

1
2
3
4 **NCBI will no longer make taxonomy identifiers for**
5 **individual influenza strains on January 15, 2018**
6
7

8 Authors: Hatcher, Eneida L^{1*}; Bao, Yiming^{2*}; Amedeo, Paolo³; Blinkova, Olga¹;
9 Cochrane, Guy⁴; Fedorova, Nadia B³; Gruner, William E⁵; Leipe, Detlef D¹; Nakamura,
10 Yasukazu⁶; Ostapchuk, Yuri¹; Palanigobu, Vasuki¹; Sanders, Robert¹; Schoch, Conrad¹;
11 Smith, Catherine⁷; Wentworth, David E⁷; Yankie, Linda¹; Zhdanov, Sergey A; Karsch-
12 Mizrahi, Ilene¹; Brister, J. Rodney¹.

- 13
14 1. Information Engineering Branch, National Center for Biotechnology Information,
15 National Library of Medicine, National Institutes of Health. Bethesda, MD, USA.
16 2. Beijing Institute of Genomics, Chinese Academy of Sciences. Beijing, China. J.
17 Craig Venter Institute. Rockville, MD, USA.
18 3. J. Craig Venter Institute. Rockville, MD, USA.
19 4. European Molecular Biology Laboratory, European Bioinformatics Institute,
20 Wellcome Trust Genome Campus. Hinxton, Cambridge, UK.
21 5. United States Air Force Materiel Command, School of Aerospace Medicine.
22 Wright-Patterson AFB, OH, USA.
23 6. DNA Databank of Japan, National Institute of Genetics. Shizuoka 411-8540,
24 Japan.
25 7. National Center for Immunization and Respiratory Diseases, Centers for Disease
26 Control and Prevention. Atlanta, GA, USA.

27 *Co-first authors

28
29 Corresponding Author: Eneida Hatcher, Eneida.Hatcher@NIH.gov
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

46 Abstract

47
48 Currently the National Center of Biotechnology Information (NCBI) assigns individual
49 taxonomy identifiers to each distinct influenza virus isolate submitted to GenBank. To
50 support this practice, individual flu isolates must be manually added to the NCBI
51 taxonomy database and unique taxonomy identifiers generated. This added layer of
52 manual processing is unique to influenza virus and prevents automatization of the flu
53 sequence submission process. Here we outline a new NCBI policy that normalizes
54 Influenza virus taxonomy processing but maintains features supported by the previous
55 approach. This change will reduce the amount of manual handling necessary for flu
56 submissions and pave the way for increased automation of the submissions process.
57 While this automation may disrupt some historic practices, it will better align influenza
58 virus data processing with other viruses and ultimately lower the submission burden on
59 data providers.

60
61
62
63

64 Introduction

65
66 GenBank is a member of the International Nucleotide Sequence Database Collaboration
67 (INSDC) (Cochrane et al. 2016) - data repositories dedicated to providing public access
68 to biological sequence data. Viral taxonomy within INSDC databases follows the
69 guidelines provided by the International Committee on the Taxonomy of Viruses (ICTV).
70 The scope of the ICTV mandate extends from species to higher level taxa, and no
71 subspecific taxa are maintained by the ICTV (Adams et al. 2017).

72
73 All viral sequences submitted to GenBank and other INSDC repositories are assigned to
74 a species. Sequences from characterized viruses are assigned to their pre-existing
75 species. Sequences from novel viruses are assigned to newly created, unclassified
76 species. Typically, subspecific taxonomic ranks are not created at the time of submission,
77 though some formally unranked subspecific taxa are made during post-submission
78 taxonomic revisions. Creation of new viral taxa within the NCBI taxonomy database -
79 whether families, species, or subspecific ranks - requires manual validation and database
80 operations.

