

# **OSAC 2021-S-0006 Standard for the Use of GenBank for Taxonomic Assignment of Wildlife**

*Wildlife Forensic Biology Subcommittee  
Biology Scientific Area Committee  
Organization of Scientific Area Committees (OSAC) for Forensic Science*



## **Draft OSAC Proposed Standard**

# **OSAC 2021-S-0006 Standard for the Use of GenBank for Taxonomic Assignment of Wildlife**

Prepared by  
Wildlife Forensic Biology Subcommittee  
Version: 2.0  
July 2021

---

### **Disclaimer:**

This OSAC Proposed Standard was written by the Wildlife Forensic Biology Subcommittee of the Organization of Scientific Area Committees (OSAC) for Forensic Science following a process that includes an [open comment period](#). This Proposed Standard will be submitted to a standards developing organization and is subject to change.

There may be references in an OSAC Proposed Standard to other publications under development by OSAC. The information in the Proposed Standard, and underlying concepts and methodologies, may be used by the forensic-science community before the completion of such companion publications.

Any identification of commercial equipment, instruments, or materials in the Proposed Standard is not a recommendation or endorsement by the U.S. Government and does not imply that the equipment, instruments, or materials are necessarily the best available for the purpose.

To be placed on the OSAC Registry, certain types of standards first must be reviewed by a Scientific and Technical Review Panel (STRP). The STRP process is vital to OSAC's mission of generating and recognizing scientifically sound standards for producing and interpreting forensic science results. The STRP shall provide critical and knowledgeable reviews of draft standards or of proposed revisions of standards previously published by standards developing organizations (SDOs) to ensure that the published methods that practitioners employ are scientifically valid, and the resulting claims are trustworthy.

The STRP panel will consist of an independent and diverse panel, including subject matter experts, human factors scientists, quality assurance personnel, and legal experts, which will be

tasked with evaluating the proposed standard based on a comprehensive list of science-based criteria.

For more information about this important process, please visit our website at:

<https://www.nist.gov/topics/organization-scientific-area-committees-forensic-science/scientific-technical-review-panels>

DRAFT

# 1 Standard for the Use of GenBank for Taxonomic Assignment of Wildlife

## 2 Foreword

3 This standard defines the requirements that shall be met when comparing evidentiary sequences to  
4 those in GenBank for taxonomic assignment of non-human samples. The aim is to provide a framework  
5 that will result in consistency in the wildlife forensic DNA community. Use of these standards is  
6 expected for forensic scientists with a working understanding of DNA sequencing.

7 This standard was developed by the Biology/ Wildlife Forensic Biology Subcommittee of the  
8 Organization of Scientific Area Committees. This standard is intended to assist those using GenBank for  
9 the taxonomic identification of wildlife in forensic casework.

10 All hyperlinks and web addresses shown in this document are current as of the publication date of this  
11 standard.

12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

**Keywords:** GenBank, BLAST, DNA, Public sequence databases, Taxonomic identification, Wildlife

34	<b>Table of Contents</b>	
35		
36		
37	<b>1. Scope</b>	<b>7</b>
38		
39	<b>2. Normative References</b>	<b>7</b>
40		
41	<b>3. Terms and Definitions</b>	<b>7</b>
42		
43	<b>4. Requirements</b>	<b>9</b>
44		
45	<b>Annex A (informative) Bibliography</b>	<b>14</b>
46		

**DRAFT**

## 47 **Standard for the Use of GenBank for Taxonomic Assignment of Wildlife**

### 48 **1. Scope**

49 This standard covers the requirements and recommendations for analysis and selection of DNA  
50 sequences retrieved from the National Center for Biotechnology Information's GenBank and  
51 their subsequent use as reference material for taxonomic identification of wildlife<sup>1</sup>. This standard does  
52 not cover the use of DNA sequences from other public sequence databases (*e.g.*, BOLD, UNITE), the  
53 protocol for downloading sequences from GenBank for inclusion in in-house databases, or the use of  
54 custom BLAST searches against GenBank. However, the criteria can be conceptually applied to other  
55 sequence databases.

### 56 **2. Normative References**

57 NCBI Field Guide Glossary available at  
58 <https://www.ncbi.nlm.nih.gov/Class/FieldGuide/glossary.html#>

59 Madden T. (2013). "The BLAST Sequence Analysis Tool." In: *The NCBI Handbook, 2nd ed.* Bethesda,  
60 MD. Available from <https://www.ncbi.nlm.nih.gov/books/NBK153387/>

61 ANSI/ASB Standard 019, First Edition. Wildlife Forensics General Standards, 2019.

62 ANSI/ASB Standard 029, First Edition. Report Writing in Wildlife Forensics: Morphology and  
63 Genetics, 2019

### 64 **3. Terms and Definitions**

65 For purposes of this document, the following definitions and acronyms apply:

#### 66 **3.1** 67 **alignment**

68 An arrangement of two or more nucleotide or protein sequences that is used to illustrate similarity  
69 among those sequences.

