
- Boabab LIMS (South Africa)



# Microbial genome sequences can be used as a molecular fingerprint to trace the source of infectious disease.



- Public health agencies exchange information about these fingerprints



(Dramatic representation from the movie

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**

# Sequencing and sharing of SARS-CoV-2 genomes has had many benefits during the pandemic.

Cite as: X. Deng *et al.*, *Science*  
10.1126/science.abb9263 (2020).

## A SARS-CoV-2 vaccine candidate would likely match all currently circulating variants

 Bethany Dearlove,  Eric Lewitus,  Hongjun Bai,  Yifan Li,  Daniel B. Reeves,  M. Gordon Joyce, Paul T. Scott,  Mihret F. Amare,  Sandhya Vasani,  Nelson L. Michael,  Kayvon Modjarrad, and  Morgane Rolland

PNAS September 22, 2020 117 (38) 23652-23662; first published August 31, 2020;  
<https://doi.org/10.1073/pnas.2008281117>

## The proximal origin of SARS-CoV-2

Kristian G. Andersen , Andrew Rambaut, W. Ian Lipkin, Edward C. Holmes & Robert F. Garry

*Nature Medicine* 26, 450–452(2020) | [Cite this article](#)

5.03m Accesses | 706 Citations | 35003 Altmetric | [Metrics](#)

**To the Editor** – Since the first reports of novel pneumonia (COVID-19) in Wuhan, Hubei province, China<sup>1,2</sup>, there has been considerable discussion on the origin of the causative virus, SARS-CoV-2<sup>3</sup> (also referred to as HCoV-19)<sup>4</sup>. Infections with SARS-CoV-2 are now widespread, and as of 11 March 2020, 121,564 cases have been confirmed in more than 110 countries, with 4,373 deaths<sup>5</sup>.
















SARS-CoV-2 is the seventh coronavirus known to infect humans; SARS-CoV, MERS-CoV and SARS-CoV-2 can cause severe disease, whereas HKU1, NL63, OC43 and 229E are associated with mild symptoms<sup>6</sup>. Here we review what can be deduced about the origin of SARS-CoV-2 from comparative analysis of genomic data. We offer a perspective on the notable features of the SARS-CoV-2 genome and discuss scenarios by which they could have arisen. Our analyses clearly show that SARS-CoV-2 is not a laboratory construct or a purposefully manipulated virus.

## Genomic surveillance reveals multiple introductions of SARS-CoV-2 into Northern California

Xianding Deng<sup>1,2\*</sup>, Wei Gu<sup>1,2\*</sup>, Scot Federman<sup>1,2\*</sup>, Louis du Plessis<sup>3\*</sup>, Oliver G. Pybus<sup>3</sup>, Nuno Faria<sup>3</sup>, Candace Wang<sup>1,2</sup>, Guixia Yu<sup>1,2</sup>, Brian Bushnell<sup>4</sup>, Chao-Yang Pan<sup>5</sup>, Hugo Guevara<sup>5</sup>, Alicia Sotomayor-Gonzalez<sup>1,2</sup>, Kelsey Zorn<sup>6</sup>, Allan Gopez<sup>1</sup>, Venice Servellita<sup>1</sup>, Elaine Hsu<sup>1</sup>, Steve Miller<sup>1</sup>, Trevor Bedford<sup>7,8</sup>, Alexander L. Greninger<sup>7,9</sup>, Pavitra Roychoudhury<sup>7,9</sup>, Lea M. Starita<sup>8,10</sup>, Michael Famulare<sup>11</sup>, Helen Y. Chu<sup>8,12</sup>, Jay Shendure<sup>8,9,13</sup>, Keith R. Jerome<sup>7,9</sup>, Catie Anderson<sup>14</sup>, Karthik Gangavarapu<sup>14</sup>, Mark Zeller<sup>14</sup>, Emily Spencer<sup>14</sup>, Kristian G. Andersen<sup>14</sup>, Duncan MacCannell<sup>15</sup>, Clinton R. Paden<sup>15</sup>, Yan Li<sup>15</sup>, Jing Zhang<sup>15</sup>, Suxiang Tong<sup>15</sup>, Gregory Armstrong<sup>15</sup>, Scott Morrow<sup>16</sup>, Matthew Willis<sup>17</sup>, Bela T. Matyas<sup>18</sup>, Sundari Mase<sup>19</sup>, Olivia Kasirye<sup>20</sup>, Maggie Park<sup>21</sup>, Godfred Masinde<sup>22</sup>, Curtis Chan<sup>22</sup>, Alexander T. Yu<sup>5</sup>, Shua J. Chai<sup>5,15</sup>, Elsa Villarino<sup>23</sup>, Brandon Bonin<sup>23</sup>, Debra A. Wadford<sup>5</sup>, Charles Y. Chiu<sup>1,2,24†</sup>

