

DPCC Data Standard Reference for SRA Run v1.1

| Input Type       | Project_Identifier   | Contributing_Institution   | Sample_Identifier   | BioProject_Accession_Number  | Embargo_End_Date   |
|------------------|--|--|---|--|--|
| Text Field       | Text Field   | Text Field   | Text Field  | Text Field   | Date Field   |
| Definition       | A unique Project Identifier generated by the DPCC by combining the Center-generated Project Code and a random 4-digit number | The institution code for the group that PERFORMED THE ANALYSIS. Must be one of the CEIRS institution codes assigned by NIAID                           | Identifier initially assigned to each sample collected. If multiple samples are taken from the same host, each sample should have its own identifier.   | The BioProject accession number associated with the sequence submission.   | In case an embargo of the information is needed, the date that the information should be released to the public databases by the DPCC  |
| Format           | Project_Code_XXXX<br>Maximum length: 21 characters   | Center three-letter code followed by three digits.<br>Maximum length: 6 characters   | Center-specific<br>Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _<br>Maximum length: 50 characters  | BioProject ID<br>Maximum length: 15 characters   | DD-Mon-YYYY<br>DD-Mon-YY<br>NA<br>Maximum length: 11 characters  |
| Value List       | None   | None   | None  | Text   | Date<br>NA   |
| Curation         | The entry must be a Project Identifier value registered with the DPCC.   | The entry must be an Institution Code value registered with the DPCC.<br><br>Center 3-letter codes are case-sensitive and must be entered in all caps. | For submissions derived from surveillance samples, Sample_Identifier must match an identifier from a previous surveillance submission.<br><br>For non-surveillance-related submissions, enter a unique Sample_Identifier. | The entry must be a valid BioProject accession number previously registered with the DPCC.   | 1. Leading 0 in DD is optional.<br>2. Month must match the first three letters of the month. Month is NOT case-sensitive.<br>3. Years may have two or four digits.<br>4. Date must conform to NIAID data release policies. |
| Examples         | SJCProj02_4001   | SJC101   | 22258468  | PRJNA37813   | 3-Mar-2011,<br>03-Mar-2011,<br>03-MAR-2011,<br>3-MAR-11, or<br>NA  |
| Notes            |  | A different Contributing_Institution may have performed the original collection of the surveillance sample that was analyzed.                          | The Sample_Identifier initially assigned to the surveillance sample must be provided.   | The sequence must be linked to a BioProject.<br><br>If no relevant BioProject exists, one must be created using the BioProject Registration process. | An exact date must be provided.<br><br>If Embargo_End_Date is NA, information will be released without delay.<br><br>Embargo_End_Date cannot be more than 12 months after submission.                                      |
| Dependent Fields |  |  |   |  |  |
| Validation       | Project_Identifier should be a valid project identifier.   | Must be an Institution Code value registered with the DPCC   | Validate field length   | Must be a valid BioProject accession number.   | Embargo_End_Date must be exact date and cannot be a past date.   |
| Message Code     | Error_9_PROJECT_NOT_FOUND  | Error_1_INVALID_VALUE  | Error_70_INVALID_FIELD_LENGTH   | Error_104_INVALID_BIOPROJECT   | Error_7_INVALID_EMBARGO_DATE   |
| Validation       | Project_Identifier exists but user does not have permission to access or edit the project.                                   |  | If Surveillance_Sample is Y, Sample_Identifier should match a surveillance record.  | Validate field length  | Embargo_End_Date cannot be more than 12 months after submission.   |
| Message Code     | Error_4_DENIED_USER_ACCESS   |  | Error_97_SAMPLE_ID_NOT_FOUND  | Error_70_INVALID_FIELD_LENGTH  | Error_100_EMBARGO_DATE_12M   |
| Validation       |  |  | If Surveillance_Sample is N, Sample_Identifier should be unique across all DPCC data.   |  |  |
| Message Code     |  |  | Error_66_NON_UNIQUE_SAMPLE_ID   |  |  |
| Validation       |  |  |   |  |  |
| Message Code     |  |  |   |  |  |
| Validation       |  |  |   |  |  |
| Message Code     |  |  |   |  |  |
| Validation       |  |  |   |  |  |
