

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

	Project Identifier	Sequencing Study Identifier	Extract Identifier	Sample Identifier
<b>Input Type</b>	Text Field	Text Field	Text Field	Text Field
<b>Definition</b>	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.
<b>Format</b>	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
<b>Value List</b>	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
<b>Curation</b>	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.
<b>Examples</b>	DIGS-UGA_SJ	SQN_Proj_01	R16_S2	22258468
<b>Notes</b>	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.
<b>Dependent Fields</b>				
<b>Validation</b>	Project_Identifier should be a valid project identifier.	Validate field length	Validate field length	Validate field length
<b>Message Code</b>	Error_9_PROJECT_NOT_FOUND	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH
<b>Validation</b>	Project_Identifier exists but user does not have permission to access or edit the project.			
<b>Message Code</b>	Error_4_DENIED_USER_ACCESS			
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				

	Virus Identifier	Strain Name	Sequencing Technology
<b>Input Type</b>	Text Field	Text Field	Text Field
<b>Definition</b>	A unique laboratory identifier assigned to the virus by the collector or creator	The WHO strain name of the virus that was sequenced	The name of the sequencing technology used to obtain the submitted sequences
<b>Format</b>	Text  Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _  Maximum length: 50 characters	Antigenic Type/Host of Origin/Geographical Origin/Strain Number/Year of Isolation (Subtype)  rg-Antigenic Type/Host of Origin/Geographical Origin/Strain Number/Year of Isolation [rg details] (Subtype)  Maximum length: 150 characters	Text  Maximum length: 250 characters
<b>Value List</b>	Text NA	Text U	Illumina Ion Torrent Oxford Nanopore PacBio
<b>Curation</b>	Values must match values previously submitted through the SWT Sequence Request Interface.	WHO strain naming convention: Fields must be ordered as follows and separated with the '/' character: 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))  If the strain is a reverse genetic virus, prefix Strain_Name with rg-. Provide reverse genetic details in brackets between year of isolation and subtype.	The entry must be one or more comma-separated members of the Value List.
<b>Examples</b>	SWN-9816	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 (H3N2v), A/duck/Alberta/35/1976 (H1N1), B/Hong Kong/432/2014, C/Texas/19876/2011, rg-A/Puerto Rico/8/1934 (H1N1), or rg-A/Egypt/N03072/2010 [PR8 internal R 6+2] (H5N1)	Ion Torrent
<b>Notes</b>	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.	Reference: <a href="http://www.cdc.gov/flu/about/viruses/types.htm">http://www.cdc.gov/flu/about/viruses/types.htm</a>  (HxNx) can be used in cases where a partial subtype has been determined (e.g., H5Nx).  If there are mixed subtypes contained within a sample use A/chicken/Fujian/4/2002 (mixed) for Strain_Name or list multiple comma-separated Strain_Names: A/mallard/Alaska/2/2007 (H3N8),A/mallard/Alaska/2/2007 (H3N9).  For other virus types, enter the common Strain_Name.	If more than one sequencing technology is used, comma-separate individual technologies.
<b>Dependent Fields</b>			
<b>Validation</b>	Validate field length	Validate field length	Field value should be one of valid values as in list.. NOTE: User can enter other value by prefixing 'OTH-'
<b>Message Code</b>	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_1_INVALID_VALUE
<b>Validation</b>		Influenza virus strain names must begin with the single letter A, B, C, or D, or be prefixed with rg-. Other strain names are not validated.	Field length including values from the Value List and free text following 'OTH-' must be less than 250 characters.
<b>Message Code</b>		QC_136_STRAIN_NAME_NOT_VALIDATED	Error_75_INVALID_FIELD_LENGTH_OTH
<b>Validation</b>		Host of origin element of strain name must be alphanumeric, underscore, dash, period, or single quote: a-z, A-Z, 0-9, -, ., ' , _	
<b>Message Code</b>		Error_137_INVALID_HOST_SPECIES_STRAIN_TEXT	
<b>Validation</b>		Geographical origin element of strain name must be alphanumeric, underscore, dash, period, or single quote: a-z, A-Z, 0-9, -, ., ' , _	
<b>Message Code</b>		Error_138_INVALID_LOCATION_STRAIN_TEXT	
<b>Validation</b>		Strain number element of strain name must be alphanumeric, underscore, dash, or period: a-z, A-Z, 0-9, -, .	
<b>Message Code</b>		Error_139_INVALID_STRAIN_NUMBER_STRAIN_TEXT	
<b>Validation</b>		Year of isolation element of strain name must be 4-digit year.	
<b>Message Code</b>		Error_110_INVALID_STRAIN_YEAR	
<b>Validation</b>		If strain is influenza A, subtype must be present as text between parentheses.	
<b>Message Code</b>		Error_112_INVALID_STRAIN_SUBTYPE	
<b>Validation</b>		Brackets cannot be present if strain does not have rg- prefix.	
<b>Message Code</b>		Error_119_INVALID_STRAIN_BRACKETS	