81
82 There are currently more than 550,000 *Influenzavirus A*, *B*, and *C* nucleotide sequences
83 in GenBank - nearly twenty percent of the entire viral nucleotide sequence content of this
84 database (see Table 1). These sequences represent a coordinated effort by the
85 international scientific community to share critical public health data (Bao et al. 2008),
86 and it is imperative that GenBank provides efficient data distribution pathways to support
87 this and similar efforts. Given the number of influenza virus sequences generated by the
88 scientific community, efficient distribution to GenBank can only be sustained through
89 increased automation of the submissions process.

90
91

92 **Table 1.** Number of influenza virus nucleotide sequences submitted by year.

Year	Sequences submitted
2014	50,383
2015	59,300
2016	84,015
Total*	551,664
*Total includes all sequences submitted before October 2017	

93
94
95
96
97

Problem

98 Each newly submitted influenza isolate is given a unique strain name, which must be
99 manually validated and added to the NCBI taxonomy database. Almost 14,000 individual
100 influenza strain names were added to the NCBI taxonomy database in 2016 (see Table
101 2). With more than 120,000 strain names total, this manual step has become a significant
102 impediment to influenza virus sequence submission automation that impacts both
103 traditional GenBank submissions and large-scale submissions.

104
105

Table 2. Number of influenza virus strain names created by year.

Species	2014	2015	2016	Total*	Number of strains which include complete genome sets*
Influenza A virus	8,567	12,600	11,120	103,927	37,283
Influenza B virus	1,301	3,874	2,761	16,394	6,066
Influenza C virus	9	36	48	272	26
Influenza D virus				45	25
*Includes all sequences submitted before October 2017					

106
107

108 The burden of manual taxonomy operations falls not only on GenBank staff, but
109 submitters also face delays and/or must complete extra steps to prepare submissions.
110 For large scale submitters, this means submitting lists of proposed strain names and
111 waiting for their approval before being able to submit sequence data. As it stands, this
112 situation is at odds with both the need for GenBank to introduce new, more efficient
113 submissions pathways and the desire of data providers to reduce submissions burden
114 and timelines.

115

116 **Proposed solution**

117

118 To facilitate more efficient submissions, NCBI will stop assigning strain level taxonomy
119 designations to influenza virus sequences on January 15, 2018. From that point forward,
120 influenza sequences will be assigned to the relevant species and will be associated with
121 a species-level taxonomy identifier. The taxonomy identifiers shown in Table 3 will be
122 automatically assigned to sequences as part of the GenBank submission process after
123 the species of the sequence is verified by homology.

124

125 **Table 3.** Influenza virus species NCBI Taxonomy identification numbers.

Species	Taxonomy identifier
Influenza A virus	11320
Influenza B virus	11520
Influenza C virus	11552
Influenza D virus	1511084

126

127

128 Once this change is made, both species and strain names will continue to be visible in
129 the DEFINITION lines of GenBank records, BLAST results, NCBI's ftp sites, and the NCBI
130 Influenza Virus Resource. However, there will be changes to influenza GenBank records.
131 The organism name will no longer include both species and strain information. Only the
132 species name will be listed, and strain information will be included in the "/strain" field as
133 illustrated in Fig. 1 (Please see Supp. Fig. 1 to see an example of a complete GenBank
134 record).

135

136

Current format:

DEFINITION Influenza A virus (A/alien/Mars/1/2033(H20N15)) segment 4, complete sequence.

SOURCE Influenza A virus (A/alien/Mars/1/2033(H20N15))
ORGANISM Influenza A virus (A/alien/Mars/1/2033(H20N15))

```

source 1..1740
       /organism="Influenza A virus (A/alien/Mars/1/2033(H20N15))"
       /mol_type="viral cRNA"
       /strain= "A/alien/Mars/1/2033"
       /serotype="H20N15"
       /host="alien"
       /segment="4"

```

Proposed format after December 1, 2017:

DEFINITION Influenza A virus (A/alien/Mars/1/2033(H20N15)) segment 4, complete sequence.