| Message Code     |  |  |   |  |  |
| Validation       |  |  |   |  |  |
| Message Code     |  |  |   |  |  |
| Validation       |  |  |   |  |  |
| Message Code     |  |  |   |  |  |

|                         | Publication PMID   | Submission Title   | Laboratory Name                                   | Study Type  | Study Abstract   |
|-------------------------|--|--|---|---|--|
| <b>Input Type</b>       | Text Field   | Text Field   | Text Field  | Text Field  | Text Field   |
| <b>Definition</b>       | The PubMed Unique Identifier (PMID) for the publication in which the sequence was published  | Short text that can be used to call out submission records in searches or displays | Laboratory name within the submitting institution | A controlled vocabulary for expressing the overall purpose of the study   | Briefly describes the goals, purpose, and scope of the study   |
| <b>Format</b>           | PMID<br>Maximum length: 50 characters  | Text<br>Maximum length: 200 characters   | Text<br>Maximum length: 100 characters            | Text<br>Maximum length: 26 characters   | Text<br>Maximum length: 500 characters   |
| <b>Value List</b>       | Text<br>NA   | None   | Text<br>NA  | Cancer Genomics<br>Epigenetics<br>Exome Sequencing<br>Forensic or Paleo-genomics<br>Gene Regulation Study<br>Metagenomics<br>Pooled Clone Sequencing<br>Population Genomics<br>Resequencing<br>RNASeq<br>Synthetic Genomics<br>Transcriptosome Analysis<br>Whole Genome Sequencing<br>Other | Text<br>NA   |
| <b>Curation</b>         | The entry must be a valid PMID number; 7 or 8-digit number with no leading zeros.<br><a href="http://www.nlm.nih.gov/bsd/mms/medlineelements.html#pmid">http://www.nlm.nih.gov/bsd/mms/medlineelements.html#pmid</a> | None   | None  | The entry must be one and only one member of the Value List.  | None   |
| <b>Examples</b>         | 19465683, or<br>NA   | SRA_Submission_13_Apr_4_2014   | Influenza Genome Sequencing Center                | Exome Sequencing  | The HA gene from the influenza strain A/WSN/1933 (H1N1) was mutagenized, and then the mutant libraries were used to grow mutant viruses. Both the pre and post-selection libraries were analyzed by Illumina sequencing along with unmutated controls. |
| <b>Notes</b>            | If the sequence was referenced in multiple publications, comma-separate the individual PMIDs.<br><br>Enter NA if the sequence was not referenced in a publication.   |  |   | Other signifies a value that is not present on the Value List. This field does not accept free text.  | If NA is entered and a PMID is provided, Study_Abstract will be filled using the publication abstract.<br><br>If NA is entered and no publication is available, Study_Abstract will be filled using DPCC registered project description.               |
| <b>Dependent Fields</b> | Study_Abstract   |  |   |   | Publication_PMIID  |
| <b>Validation</b>       | Validate field length  | Validate field length  | Validate field length                             | Field value should be one of valid values as in list.   | Validate field length  |
| <b>Message Code</b>     | Error_70_INVALID_FIELD_LENGTH  | Error_70_INVALID_FIELD_LENGTH  | Error_70_INVALID_FIELD_LENGTH                     | Error_1_INVALID_VALUE   | Error_70_INVALID_FIELD_LENGTH  |
| <b>Validation</b>       | Must be a valid PMID.  |  |   |   |  |
| <b>Message Code</b>     | Error_96_INVALID_PMIID   |  |   |   |  |
| <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>     |  |  |   |   |  |

|                         | Surveillance_Sample  | Sample_Description  | Virus_Identifier  | Strain_Name  |
|-------------------------|--|---|---|--|
| <b>Input Type</b>       | Text Field   | Text Field  | Text Field  | Text Field   |
| <b>Definition</b>       | Is the sequence derived from a surveillance sample?  | Short text that can be used to identify the sample  | A unique laboratory identifier assigned to this virus by the collector  | The WHO strain name of the virus that was sequenced  |
| <b>Format</b>           | Text<br>Maximum length: 1 character  | Text<br>Maximum length: 300 characters  | Text<br>Allowed characters include alphanumeric, hyphen, and underscore: a-z, A-Z, 0-9, -, _<br>Maximum length: 50 characters                         | Antigenic Type/Host of Origin/Geographical Origin/Strain Number/Year of Isolation (Subtype)<br>rg-Antigenic Type/Host of Origin/Geographical Origin/Strain Number/Year of Isolation [rg details] (Subtype)<br>Maximum length: 150 characters   |