	Forward Primer	Reverse Primer	Assembly Method	Assembler Version
<b>Input Type</b>	Text Field	Text Field	Text Field	Text Field
<b>Definition</b>	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	The version of the assembly program used or, if not available, the date the assemblies were made
<b>Format</b>	Name:Sequence Maximum length: 500 characters	Name:Sequence Maximum length: 500 characters	Text Maximum length: 150 characters	Text DD-Mon-YYYY DD-Mon-YY D-Mon-YYYY D-Mon-YY Maximum length: 50 characters
<b>Value List</b>	Text U	Text U	None	Text Date
<b>Curation</b>	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.	The entry must be the version number of the assembly program used in format v.x.x or the date the assemblies were created.
<b>Examples</b>	fwd_seq:catgtttcacaaaggaaa, or U	rev_seq:atgtggatgcagattgtgga, or U	IonTorrent Assembler, or BioEdit	v.3.2, 03-Mar-2011, 3-Mar-2011, or 3-Mar-11
<b>Notes</b>	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.	If more than one assembly method is used, comma-separate individual versions.
<b>Dependent Fields</b>				
<b>Validation</b>	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	Validate field length
<b>Message Code</b>	Error_83_INVALID_PRIMER_NAME	Error_83_INVALID_PRIMER_NAME	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH
<b>Validation</b>	Validate field length	Validate field length		Validate assembler version format or date format.
<b>Message Code</b>	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH		Error_64_VALIDATION_ASSEMBLER_VERSION
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				
<b>Message Code</b>				
<b>Validation</b>				

	Coverage	File Name	Comments
<b>Input Type</b>	Text Field	Text Field	Text Field
<b>Definition</b>	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
<b>Format</b>	Text Maximum length: 50 characters	Text Maximum length: 2000 characters	Text Maximum length: 2000 characters
<b>Value List</b>	Number U	Text	Text NA
<b>Curation</b>	The entry must be a number or enter U if unknown.	The entry must be the full file name, with extension.	None
<b>Examples</b>	25.47	22258468_1.fasta	NA
<b>Notes</b>	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.
<b>Dependent Fields</b>			
<b>Validation</b>	Validate field length	Validate field length	Validate field length
<b>Message Code</b>	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH	Error_70_INVALID_FIELD_LENGTH
<b>Validation</b>	Coverage has an invalid value. Must be a number or U	File not found in ZIP file	
<b>Message Code</b>	Error_65_VALIDATION_COVERAGE	Error_87_DATA_FILE_MISSING_FROM_ZIP	
<b>Validation</b>			
<b>Message Code</b>			
<b>Validation</b>			
<b>Message Code</b>			
<b>Validation</b>			
<b>Message Code</b>			
<b>Validation</b>			
<b>Message Code</b>			
<b>Validation</b>			
<b>Message Code</b>			