 [Comment on this paper](#)

## Large scale sequencing of SARS-CoV-2 genomes from one region allows detailed epidemiology and enables local outbreak management

 Andrew J Page, Alison E Mather,  Thanh Le Viet, Emma J Meader,  Nabil-Fareed J Alikhan,  Gemma L Kay,  Leonardo de Oliveira Martins,  Alp Aydin, David J Baker, Alexander J. Trotter, Steven Rudder,  Ana P Tedim, Anastasia Kolyva, Rachael Stanley,  Maria Diaz, Will Potter, Claire Stuart, Lizzie Meadows, Andrew Bell, Ana Victoria Gutierrez,  Nicholas M Thomson,  Evelien M Adriaenssens, Tracey Swingler, Rachel AJ Gilroy, Luke Griffith, Dheeraj K Sethi, Rose K Davidson,  Robert A Kingsley, Luke Bedford, Lindsay J Coupland, Ian G Charles, Ngozi Elumogo,  John Wain, Reenesh Prakash,  Mark A Webber, SJ Louise Smith,  Meera Chand, Samir Dervisevic,  Justin O'Grady, The COVID-19 Genomics UK (COG-UK) consortium

doi: <https://doi.org/10.1101/2020.09.28.20201475>

# 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

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)

Public databases:

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

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

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

<input type="checkbox"/> NPS in UTM	<b>If possible:</b>
<input type="checkbox"/> Throat Swab in UTM	<input type="checkbox"/> BAL
<input type="checkbox"/> Other (Specify):	<input type="checkbox"/> Sputum

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

# The SARS-CoV-2 Contextual Data Specification

## SARS-CoV-2 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

## 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**





# Supporting documentation

pha4ge / SARS-CoV-2-Contextual-Data-Specification

<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

master 1 branch 3 tags Go to file Add file Code

fmaguire Make sure all INSDC are represented 2c69e91 on Jul 15 48 commits

File	Commit Message	Date
PHA4GE Contextual Data SOP.docx	Make sure all INSDC are represented	last month
PHA4GE SARS-CoV-2 Contextual D...	Make sure all INSDC are represented	last month
PHA4GE SARS-CoV-2 EBI assembly...	Add EBI protocols	last month
PHA4GE SARS-CoV-2 EBI submissi...	Add EBI protocols	last month
PHA4GE SARS-CoV-2 GISAID Subm...	Add GISAID submission protocol	last month
PHA4GE SARS-CoV-2 NCBI assemb...	Add NCBI protocols	last month
PHA4GE SARS-CoV-2 NCBI submis...	Add NCBI protocols	last month
PHA4GE SARS-CoV-2 Standardised...	Make sure all INSDC are represented	last month
PHA4GE SOP for populating EBI su...	Add EBI protocols	last month
PHA4GE SOP for populating NCBI s...	Add NCBI protocols	last month
PHA4GE to Sequence Repository Fi...	update filnemaes in readme; remove version from filenames	last month
PHA4GE_SARS-CoV-2_Contextual_...	Make sure all INSDC are represented	last month
README.md	Merge pull request #4 from pha4ge/json_update	last month