SOURCE Influenza A virus
ORGANISM Influenza A virus

```

source 1..1740
       /organism="Influenza A virus"
       /mol_type="viral cRNA"
       /strain= "A/alien/Mars/1/2033"
       /serotype="H20N15"
       /host="alien"
       /segment="4"

```

137

138 **Figure 1.** Differences in influenza virus GenBank records after NCBI no longer makes strain-level
139 organism names for flu sequences.

140

141 GenBank nucleotide (<https://www.ncbi.nlm.nih.gov/nucleotide/>) and protein
142 (<https://www.ncbi.nlm.nih.gov/protein/>) databases will continue to support searches
143 based on components of strain names (e.g. A/New York/61A/2003), as will the NCBI
144 Influenza Virus Resource (<https://www.ncbi.nlm.nih.gov/genome/viruses/variation/flu/>).
145 This resource also supports sequence downloads that include user-defined DEFINITION
146 lines derived from isolate descriptors and other sequence metadata.

147

148 New strain-level taxonomy names will not be created in NCBI's taxonomy database. For
149 example, pages such as this one
150 <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=311639> will not be
151 made for individual influenza virus strains after January 15, 2018. All new submissions
152 will point to species level pages such as this one for Influenza A virus -
153 <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Info&id=11320>.

154

155 Influenza virus strains will not be found by searching the NCBI Taxonomy database
156 (<https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi>) using new strain names.
157 However, previously existing strain-level names will not be removed from the NCBI
158 taxonomy database and will still be found using strain name searches. The NCBI
159 Taxonomy pages for the *Orthomyxoviridae* family, *Influenzavirus* genera, influenza virus
160 species, and legacy strain rank pages will include a highlighted banner to make users
161 aware of the changes in policy.

162
163 The NCBI BioSample resource (Barrett et al. 2012) will continue to provide unique
164 accessions for individual samples/isolates if requested. This accession will be linked to
165 submitted sequences in SRA and GenBank. BioSample records also provide storage for
166 highly detailed sample descriptors, providing a much richer biological context to
167 sequences compared to the source descriptors in the GenBank record. More information
168 about BioSample can be found here: <https://www.ncbi.nlm.nih.gov/biosample>.

169
170

171 **Submitting Influenza sequences to GenBank**

172

173 As of February 1, 2018, there will be two options to submit Influenza A, B, or C virus
174 sequence data: an interactive web wizard and a programmatic interface. Both options will
175 facilitate submissions by including an automated process for annotation, so users will not
176 have to include annotation files. To see how your sequence will be annotated, please use
177 the NCBI Influenzavirus Annotation Tool at
178 <https://www.ncbi.nlm.nih.gov/genomes/FLU/annotation/>. We expect to add the capability
179 to accept Influenza D virus sequences through these submission tools in early 2018.

180

181 The new interactive web wizard for submitting influenza sequences is at
182 <https://submit.ncbi.nlm.nih.gov/subs/genbank/>. A description of the submission process
183 and details on the required files are provided at
184 <https://submit.ncbi.nlm.nih.gov/genbank/help/> (see also Table 4 and the Supplemental
185 text documents). Data providers will need to submit FASTA-formatted sequence files
186 and a tab-separated table with source information - sequence_ID, isolate, collection-
187 date, host, collection country/geographic origin, isolation-source, and serotype for
188 Influenza virus A. We encourage submitters to include passage history in a “note”
189 column if known, however it is not required at this time. Influenza strain names will be
190 constructed from the source information table in the format “virus
191 type/country/isolate/year(serotype),” although serotype will only be included for
192 influenza A viruses.