#### 70 **3.2** 71 **Basic Local Alignment Search Tool** 72 **BLAST**

73 The a) BLAST algorithm, and b) a suite of database search programs that implement variations of  
74 this algorithm to generate alignments between a nucleotide or protein sequence in a query, and  
75 nucleotide or protein sequences within a database.

#### 76 **3.3** 77 **expectation value** 78 **e-value**

79 The number of distinct alignments expected by chance; the default sorting metric in BLAST search  
80 results.

81 <sup>1</sup> For the purposes of this document, "wildlife" species are defined as non-human multicellular animals and plants,  
82 whether wild, captive-bred, or domesticated.

83

84 **3.4**  
85 **GenBank**

86 A public repository of DNA sequences maintained by the National Center for Biotechnology Information,  
87 part of the U.S. National Institutes of Health.

88 **3.5**  
89 **hit(s)**

90 Sequence(s) returned from GenBank when performing a BLAST search. Also known as a “subject  
91 sequence.”

92 **3.6**  
93 **interspecific**

94 Between members of different species.

95 **3.7**  
96 **intraspecific**

97 Between members of the same species.

98 **3.8**  
99 **National Center for Biotechnology Information**  
100 **NCBI**

101 The U.S. National Center for Biotechnology Information (NCBI) is located in Bethesda, Maryland and is  
102 part of the United States National Library of Medicine (a branch of the National Institutes of Health).  
103 NCBI houses a series of databases relevant to biotechnology and biomedicine and provides several  
104 bioinformatics tools for searching and analyzing the housed data.

105 **3.9**  
106 **phylogram**

107 A branching diagram that illustrates relationships amongst organisms. Phylograms are typically  
108 generated using genetic sequences and/or morphological characters.

109 **3.10**  
110 **query**

111 (n) The nucleotide or protein sequence that has an unknown source (*i.e.*, evidence sequence), or (v) the  
112 action of searching an unknown sequence against a database.

113 **3.11**  
114 **query coverage**

115 The percent of the query sequence length that is included in the aligned segment with a hit.

116 **3.12**  
117 **sequence identity**

118 The percentage or number of nucleotides or amino acids that are identical between two sequences.

119 **3.13**  
120 **subject sequence(s)**

121 A nucleotide or protein sequence(s) returned from a GenBank BLAST search. Also known as a “hit”.

122 **3.14**

123 **taxonomic identification**

124 Analyses to establish the classification of biological evidence to family, genus, species, etc. These  
125 analyses are based on class characters (*e.g.*, morphological, genetic) that are diagnostic for the  
126 taxonomic level in question.

127

128 **3.15**

129 **topology**

130 The branching structure of a phylogram.

131 **3.16**

132 **voucher specimen**

133 Biological specimen that is representative of its species in accordance with the relevant taxonomic  
134 authority and is therefore valid for comparative purposes. Voucher specimens are of known identity,  
135 and are curated with available associated geographic, field collection, and life history data.

136

137 **4. Requirements**

138 Details about the operation of BLAST can be found in Madden (2013), and detailed information on the  
139 terms in the BLAST output can be found in the NCBI Field Guide Glossary.

140 The following requirements and recommendations address criteria for the preparation and submission  
141 of evidentiary query sequences (4.1) and evaluation and interpretation of BLAST results from GenBank  
142 (4.2, 4.3), which should take into account whether the returned hit(s) is attributed to the correct  
143 species and whether the hit(s) is a close enough match for the taxon in question, appropriate level  
144 assignment (4.4) and reporting results from GenBank (4.5).

145 **4.1** Prior to performing a BLAST search, evidentiary query sequences:

146 **4.1.1** Shall be prepared by removing non-template flanking regions (*e.g.* primer);

147 **4.1.2** Shall meet sequence quality criteria as defined by the laboratory. Thus, laboratories are  
148 responsible for having these criteria clearly defined and ensuring their analysts follow these  
149 recommendations.

150 **4.1.3** Shall be examined to ensure it does not contain premature stop codons (*e.g.* by translation).

151 **4.2** To ensure that a hit(s) on which conclusions are based are of high quality, an initial assessment of  
152 the BLAST results:

153 **4.2.1** Shall ensure the hit(s) belongs to the expected broader taxonomic group (*e.g.*, macerated  
154 plant tissue returns matches to sequences from the plant kingdom, not the bacterial  
155 kingdom).

