

DPCC Data Standard Reference for DIGS Sequencing Results-Consensus Sequence v1.0

	Project_Identifier	Sequencing_Study_Identifier	Extract_Identifier	Sample_Identifier	Virus_Identifier
Input Type	Text Field	Text Field	Text Field	Text Field	Text Field
Definition	A unique Project Identifier generated by the DPCC by combining the Center-generated Project Code and a random 4-digit number	The unique code associated with the sequencing study.	A unique identifier assigned to the extract by the DPCC.	Identifier initially assigned to each sample collected. If multiple samples are taken from the same host, each sample should have its own identifier.	A unique laboratory identifier assigned to the virus by the collector or creator
Format	Project_Code_XXXX Maximum length: 21 characters	Text Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _ Maximum length: 50 characters	Text Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _ Maximum length: 60 characters	Center-specific Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _ Maximum length: 50 characters	Text Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _ Maximum length: 50 characters
Value List	DIGS-JCVI_COLLAB DIGS-JHU_JH DIGS-MSSM_CRIP DIGS-UGA_CRIP DIGS-UGA_Emory-UGA DIGS-UGA_SJ DIGS-UT_SJ	Text NA	Text NA	None	Text NA
Curation	The entry must be a DIGS Project Identifier value registered with the DPCC.	None	Values must match values previously submitted through the SWT Sequence Request Interface.	The Sample_Identifier initially assigned to the surveillance sample must be provided.	Values must match values previously submitted through the SWT Sequence Request Interface.
Examples	DIGS-UGA_SJ	SQN_Proj_01	R16_S2	22258468	SWN-9816
Notes	Use the identifier for the DIGS facility performing the sequencing. That project may be different than the project under which the original sample was collected or the virus isolated.	If you require multiple samples to be sequenced by the same sequencing core, you may provide a Sequencing_Study_Identifier. All future submissions linked to the same Sequencing_Study_Identifier will automatically be sent to the same sequencing core. If you do not have this requirement, enter NA.		Submissions with a corresponding surveillance submission may import information already held in the DPCC database. Submissions without a corresponding surveillance submission must complete all fields in the metadata template. For non-surveillance derived submissions, enter a unique Sample_Identifier.	This field provides an additional layer of tracking if multiple viruses were isolated from the same surveillance sample or if the virus is a laboratory-created virus.
Dependent Fields					
Validation	Project_Identifier should be a valid project identifier.	Validate field length	Validate field length	Validate field length	Validate field length
Message Code	Error_9_PROJECT_NOT_FOUND	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH
Validation	Project_Identifier exists but user does not have permission to access or edit the project.				
Message Code	Error_4_DENIED_USER_ACCESS				
Validation					
Message Code					
Validation					
Message Code					
Validation					
Message Code					
Validation					
Message Code					
Validation					
Message Code					
Validation					
Message Code					

