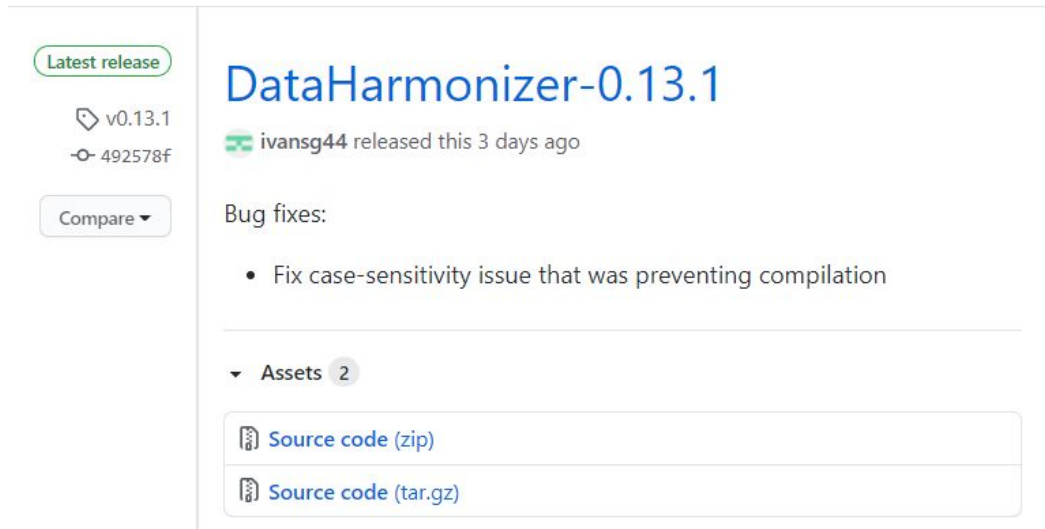


Contextual Data (Metadata) Curation

- I. **Purpose:** To harmonize SARS-CoV-2 contextual data across data providers in the CanCOGeN network.
- Data providers will extract and curate lab-specific contextual data according to the steps outlined in the procedure below.
 - Laboratories will populate the harmonized template with information from their datasets using the *DataHarmonizer* application.
 - Data providers will share the harmonized data with the national database according to the agreed upon mechanism.
- II. **Data:** The contextual data describing sample collection and processing, host information, sequencing, and bioinformatics and QC metrics as supplied by the data provider.
- III. **Procedure:**

	Action	Related docs
1	<p>Download the zip file ("Source code (zip)") containing The DataHarmonizer application from the following link: https://github.com/Public-Health-Bioinformatics/covid19ValidationGrid/releases</p>  <p>Extract the zip file's contents, and navigate into the extracted folder. Open main.html. The validator application will open in your default browser. It should look like this:</p>	

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

	Action	Related docs
	<p>Data can be entered into the validator application manually, by typing values into the application's spreadsheet, or data can be imported from local xlsx, xls, tsv and csv files.</p> <p><i>Note: Only files containing the headers expected by the DataHarmonizer can be opened in the application.</i></p> <p>To import local data, click File on the top-left toolbar, and then click Open. To enter data in a new file, click File on the top-left toolbar, and then click New. Data entered into the spreadsheet can be copied and pasted.</p>	
2	<p>Before you begin to curate sample metadata:</p> <ul style="list-style-type: none"> • Review your dataset • Review the fields in the template of the Validator application • Review the field descriptions in the SOP Appendix 	
3	<p>Familiarize yourself with DataHarmonizer functionality by reviewing the “Getting Started”. To access "Getting Started", click on the green Help button on the top-left toolbar, then click Getting Started. Definitions, examples and further guidance are available by double clicking on the field headers, or by using the “Reference Guide”. To access the “Reference Guide” click on the Help button, then click Reference Guide.</p>	



CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

	Action	Related docs
4	<p>Confirm mapping of your data fields to those in the harmonized template with the data steward (e.g. your supervisor).</p> <p><i>Note: A version of this information will be made public in GISAID and NCBI, however, another version of this data will be captured in the access controlled national database. Confirm the level of granularity of information that can be shared publicly vs in the national database, with the data steward and/or the privacy officer. The most detailed information allowable should be included here.</i></p>	
5	<p>Enter data into the validator spreadsheet.</p> <ul style="list-style-type: none"> • Hide non-required fields (colour-coded purple and white) by clicking Settings on the top-left toolbar, followed by clicking on Show Required Columns (colour-coded in yellow). • Double click in the field headers to see definitions and detailed guidance as needed (or consult Appendix A). • Jump to a specific field header by clicking Settings on the top-left toolbar, followed by clicking on Jump to, then select the field header of the column you would like to view from the drop down list. • Populate the validator template with the information from your dataset. • Use picklists when provided. • A value must be entered for every required field in each row. If data is missing or not collected, choose a null value from the picklist. <ul style="list-style-type: none"> ○ Not Applicable ○ Missing ○ Not Collected ○ Not Provided ○ Restricted Access • Free text can be provided when picklists are not available. <p>If a desired term is not present in a picklist, contact emma.griffiths@bccdc.ca.</p> <p><i>Note: Sometimes there will be constraints on what information can be shared, other times a field may not be applicable to your sample. Use the null values (controlled vocabulary indicating the reason why information is not provided) in the picklist to report missing data.</i></p> <p>Required fields are organized into subsections (see Appendix A for required field definitions and guidance, and Appendix B for examples of how to structure sample descriptions):</p>	

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

	Action	Related docs												
	<table><tr><th>Subsection</th><th>Required Fields</th></tr><tr><td>Sample Collection and Processing <i>Note: Evaluate with your supervisor whether the specimen collector sample ID is considered identifiable by your institutional policies. If not, copy the sample ID into the sample ID field in the validator spreadsheet. If yes, provide the alternative sample ID as specified by the lab. Be sure to keep a copy of the key.</i></td><td>specimen collector sample ID sample collected by sequence submitted by sample collection date geo_loc (country) geo_loc (province/territory) organism isolate</td></tr><tr><td>Describing the material and/or site sampled. <i>Note: Seven fields have been introduced to capture different kinds of anatomical and environmental samples, as well as collection methods. Populate only the fields that pertain to your sample - provide null values for the fields that are not applicable. Provide the most granular information allowable according to your organization's data sharing policies.</i></td><td>anatomical material anatomical part body product environmental material environmental site collection device collection_method</td></tr><tr><td>Host Information</td><td>host (scientific name) host disease host age host gender</td></tr><tr><td>Sequencing</td><td>sequencing instrument</td></tr><tr><td>Bioinformatics and QC Metrics</td><td>consensus sequence method</td></tr></table>	Subsection	Required Fields	Sample Collection and Processing <i>Note: Evaluate with your supervisor whether the specimen collector sample ID is considered identifiable by your institutional policies. If not, copy the sample ID into the sample ID field in the validator spreadsheet. If yes, provide the alternative sample ID as specified by the lab. Be sure to keep a copy of the key.</i>	specimen collector sample ID sample collected by sequence submitted by sample collection date geo_loc (country) geo_loc (province/territory) organism isolate	Describing the material and/or site sampled. <i>Note: Seven fields have been introduced to capture different kinds of anatomical and environmental samples, as well as collection methods. Populate only the fields that pertain to your sample - provide null values for the fields that are not applicable. Provide the most granular information allowable according to your organization's data sharing policies.</i>	anatomical material anatomical part body product environmental material environmental site collection device collection_method	Host Information	host (scientific name) host disease host age host gender	Sequencing	sequencing instrument	Bioinformatics and QC Metrics	consensus sequence method	
	Subsection	Required Fields												
	Sample Collection and Processing <i>Note: Evaluate with your supervisor whether the specimen collector sample ID is considered identifiable by your institutional policies. If not, copy the sample ID into the sample ID field in the validator spreadsheet. If yes, provide the alternative sample ID as specified by the lab. Be sure to keep a copy of the key.</i>	specimen collector sample ID sample collected by sequence submitted by sample collection date geo_loc (country) geo_loc (province/territory) organism isolate												
	Describing the material and/or site sampled. <i>Note: Seven fields have been introduced to capture different kinds of anatomical and environmental samples, as well as collection methods. Populate only the fields that pertain to your sample - provide null values for the fields that are not applicable. Provide the most granular information allowable according to your organization's data sharing policies.</i>	anatomical material anatomical part body product environmental material environmental site collection device collection_method												
	Host Information	host (scientific name) host disease host age host gender												
	Sequencing	sequencing instrument												
Bioinformatics and QC Metrics	consensus sequence method													
6	Validate the entered data by clicking on the Validate button on the top-left toolbar.													