193

194 A programmatic submission interface will be available for centers submitting large-scale
195 surveillance data. Submitters are encouraged to contact NCBI at [gb-
196 admin@ncbi.nlm.nih.gov](mailto:gb-admin@ncbi.nlm.nih.gov) prior to submission to ensure that the information included in
197 the submission is properly formatted. An archive file (.zip or tar.gz) containing the
198 sequence(s) in FASTA format, source information table in tab-delimited format, a
199 submission template file, and a submission form which includes information necessary
200 for processing should be submitted. Each file must have a specific extension, as shown

201 in Table 4. Users can generate submission template files by entering their information at
 202 <https://submit.ncbi.nlm.nih.gov/genbank/template/submission/>, and saving a template
 203 that can be submitted with each of their submissions. Details on the submission process
 204 through the programmatic interface will be provided soon.

205
 206 Reassortant influenza strains play an important role in pandemics and vaccine
 207 research. Currently, reassortant viruses cannot be submitted through the web wizard,
 208 and users should contact gb-admin@ncbi.nlm.nih.gov if they would like to submit
 209 reassortant influenza viruses. Although the programmatic interface can be used to
 210 submit reassortants, we strongly encourage submitters to send an email first to ensure
 211 that their submission files are in the correct format. The strain name should follow the
 212 accepted format, e.g., Influenza A virus (A/reassortant/VG-552(Japan/5685/2016 x New
 213 York/5689/2016)(H3N2).

214

215 **Table 4.** Files required for submission of influenza virus sequences to NCBI.

File	Requirement - Web wizard	Requirement - Programmatic Interface (mandatory file extension)	Description
Sequence data	Required	Required (.fsa)	Nucleotide sequences in FASTA format. The sequence identifier must match the Sequence_ID used in the source information table. Spaces are not allowed in the sequence identifier.
Source information table	Required	Required* (.src)	Tab-delimited table which must include: Wizard submissions: sequence_ID, isolate, collection-date, host, collection country, isolation-source, serotype. Programmatic submissions: sequence_ID, strain, collection-date, host, collection country, isolation-source, serotype and any other source metadata that fits the submission. See here for a list of all source modifiers: https://www.ncbi.nlm.nih.gov/WebSub/html/help/genbank-source-table.html#modifiers
Submission template	N/A	Required (.sbt)	Submitter names and organizations, and publications associated with or describing the sequence. Can be generated at https://submit.ncbi.nlm.nih.gov/genbank/template/submission/
Submission form	N/A	Required (.xml)	Includes instructions for the submission pipeline.
Structured comment	N/A	Optional (.cmt)	Additional metadata that does not have designated fields in GenBank records. For more information, please see https://www.ncbi.nlm.nih.gov/genbank/structuredcomment/
*Required unless all required source information is provided in the FASTA file as part of the defines			

216
217 After successful submission, an email will be sent to inform the submitter of the new
218 accession numbers associated with the sequence IDs, the expected release date, and a
219 preview of the records in GenBank format.

220
221 In order to facilitate this transition for providers and users of flu data, we will maintain a
222 webpage with the changes outlined in this paper and include links to it from the Influenza
223 virus resource at <https://www.ncbi.nlm.nih.gov/genome/viruses/variation/flu/help-center/>.

224 225 **Summary**

226
227 The change in GenBank influenza virus strain name processing will support improved
228 submissions efficiency. While individual strain names will no longer be added to the NCBI
229 taxonomy database after January 15, 2018, virus strain names will remain available and
230 searchable through NCBI resources. GenBank staff will provide assistance to data
231 providers to ensure a smooth transition to this new policy.

232 233 **Funding**

234 This research was supported in part by the Intramural Research Program of the NIH,
235 National Library of Medicine.

236 This project has been funded in part with federal funds from the National Institute of
237 Allergy and Infectious Diseases, National Institutes of Health, Department of Health and
238 Human Services under Award Number U19AI110819.