	Strain_Name
Input Type	Text Field
Definition	The WHO or ICTV strain name of the virus that was sequenced
Format	<p>Influenza A virus: Antigenic Type/Host of Origin/Geographical Origin/Strain Number/Year of Isolation (Subtype)</p> <p>Influenza B, C, or D virus: Antigenic Type/Host of Origin/Geographical Origin/Strain Number/Year of Isolation</p> <p>SARS-CoV-2 and other viruses: Virus Name/Host of Origin/Geographical Origin/Strain Number/Year of Isolation</p> <p>Maximum length: 150 characters</p>
Value List	Text U
Curation	<p>For influenza viruses, strain names most follow WHO naming convention: Fields must be ordered as follows and separated with the '/' character:</p> <ol style="list-style-type: none"> 1. The antigenic type (e.g., A, B, C, D) 2. The host of origin (e.g., swine, equine, chicken. For human-origin viruses, no host of origin designation is given.) 3. Geographical origin (e.g., Denver, Taiwan) 4. Strain number (e.g., 15, 7) 5. Year of isolation (e.g., 2009, 1934) 6. For influenza A viruses, the hemagglutinin and neuraminidase antigen description in parentheses (e.g., (H1N1), (H3N2)) <p>For SARS-CoV-2 viruses, strain names most follow ICTV naming convention: Fields must be ordered as follows and separated with the '/' character:</p> <ol style="list-style-type: none"> 1. Virus name (e.g., SARS-CoV-2) 2. The host of origin (e.g., human. Human-origin viruses must include the origin designation.) 3. Country of geographical origin as a three-letter code from the DPCC's Country Codes list (e.g., USA, MEX, CAN) 4. Strain number (e.g., 15, 7) 5. Year of isolation (e.g., 2019, 2020) <p>For all other viruses, please use the following convention: Fields must be ordered as follows and separated with the '/' character:</p> <ol style="list-style-type: none"> 1. Virus name (e.g., MERS-CoV, Bat-CoV, etc.) 2. The host of origin (e.g., human, bat, camel. Human-origin viruses must include the origin designation.) 3. Geographical origin, either regional locality or country (e.g., Denver, Taiwan) 4. Strain number (e.g., 15, 7) 5. Year of isolation (e.g., 2019, 2020)
Examples	<p>For influenza:</p> <p>A/Hong Kong/1/1968 (H3N2), A/chicken/Fujian/4/2002 (H3N6), A/chicken/Fujian/4/2002 (HxNx), A/chicken/Fujian/4/2002 (mixed), A/swine/Iowa/233-56/2011 (H3N2), A/duck/Alberta/35/1976 (H1N1), B/Hong Kong/432/2014, C/Texas/19876/2011, or D/swine/Oklahoma/1334/2011</p> <p>For SARS-CoV-2:</p> <p>SARS-CoV-2/human/USA/NY-PV08486/2020</p> <p>For other viruses:</p> <p>PHEV-CoV/swine/USA/15TOSU25049/2015</p>
Notes	<p>WHO Reference for influenza viruses: http://www.cdc.gov/flu/about/viruses/types.htm ICTV Reference for SARS-CoV-2 viruses: https://pubmed.ncbi.nlm.nih.gov/32123347</p> <p>(HxNx) can be used in cases where a partial subtype has been determined (e.g., H5Nx, HxN2).</p> <p>If there are mixed subtypes contained within a sample use A/chicken/Fujian/4/2002 (mixed) for Strain_Name and enter additional subtype information in the Comments field.</p>
Dependent Fields	
Validation	Validate field length
Message Code	Error_70_INVALID_FIELD_LENGTH
Validation	Geographical origin element of strain name must be alphanumeric, underscore, dash, period, or single quote: a-z, A-Z, 0-9, _, -, ., ''
Message Code	Error_138_INVALID_LOCATION_STRAIN_TEXT
Validation	Strain number element of strain name must be alphanumeric, underscore, dash, or period: a-z, A-Z, 0-9, _, -, .
Message Code	Error_139_INVALID_STRAIN_NUMBER_STRAIN_TEXT
Validation	Year of isolation element of strain name must be 4-digit year.
Message Code	Error_110_INVALID_STRAIN_YEAR
Validation	If strain is influenza A, subtype must be present as text between parentheses.
Message Code	Error_112_INVALID_STRAIN_SUBTYPE
Validation	Brackets cannot be present if strain does not have rg- prefix.
Message Code	Error_119_INVALID_STRAIN_BRACKETS
Validation	If included, the strain must have a matched pair of opening and closing brackets or parentheses.
Message Code	Error_158_INVALID_STRAIN_OPEN_BRACKETS
Validation	Brackets cannot be present if strain does not have rg- prefix.
Message Code	Error_159_INVALID_NUMBER_OF_STRAIN_ELEMENTS
Validation	If included, the strain must have a matched pair of opening and closing brackets or parentheses.
Message Code	Error_158_INVALID_STRAIN_OPEN_BRACKETS
Validation	Brackets cannot be present if strain does not have rg- prefix.
Message Code	Error_159_INVALID_NUMBER_OF_STRAIN_ELEMENTS