- **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>

# Protocols to mobilize harmonized data

The screenshot displays the PHA4GE workspace on Protocols.io. The workspace name is PHA4GE, with the tagline 'The Public Health Alliance for Genomic Epidemiology'. The interests listed are Public Health, Pathogen Genomics, Bioinformatics, Open Data, Open Source, Interoperability, Reproducibility, Standards, and Metadata. The navigation menu includes Administration, New, Upgrade, Workspace Folder (8), Tasks, Export Group Publications, and Contact Admin. The main content area shows a list of publications under the 'Publications' tab, which has 7 items. Two publications are visible, both dated Jul 09, 2020, and both associated with the ENA (European Nucleotide Archive) logo. The first publication is 'SARS-CoV2 EBI assembly submission protocol' by Nabil-Fareed Alikhan, Emma Griffiths, Ruth Timme, and Duncan MacCannell. The second is 'SOP for populating EBI submission templates (ENA)' by the same authors. Both publications are associated with the 'Coronavirus Method Development Community' and 'PHA4GE' tags. The contact for both 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****

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



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

<https://www.preprints.org/manuscript/202008.0220/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<sup>\*</sup>, Ruth E. Timme<sup>ID</sup>, Andrew J. Page, Nabil-Fareed Alikhan, Dan Fornika, Finlay Maguire, Catarina Inês Mendes, Simon H. Tausch<sup>ID</sup>, Allison Black, Thomas R. Connor, Gregory H. Tyson, David M. Aanensen, Brian Alcock, Josefina Campos, Alan Christoffels<sup>ID</sup>, 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

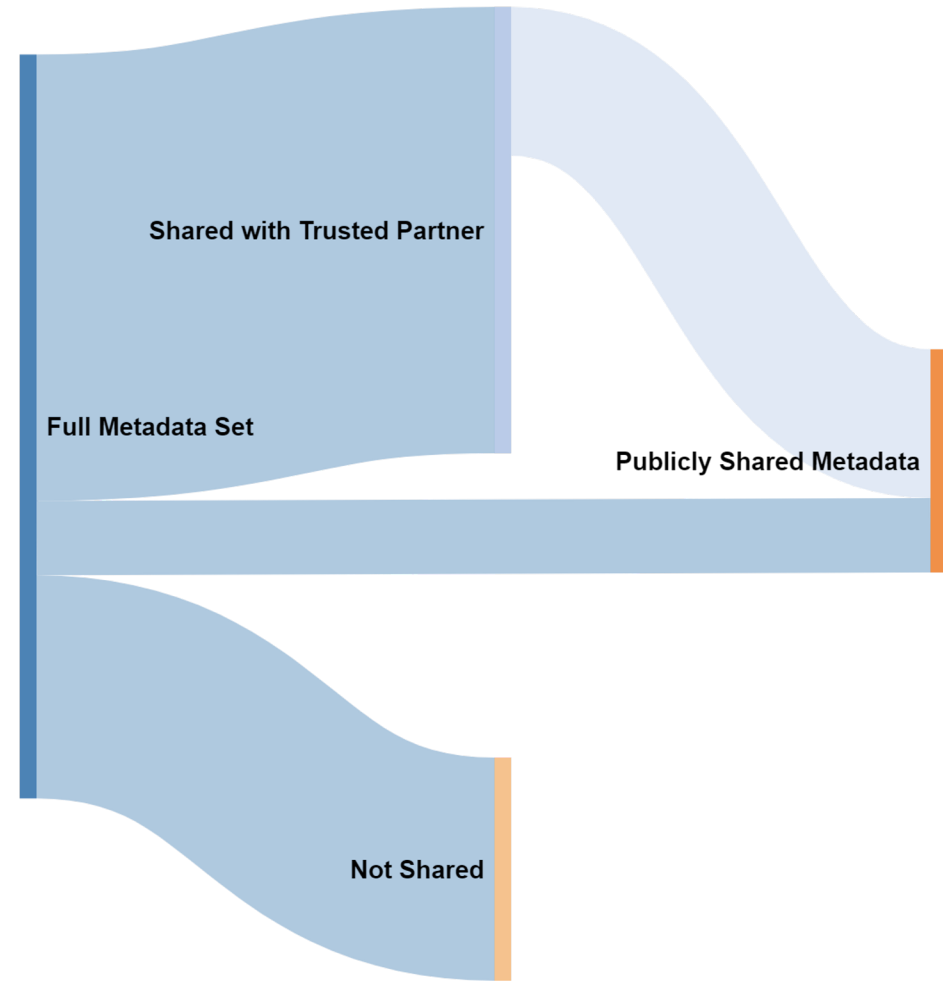
Version 1 : Received: 7 August 2020 / Approved: 9 August 2020 / Online: 9 August 2020 (15:53:58 CEST)