239 240 241 **References**

- 242 1. Adams MJ, Lefkowitz EJ, King AM, Harrach B, Harrison RL, Knowles NJ, Kropinski AM, Krupovic
243 M, Kuhn JH, Mushegian AR, Nibert ML, Sabanadzovic S, Sanfacon H, Siddell SG, Simmonds P,
244 Varsani A, Zerbini FM, Orton RJ, Smith DB, Gorbalenya AE, and Davison AJ. 2017. 50 years of the
245 International Committee on Taxonomy of Viruses: progress and prospects. *Arch Virol* 162:1441-
246 1446. 10.1007/s00705-016-3215-y
- 247 2. Bao Y, Bolotov P, Dernovoy D, Kiryutin B, Zaslavsky L, Tatusova T, Ostell J, and Lipman D. 2008.
248 The influenza virus resource at the National Center for Biotechnology Information. *J Virol*
249 82:596-601. 10.1128/JVI.02005-07
- 250 3. Barrett T, Clark K, Gevorgyan R, Gorelenkov V, Gribov E, Karsch-Mizrachi I, Kimelman M, Pruitt
251 KD, Resenchuk S, Tatusova T, Yaschenko E, and Ostell J. 2012. BioProject and BioSample
252 databases at NCBI: facilitating capture and organization of metadata. *Nucleic Acids Res* 40:D57-
253 63. 10.1093/nar/gkr1163
- 254 4. Cochrane G, Karsch-Mizrachi I, Takagi T, and International Nucleotide Sequence Database C.
255 2016. The International Nucleotide Sequence Database Collaboration. *Nucleic Acids Res* 44:D48-
256 50. 10.1093/nar/gkv1323

257

Influenza A virus (A/alien/Mars/1/2033(H20n15)) HA gene for hemagglutinin, complete sequence