| <b>Value List</b>       | Y<br>N   | Text<br>NA  | Text<br>NA  | Text<br>U  |
| <b>Curation</b>         | The entry must be one and only one member of the Value List.                               | None  | None  | WHO strain naming convention: Fields must be ordered as follows and separated with the ' ' character:<br>1. The antigenic type (e.g., A, B, C, D)<br>2. The host of origin (e.g., swine, equine, chicken. For human-origin viruses, no host of origin designation is given.)<br>3. Geographical origin (e.g., Denver, Taiwan)<br>4. Strain number (e.g., 15, 7)<br>5. Year of isolation (e.g., 2009, 1934)<br>6. For influenza A viruses, the hemagglutinin and neuraminidase antigen description in parentheses (e.g., (H1N1), (H3N2))<br><br>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. |
| <b>Examples</b>         | Y  | Run_1_060315  | 100245  | A/Hong Kong/1/1968 (H3N2),<br>A/chicken/Fujian/4/2002 (H3N6),<br>A/chicken/Fujian/4/2002 (HxNx),<br>A/chicken/Fujian/4/2002 (mixed),<br>A/swine/Iowa/233-56/2011 (H3N2v),<br>A/duck/Alberta/35/1976 (H1N1),<br>B/Hong Kong/432/2014,<br>C/Texas/1987/6/2011,<br>rg-A/Puerto Rico/8/1934 (H1N1), or<br>rg-A/Egypt/N03072/2010 [PR8 internal R 6+2] (H5N1)<br>Reference: <a href="http://www.cdc.gov/flu/about/viruses/types.htm">http://www.cdc.gov/flu/about/viruses/types.htm</a>   |
| <b>Notes</b>            | Y = Yes<br>N = No<br><br>Enter N if the sequence is a lab strain or reverse genetic virus. | This field provides an additional layer of tracking if multiple sequence samples were isolated from a single surveillance sample. | This field provides an additional layer of tracking if Strain_Name is not unique or multiple viruses were isolated from a single surveillance sample. | (HxNx) can be used in cases where a partial subtype has been determined (e.g., H5Nx).<br><br>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).<br><br>For other virus types, enter the common Strain_Name.   |
| <b>Dependent Fields</b> |  |   |   |  |
| <b>Validation</b>       | Field value should be one of valid values as in list.                                      | Validate field length   | Validate field length   | Validate field length  |
| <b>Message Code</b>     | Error_1_INVALID_VALUE  | Error_70_INVALID_FIELD_LENGTH   | Error_70_INVALID_FIELD_LENGTH   | Error_70_INVALID_FIELD_LENGTH  |
| <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.   |
| <b>Message Code</b>     |  |   |   | QC_136_STRAIN_NAME_NOT_VALIDATED   |
| <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  |



|                         | Insert_Size  | Insert_Size_Standard_Deviation  | Probe_Set  | Locus_Name  | Library_Construction_Protocol   |
|-------------------------|--|---|--|---|---|
| <b>Input Type</b>       | Text Field   | Text Field  | Text Field   | Text Field  | Text Field  |
| <b>Definition</b>       | Size of the insert for paired reads in KB                                    | Standard deviation of insert size   | Reference to an archived primer or probe set                             | Name(s) the gene(s) or locus(oci) or other genomic feature(s) targeted by the sequence  | Additional details about the experimental design and molecular strategies including hybrid selection and affinity capture reagents or any detail that distinguishes the experiment from other similar experiments |
| <b>Format</b>           | Number<br>Maximum length: 20 characters                                      | Number<br>Maximum length: 20 characters   | Text<br>Maximum length: 50 characters                                    | Text<br>Maximum length: 30 characters   | Text<br>Maximum length: 500 characters  |
| <b>Value List</b>       | Number<br>U<br>NA  | Number<br>U<br>NA   | Text<br>NA   | 16S rRNA<br>18S rRNA<br>COX1<br>exome<br>ITS1-5.8S-ITS2<br>matK<br>OTH-<br>NA   | Text<br>NA  |
| <b>Curation</b>         | None   | None  | None   | The entry must be one and only one member of the Value List.  | None  |
| <b>Examples</b>         | 165  | 92.8  | dbProbe  | OTH-RBCL  | 50-nt overlapping paired-end reads of influenza nucleoprotein PCR amplicons processed with Nextera.   |