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

	Action	Related docs
	<p>Missing information and invalid entries in required fields will be highlighted in red.</p> <ul style="list-style-type: none"> • Observe invalid rows by clicking Settings in the top-left toolbar, and then clicking on Show invalid rows. • Observe valid rows by clicking Settings in the top-left toolbar, and then clicking on Show valid rows. • Return view to all rows by clicking Settings in the top-left toolbar, and then clicking on Show all rows. <p><i>Note: Row viewing options only appear after a validation attempt has been made.</i></p>	
7	<p>Address any invalid data that was flagged in red in the template.</p> <ul style="list-style-type: none"> •  Pale Red = Incorrect data format •  Dark Red = Required data missing 	
8	<p>Export validated data by clicking File on the top-left toolbar, and then clicking on Save as. Enter the file name and press Save. Export to IRIDA, GISAID, or CNPHI formats by clicking File on the top-left toolbar, and then clicking Export to.</p> <ul style="list-style-type: none"> • Have the validated data reviewed by the data steward (i.e. your supervisor) 	
9	<p>Submit validated data to the national database.</p> <p>You can submit either by i) emailing the validated data to your NML contact, or ii) uploading the validated data directly to through the CNPHI interface via the CNPHI Metadata Uploader.</p> <ul style="list-style-type: none"> • Before uploading to CNPHI, export your data in “CNPHI” format by clicking File on the top-left toolbar, then clicking Export To. Type in the file name, and select “CNPHI” from the Format picklist. Then click Export. • See CNPHI documentation for more information regarding Metadata Upload. 	
10	<p>Optional: Format validated data for GISAID submission.</p> <p>The DataHarmonizer will automate the preparation of a GISAID submission form from the entered data by exporting the data in GISAID format.</p> <ul style="list-style-type: none"> • Export your data in “GISAID” format by clicking File on the top-left toolbar, then clicking Export To. Type in the file name, and select “GISAID” from the Format picklist. Then click Export. 	

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

	Action	Related docs
11	<p>Optional: Format validated data for IRIDA submission.</p> <p>The DataHarmonizer will automate the preparation of an IRIDA submission form from the entered data by exporting the data in IRIDA format.</p> <ul style="list-style-type: none"> Export your data in “IRIDA” format by clicking File on the top-left toolbar, then clicking Export To. Type in the file name, and select “IRIDA” from the Format picklist. Then click Export. <p><i>Note: If the top row containing the broad headings (Database identifiers, Sample collection and processing, Host information, Sequencing, Bioinformatics and QC, Authors) is not removed, the IRIDA metadata upload will fail. By exporting the sheet in IRIDA format, the DataHarmonizer completes this formatting for you.</i></p>	<p>Adding a new sample to IRIDA:</p> <p>https://irida.cefacility.ca/documentation/user/samples/#adding-a-new-sample</p>
12	<p>Additional Information:</p> <p>A local copy of the Standard Operating Procedure (SOP) is included in every download of the DataHarmonizer. To access it, click on the green Help button on the top-left toolbar, then click SOP.</p> <p>The latest version of the SOP is published online and accessible via a web browser at all times.</p>	<p>Latest SOP for DataHarmonizer:</p> <p>https://docs.google.com/document/d/e/2PACX-1vR4UkqrLaj1-9jxmrNk9mZ4S4Siim8onPHqgdXKd9m1lOr oXmekClfPsXlqgFDio1rWZW7lHArSAbOg/pub</p>

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

IV. Appendix A: Required Field Definitions and Guidance

Field definitions for required fields, as well as guidance and examples, are provided below. This information has been sourced from the DataHarmonizer reference guide. Guidance for strongly recommended and optional fields can be found in the reference guide. For access to information on non-required fields, refer to “Procedure - Action 3”.