GenBank: CY999999.1

[FASTA](#) [Graphics](#)Go to:

```

LOCUS       CY999999          1714 bp    RNA        linear    VRL 09-MAR-2033
DEFINITION  Influenza A virus (A/alien/Mars/1/2033(H20n15)) HA gene for
             hemagglutinin, complete sequence.
ACCESSION   CY999999
VERSION     CY999999.1
KEYWORDS    .
SOURCE      Influenza A virus
  ORGANISM  Influenza A virus
             Viruses; ssRNA viruses; ssRNA negative-strand viruses;
             Orthomyxoviridae; Influenzavirus A.
COMMENT     EXAMPLE OF A GENBANK RECORD FOR INFLUENZA VIRUSES SUBMITTED
             AFTER DECEMBER 1, 2017.
             COMPLETENESS: full length.
FEATURES             Location/Qualifiers
     source           1..1714
                     /organism="Influenza A virus"
                     /mol_type="genomic RNA"
                     /strain="A/alien/Mars/1/2033"
                     /serotype="H20N15"
                     /segment="4"
     gene             32..1714
                     /gene="HA"
     CDS              32..1714
                     /gene="HA"
                     /codon_start=1
                     /product="Hemagglutinin"
                     /translation="METISLITILLVVTASNADKICIGHQSTNSTETVDTLTETNVPV
                     THAKELLTHEHNGMLCATSLGHPLILDTCIEGLVYGNPSCDLLLGGREWSYIVERSS
                     AVNGTCYPGNVENLEELRRLFSSASSYQRIQIFPDTTWNVTYTGTSRACSGSFYRSMR
                     WLTKQSGFYFVQDAQYTNNRGKSIILFVWGIHHPPTYEQTNLYIRNDTTTSVTTEDLN
                     RTFKPVIQPRPLVNLQGRIDYYSVLKPGQTLRVRNGLIAPWYGHVLSGGSHGRI
                     LKTDLKGGNVCVQCQTEKGLNSTLPHFNISKYAFGTCPKYVRVNSLKLAVGLRNVPA
                     RSSRGLFGAIAGFIEGGWPLVAGWYGFQHSNDQGVGMAADRSTQKAIKDKITSKVNN
                     IVDKMNKQYEI IDHEFSEVETRLNMINNKIDDQIQDVWAYNAELLVLENQKTLDEHD
                     ANVNNLYNKVKRALGNSAMEDGKGC FELYHKDDQCMETIRNGTYNRRKYREESRLER
                     QKIEGVKLESEGTYKILTYSTVASSLVLAMGFAAFLFWAMSNNGSRCRNICI"
ORIGIN
1  gcaaaagcag  ggaattact  taactagcaa  aatggaaaca  atactactaa  taactatact
61  actagtagta  acagcaagca  atgcagataa  aatctgcatac  ggccaccagt  caacaaactc
121  cacagaaact  gtggacacgc  taacagaaac  caatgttcct  gtgacacatg  ccaaaagaatt
181  gctccacaca  gaggcataatg  gaatgctgtg  tgcaacaagc  ctgggacatc  ccctcattct
241  agacacatgc  actattgaag  gactagtcta  tggcaaccct  tcttgtagcc  tgctgttggg
301  aggaagagaa  tggctcctaca  tcgtcgaag  atcatcagct  gtaaatggaa  cgtgttacc
361  tgggaatgta  gaaaacctag  aggaactcag  gacactttt  agttccgta  gttcctacca
421  aagaatccaa  atcttcccag  acacaacctg  gaatgtgact  tacactggaa  caagcagagc
481  atgttcagg  tcattctaca  ggagtatgag  atggctgact  caaaagagcg  gttttacc
541  tgttcaagac  gcccaataca  caaataacag  gggaaagagc  attcttttcg  tgtggggcat
601  acatcaccca  cccacctata  cggagcaaac  aaattgtac  ataagaaacg  acacaacaac
661  aagcgtgaca  acagaagatt  tgaataggac  cttcaaacca  gtgataggc  caaggcccct
721  tgtcaatggt  ctgcagggaa  gaattgatta  ttattggtcg  gtactaaaac  caggccaaac
781  attgcgagta  cgatccaatg  ggaatcta  tgctccatgg  tatggacacg  ttcttccagg
841  agggagccat  ggaagaatcc  tgaagactga  tttaaaagg  ggtaattgtg  tagtgcaatg
901  tcagactgaa  aaaggtggct  taaacagtac  attgccattc  cacaatatca  gtaaatatgc
961  atttggaacc  tgccccaaat  atgtaagagt  taatagtctc  aaactggcag  tcggtctgag
1021  gaactgact  gctagatcaa  gtagaggact  atttggagcc  atagctgga  tcatagaagg
1081  aggttggcca  ggactagtgc  ctggctggt  tggttccag  cattcaaatg  atcaaggggt
1141  tggatggtc  gcagatagg  attcaactca  aaaggcaatt  gataaaaata  catccaaggt
1201  gaataatata  tgcgacaaga  tgaacaagca  atatgaaata  attgatcatg  aattcagtg
1261  ggttgaact  agactcaata  tgatcaata  taagattgat  gaccaaaac  aagacgtatg
1321  ggcataata  gcagaattgc  tagtactact  tgaaaatcaa  aaaaactcg  atgagcatga
1381  tgcgaactg  aacaatctat  ataacaaggt  gaagagggca  ctgggctcca  atgctatgga
1441  agatgggaaa  ggctgtttcg  agctatacca  taatgtgat  gatcagtgca  tggaaaacat
1501  tcggaacgg  acctataata  ggagaaagta  tagagaggaa  tcaagactag  aaagcgagaa
1561  aatagaggg  gttaaagctg  aatctgagg  aacttaca  atctcacc  tttattcgac
1621  tgtgcctca  tctctgtgc  ttgcaatgg  gtttctgccc  ttctgttct  gggccatgct
1681  caatggatct  tgcagatgca  acattgtgat  ataa
//

```

258
259
260

Supplemental Figure 1. Example of a GenBank record for after new Influenza taxonomy policy

261 **Supplemental Text Files**

- 262 A. Submission template file
- 263 B. Source information table including strain for programmatic interface, .src
- 264 C. Source information table including isolate for web wizard, .txt
- 265 D. Sequence data
- 266 E. Submission form