| <b>Notes</b>            | If Paired is selected for Library_Layout, Insert_Size must be a number or U. | If Paired is selected for Library_Layout, Insert_Size_Standard_Deviation must be a number or U. | If Targeted loci is selected for Library_Layout, Probe_Set must be text. | If Targeted loci is selected for Library_Layout, Locus_Name must be text.<br><br>If OTH- is selected, append free text to describe. |   |
| <b>Dependent Fields</b> | Library_Layout   | Library_Layout  | Library_Layout   | Library_Layout  |   |
| <b>Validation</b>       | Validate field length  | 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-'                          | Validate field length   |
| <b>Message Code</b>     | Error_70_INVALID_FIELD_LENGTH  | Error_70_INVALID_FIELD_LENGTH   | Error_70_INVALID_FIELD_LENGTH  | Error_1_INVALID_VALUE   | Error_70_INVALID_FIELD_LENGTH   |
| <b>Validation</b>       | Value must be number, U, or NA   | Value must be number, U, or NA  |  | Validate field length   |   |
| <b>Message Code</b>     | Error_86_INVALID_NUM_U_NA  | Error_86_INVALID_NUM_U_NA   |  | Error_75_INVALID_FIELD_LENGTH_OTH   |   |
| <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>     |  |   |  |   |   |



| Input Type       | Internal Field Only  |  |  |   |   |
|------------------|--|--|--|---|---|
|                  | Run Title  | Run Date   | Run File Name  | Comments  | SRA Study Accession Numbers   |
| Definition       | Text Field<br>Short text that can be used to distinctly call out individual runs in searches or displays | Date Field<br>Date when the run took place   | Text Field<br>The name(s) of the run data file(s).   | Text Field<br>Text describing anything else of interest related to the submission | Text Field<br><b>Internal Field Only</b><br>Accession numbers associated with the SRA Study submission. |
| Format           | Text<br>Maximum length: 100 characters   | DD-Mon-YYYY<br>DD-Mon-YY<br>Mon-YYYY<br>Mon-YY<br>YYYY<br>YY<br>Maximum length: 11 characters  | Text<br>Maximum length: 2000 characters  | Text<br>Maximum length: 2000 characters   | Text<br>Maximum length: 200 characters  |
| Value List       | Text<br>NA   | Date<br>U  | Text   | Text<br>NA  | Text  |
| Curration        | None   | 1. Leading 0 in DD is optional.<br>2. Month must match the first three letters of the month. Month is NOT case-sensitive.<br>3. Years may have two or four digits.<br>4. Use U (Unknown) if date is not known. | The file or files included in the submission must have one of the following extensions:<br>BAM<br>FA, FAS, or FASTA<br>FASTQ<br>H5 or HDF5<br>SFF<br>SRA   | None  | None  |
| Examples         | Run 1  | 3-Mar-2011,<br>03-Mar-2011,<br>03-MAR-2011,<br>Mar-2011,<br>3-Mar-11,<br>2011,<br>11, or<br>U  | 707971081216.1.TCAG.bam  | NA  | None  |
| Notes            |  |  | The entry must be the full file name, with extension (e.g., sequence_data.fastq).<br><br>Users can submit multiple comma-separated files. File name should be unique in the zip file.<br><br>BAM = binary text file that contains sequence alignment data<br>FA, FAS, or FASTA = text-based format for nucleotide sequences<br>FASTQ = text-based format for storing sequence data and corresponding quality scores<br>H5 or HDF5 = output files produced by PacBio and Oxford Nanopore sequencers<br>SFF = standard flowgram format<br>SAR = Sequence Read Archive format | If there are no comments, enter NA.   |   |
| Dependent Fields |  |  |  |   |   |
| Validation       | Validate field length  | Validate date format   | Validate field length  | Validate field length   |   |
| Message Code     | Error_70_INVALID_FIELD_LENGTH  | Error_67_INVALID_DATE_WITH_U   | Error_70_INVALID_FIELD_LENGTH  | Error_70_INVALID_FIELD_LENGTH   |   |
| Validation       |  | Date cannot be future date   | Validate file extension  |   |   |
| Message Code     |  | Error_68_INVALID_FUTURE_DATE   | Error_132_SRA_RUN_UNSUPPORTED_FILE_EXT   |   |   |
| Validation       |  |  | File is missing or is empty  |   |   |
| Message Code     |  |  | Error_137_FILE_MISSING_OR_EMPTY  |   |   |
| Validation       |  |  |  |   |   |
| Message Code     |  |  |  |   |   |
| Validation       |  |  |  |   |   |
| Message Code     |  |  |  |   |   |
| Validation       |  |  |  |   |   |
| Message Code     |  |  |  |   |   |
| Validation       |  |  |  |   |   |
| Message Code     |  |  |  |   |   |
| Validation       |  |  |  |   |   |
| Message Code     |  |  |  |   |   |
| Validation       |  |  |  |   |   |