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

The screenshot shows the SoundCloud interface for a podcast episode. The top navigation bar includes 'SOUNDCLOUD', 'Home', 'Stream', 'Library', a search bar, and buttons for 'Sign in', 'Create account', and 'Upload'. The main content area features a play button, the title '26 SARS-CoV-2 contextual data specification for open genomic epidemiology', and a timestamp of '5 days ago'. A '# Science' tag is visible. To the right is a podcast cover for 'Micro Binfie Podcast' featuring a purple bean character with a microphone. At the bottom, there is a waveform player with a progress bar from 0:06 to 37:46.

# How do you use it?

- as much or as little as you want, it's up to you!
- structure metadata **consistently across labs**
- share with **public** repos, **trusted partners**, use for more **efficient private analyses**
- **future-proof** metadata



# How does the PHA4GE Spec make public health genomics contextual data **FAIR**?

**Findable** – *every piece of information has a home, one stop shop*

- data elements standardized, not buried in methods
- ontologies offer URIs (unique, persistent identifiers)

**Accessible** – *understandable by humans/computers*

- spreadsheet and JSON
- protocols for storage in trusted repositories

**Interoperable** – *harmonization across users/standards*

- defines data structures for streamlined communication, data integration

**Reusable** - *enriched datasets*

- genomic information has many uses, enriched contextual data makes data fit for more purposes
- spec usage license (CCBY 4.0)

# Putting standards into practice:

## How to make data FAIR using the PHA4GE spec

### Practical examples

- a) Harmonizing variable contextual data
- b) How to submit harmonized data to

NCBI

### Examples of implementation at organizations

- a) DataHarmonizer (Canada)
- b) Austrakka (Australia)
- c) Baobab LIMS (South Africa)

---Quick Q&A---

- Follow us on twitter
  - @BaobabLIMS
- Online documentation
  - <https://media.readthedocs.org/pdf/baobab-lims/latest/baobab-lims.pdf>
- *Website*
  - [www.baobablims.org](http://www.baobablims.org)
- Get the code (and more)
  - <https://github.com/BaobabLims>
- Send us an email
  - Training – [dominique@sanbi.ac.za](mailto:dominique@sanbi.ac.za)
  - Helpdesk – [help@baobablims.org](mailto:help@baobablims.org)

# Summary

- spec for SARS-CoV-2 **public health contextual data** for **harmonization** across labs and datasets
- future-proof data
- **FAIR**: providing **consistent structure, human/machine-readable**, encourages **data sharing** in responsible way, linking information using **ontologies**
- used by members of **sequencing consortia**
- implemented in **different tools/platforms**



# Thank you!

## 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

# Thank you!



## Welcome to AusTrakka

From genomics to public health decisions for Australia

Combining Genomics & Epidemiological Data	Promoting Data Sharing Across Public Health Labs	Ensuring Better Health Outcomes for Australians
 <p>Human, Environment or Food Sample</p>  <p>Genomic and Epidemiological Data</p>		 <p>More Efficient Public Health Outcomes</p>  <p>Healthier Australia</p>





Thank you for listening and participating!

Get the PHA4GE spec here

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

Get the preprint here

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

Get the DataHarmonizer here

<https://github.com/Public-Health-Bioinformatics/DataHarmonizer/releases/>

Learn about AusTrakka

<https://portal.austrakka.net.au/>

Learn about Baobab LIMS

<https://github.com/BaobabLims>

Contact us

[datastructures@pha4ge.org](mailto:datastructures@pha4ge.org)