	Sequencing_Technology	Forward_Primer	Reverse_Primer	Assembly_Method
Input Type	Text Field	Text Field	Text Field	Text Field
Definition	The name of the sequencing technology used to obtain the submitted sequences	The forward PCR primer that was used to amplify the nucleic acid that was sequenced	The reverse PCR primer that was used to amplify the nucleic acid that was sequenced	The name of the program used to assemble next-generation or sanger sequencing reads
Format	Text Maximum length: 250 characters	Name:Sequence Maximum length: 500 characters	Name:Sequence Maximum length: 500 characters	Text Maximum length: 150 characters
Value List	Illumina Ion Torrent Oxford Nanopore PacBio	Text U	Text U	None
Curation	The entry must be one or more comma-separated members of the Value List.	The entry must include the forward primer name and nucleotide sequence separated by a colon.	The entry must include the reverse primer name and nucleotide sequence separated by a colon.	The entry must be the name of a valid sequence assembly program.
Examples	Ion Torrent	fwd_seq:catgttgcacaaggga, or U	rev_seq:atgggatgcagattgtgga, or U	IonTorrent Assembler, or BioEdit
Notes	If more than one sequencing technology is used, comma-separate individual technologies.	If multiple forward primers were used, comma-separate individual forward primers. Enter U if the forward primer is unknown.	If multiple reverse primers were used, comma-separate individual reverse primers. Enter U if the reverse primer is unknown.	Sequences must be pre-assembled. Raw sequence reads from next generation sequencing technologies should not be submitted to GenBank. If more than one assembly method is used, comma-separate individual methods.
Dependent Fields				
Validation	Field value should be one of valid values as in list.. NOTE: User can enter other value by prefixing 'OTH.'	Validate primer name and sequence format of primename:primer sequence. Multiple primers will be separated by comma.	Validate primer name and sequence format of primename:primer sequence. Multiple primers will be separated by comma.	Validate field length
Message Code	Error_1_INVALID_VALUE	Error_83_INVALID_PRIMER_NAME	Error_83_INVALID_PRIMER_NAME	Error_70_INVALID_FIELD_LENGTH
Validation	Field length including values from the Value List and free text following 'OTH-' must be less than 250 characters.	Validate field length	Validate field length	
Message Code	Error_75_INVALID_FIELD_LENGTH_OTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				

	Assembler_Version	Coverage	File_Name	Comments
Input Type	Text Field	Text Field	Text Field	Text Field
Definition	The version of the assembly program used or, if not available, the date the assemblies were made	The average number of reads representing a given nucleotide in the sequence	The name of the sequence data file.	Text describing anything else of interest related to the submission
Format	Text DD-Mon-YYYY DD-Mon-YY D-Mon-YYYY D-Mon-YY Maximum length: 50 characters	Text Maximum length: 50 characters	Text Maximum length: 2000 characters	Text Maximum length: 2000 characters
Value List	Text Date	Number U	Text	Text NA
Curation	The entry must be the version number of the assembly program used in format v.x.x or the date the assemblies were created.	The entry must be a number or enter U if unknown.	The entry must be the full file name, with extension.	None
Examples	v.3.2, 03-Mar-2011, 3-Mar-2011, or 3-Mar-11	25.47	22258468_1.fasta	NA
Notes	If more than one assembly method is used, comma-separate individual versions.	If more than one coverage is used, comma-separate individual coverages.	There must be one sequence data file per extract. All sequence data files and the Sequencing Results Metadata file must be submitted in one zip file.	If there are no comments, enter NA.
Dependent Fields				
Validation	Validate field length	Validate field length	Validate field length	Validate field length
Message Code	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH
Validation	Validate assembler version format or date format.	Coverage has an invalid value. Must be a number or U	File not found in ZIP file	
Message Code	Error_64_VALIDATION_ASSEMBLER_VERSION	Error_65_VALIDATION_COVERAGE	Error_87_DATA_FILE_MISSING_FROM_ZIP	
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				
Validation				
Message Code				