Sample Collection and Processing

specimen collector sample ID

The user-defined name for the sample.

Store the collector sample ID. If this number is considered identifiable information, provide an alternative ID. Make sure to store the key between this alternative ID and the original ID for traceability. Every collector sample ID from a single submitter must be unique. It can have any format, but we suggest that you make it concise, unique and consistent within your lab.

e.g. prov_rona_99

sample collected by

The name of the agency that collected the original sample.

The name of the sample collector should be written out in full, (with minor exceptions) and be consistent across multiple submissions e.g. Public Health Agency of Canada, Public Health Ontario, BC Centre for Disease Control. The sample collector specified is at the discretion of the data provider (i.e. may be hospital, provincial public health lab, or other).

e.g. BC Centre for Disease Control

sequence submitted by

The name of the agency that generated the sequence.

The name of the agency should be written out in full, (with minor exceptions) and be consistent across multiple submissions e.g. Public Health Agency of Canada, Public Health Ontario, BC Centre for Disease Control.

e.g. Public Health Ontario

sample collection date

The date on which the sample was collected.

Sample collection date is critical for surveillance and many types of analyses. Required granularity includes year, month and day. Record the collection date accurately in the template. Before sharing this data, ensure you have consulted the data steward and/or your privacy officer regarding whether they consider this date to be identifiable information. If this date is considered identifiable, it is acceptable to add "jitter" to the collection date you share by adding or subtracting a calendar day (acceptable by GISAID). Do not change the collection date in your original records. Alternatively, "received date" may be used as a substitute in the data you share. The date should be provided in ISO 8601 standard format "YYYY-MM-DD".

e.g. 2020-03-16

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

geo_loc_name (country)

The country where the sample was collected.

Provide the country name from the controlled vocabulary provided.

e.g. Canada

geo_loc_name (province/territory)

The province/territory where the sample was collected.

Provide the province/territory name from the controlled vocabulary provided.

e.g. Saskatchewan

organism

Taxonomic name of the organism.

Use Severe acute respiratory syndrome coronavirus 2. This value is provided in the template.

e.g. Severe acute respiratory coronavirus 2

isolate

Identifier of the specific isolate.

Provide the isolate name. This identifier should be an unique, indexed, alpha-numeric ID within your laboratory. The isolate name is often the same as the specimen collector sample ID.

Suggested: Isolate name should be identical to the GISAID virus name, which should be written in the format "hCov-19/CANADA/xxxxx/2020".

e.g. hCov-19/CANADA/prov_rona_99/2020

Describing the material and/or site sampled.

anatomical material

A substance obtained from an anatomical part of an organism e.g. tissue, blood.

Provide a descriptor if an anatomical material was sampled. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Blood

anatomical part

An anatomical part/location of an organism e.g. oropharynx.

Provide a descriptor if an anatomical part was sampled. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Nasopharynx (NP)

body product

A substance excreted/secreted from an organism e.g. feces, urine, sweat.

Provide a descriptor if a body product was sampled. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Feces

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

environmental material

A substance or object obtained from the natural or man-made environment e.g. soil, water, sewage.

Provide a descriptor if an environmental material was sampled. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Face Mask

environmental site

An environmental location may describe a site in the natural or built environment e.g. contact surface, metal can, hospital, wet market, bat cave.

Provide a descriptor if an environmental site was sampled. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Building floor

collection device

The instrument or container used to collect the sample e.g. swab.

Provide a descriptor if a device was used for sampling. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Swab

collection_method

The process used to collect the sample e.g. phlebotomy, necropsy.

Provide a descriptor if a collection method was used for sampling. Use the picklist provided in the template. If a desired term is missing from the picklist, contact emma.griffiths@bccdc.ca. If not applicable, do not leave blank. Choose a null value.

e.g. Bronchoalveolar Lavage (BAL)

Host Information

host (scientific name)

The taxonomic, or scientific name of the host.

Common name or scientific name are required if there was a host. Both can be provided, if known. Use terms from the pick lists in the template. Scientific name e.g. *Homo sapiens*. If the sample was environmental, put "not applicable".

e.g. *Homo sapiens*

host disease

The name of the disease experienced by the host.

This field is only required if there was a host. If the host was a human select COVID-19 from the pick list. If the host was asymptomatic, this can be recorded under "host health state details". If the host is not human, and the disease state is not known or the host appears healthy, put "not applicable".

e.g. COVID-19

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

host age

Age of host at the time of sampling.

Enter the age of the host in years. If not available, provide a null value. If there is not host, put "Not Applicable".

e.g. 79

host age bin

Age of host at the time of sampling, expressed as an age group.

Select the corresponding host age bin from the pick list provided in the template. If not available, provide a null value. The "host age bin" field will automatically propagate with the bin that corresponds to the input in "host age". If not available or you are not permitted to share, put a null value

Age Bins:

0 - 9

10 - 19

20 - 29

30 - 39

40 - 49

50 - 59

60 - 69

70 - 79

80 - 89

90 - 99

100+

host gender

The gender of the host at the time of sample collection.

Select the corresponding host gender from the pick list provided in the template. If not available, choose a null value. If there is no host, put "Not Applicable".

e.g. Male

Sequencing

sequencing instrument

The model of the sequencing instrument used.

Select a sequencing instrument from the picklist provided in the template.

e.g. Minlon

Bioinformatics and QC Metrics

consensus sequence method

The name and version number of the protocol used to produce the consensus sequence.

Provide the software name followed by the version.

e.g. iVar 1.2

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

V. **Appendix B: Structuring Sample Descriptions (Examples)**

Several examples are provided below which illustrate how to structure common sample descriptions.

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 (NP)	Swab

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

host (scientific name)	host (common name)	host disease	anatomical part	collection device
Homo sapiens	Human	COVID-19	Oropharynx (OP)	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. *swab of a hospital bed rail* should be recorded:

environmental site	environmental material	collection device
Hospital	Bed Rail	Swab

e.g. *tissue from a bat (Chiroptera) in a cave* should be recorded:

Host (common name)	Host (scientific name)	host disease	anatomical_part	environmental_site
Bat	Chiroptera	Not applicable	Tissue	Cave

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

e.g. *particulates from air filter* should be recorded:

environmental material	collection method
Particulate Matter	Air Filtration

VI. **Appendix C: Null Value Definitions**

Not Applicable

Information is inappropriate to report, can indicate that the standard itself fails to model or represent the information appropriately.

Missing

Information was known to be recorded in the past, but the observed value cannot be located or retrieved for some reason.

Not Collected

Information of an expected format was not given because it has not been collected.

Not Provided

Information of an expected format was not given, a value may be given at the later stage.

Restricted Access

Information exists but can not be released openly because of privacy concerns.

Source:

International Nucleotide Sequence Database Collaboration (INSDC) Missing Value Reporting Terms (2017-2018). *ENA Training Modules*:

<https://ena-docs.readthedocs.io/en/latest/submit/samples/missing-values.html>

CanCOGeN – SARS-CoV-2
CanCOGeN_1.3 Contextual Data Curation

Revision History

Version	Date	Writer	Description of Change
0.0	May 25, 2020	Lauren Tindale, Emma Griffiths	Created protocol
1.0	June 8, 2020	Emma Griffiths	Protocol edited
1.1	June 16, 2020	Emma Griffiths	Protocol edited
1.2	October 05, 2020	Rhiannon Cameron	Protocol edited
1.3	October 06, 2020	Rhiannon Cameron	Protocol edited