156 NOTE: In situations involving a complete unknown, it may not be possible to complete this assessment.

157 **4.2.2** Shall ensure that any hit(s) that is an anomaly among the returned results is not used.  
158 This would be indicated by being the only representative of its species interleaved  
159 among many in a different taxonomic group. This could be an indication of human error  
160 in sequence labeling during sequence preparation prior to GenBank upload.



161 **4.2.3** Shall ensure the hit(s) does not originate from an environmental sample (e.g., bulk soil  
162 extraction, bacterial swab) or low copy sample.

163 NOTE: The original publication can often be consulted to determine the source of the sequence. In some  
164 instances, this determination may not be possible.

165 **4.2.4** Should include a review for descriptors or characteristics that indicate the sequence was  
166 not reviewed prior to uploading in GenBank.

167 NOTE: Sequences that have not been reviewed for quality may include descriptors such as “NGS”, “MPS”,  
168 “EST”, “shotgun”, “library”, and “WGS”; these may have been batch uploaded directly from the sequencing  
169 platform. Unedited sequences may also have a higher number of “Ns” or degenerate bases at the ends, or  
170 contain non-template flanking (e.g., primer, adapter) sequences.

171 **4.2.5** Should include a review for ambiguous bases.

172 NOTE: Ambiguous bases should be treated with caution, as they can indicate poor-quality sequence, but  
173 they can also indicate heteroplasmic sites within a high-quality sequence.

174 **4.2.6** Shall ensure the hit(s) from a protein coding region does not contain premature stop  
175 codons.

176 **4.3** Any hit(s) on which conclusions are based shall be evaluated to determine if the returned  
177 sequence is attributed to the correct species based on the criteria listed below. This section is to  
178 determine if returned sequences are appropriate for interpretations as outlined in Section 4.4.  
179 These criteria confer either strong or moderate support to the attribution. If the returned  
180 sequence(s) does not meet at least the moderate criteria, they shall not be used for taxonomic  
181 assignment to the species level. :

182 **4.3.1** Strong criteria (not all of these criteria have to be met, see section 4.5 for more information  
183 about how to evaluate relevant criteria):

184 a) Sequence(s) is derived from a voucher specimen that bears a unique identifier.

185 b) Sequence(s), when downloaded, aligned with sequences from closely-related species  
186 and used to construct a phylogram, results in a species-level topology concordant with  
187 expectations from the peer-reviewed literature.

188 c) Sequence(s) is from a study published in a peer-reviewed journal; the study addresses  
189 the phylogeny or taxonomy of the taxon of interest and the publication or accompanying  
190 metadata makes it clear that the source specimen(s) was morphologically identified by a  
191 taxonomic expert.

192 d) Sequence(s) is part of a population genetic study for the given species published in a  
193 peer-reviewed journal.

194 NOTE: Typically a population genetic study characterizes numerous individuals from the studied  
195 species in order to explore intraspecific variation (sample sizes will vary based on genetic  
196 variability and rareness of the species in question; published studies will have sample sizes that  
197 are appropriate for the species in question). The individuals may either be from the same  
198 geographic region, or from distinct populations within the known distributional range.

199 **4.3.2** Moderate criteria (not all of these criteria have to be met, see section 4.5 for more  
200 information about how to evaluate relevant criteria):

201 a) Sequence(s) is from a study published in a peer-reviewed journal; the study includes  
202 additional data establishing species identity (*e.g.*, morphological evidence, museum  
203 specimen), but it is not clear that the source specimen was a voucher (4.3.1a) or was  
204 morphologically identified by a taxonomic expert (4.3.1c).

205 b) Sequence(s) is from a phylogenetic study in a peer-reviewed journal; the study  
206 addresses phylogeny or taxonomy of the taxon of interest and:

207 i. includes most or all members of the genus in question, and

208 ii. the locus shows resolution at the species level (see 4.4.2).

209 c) Sequence(s) is one of multiple identical or near-identical sequences for the same  
210 locus and species from different submitters or geographic locations.

211 d) Sequence(s) is not from a peer-reviewed study on the taxon of interest, but is  
212 accompanied by additional metadata concerning the source individual (*e.g.*, location  
213 life history stage, name of collector, name of taxonomic expert who rendered the  
214 source individual's identification).

215 **4.4** The following should be evaluated to determine the appropriate level for taxonomic  
216 assignment:

217 **4.4.1** Whether all likely candidate species in the taxonomic group in question are  
218 represented amongst the returned hit(s).

219 NOTE: Complete taxon sampling is ideal, but often not feasible. If relevant taxa are missing,  
220 other loci or additional reference material should be considered. Species that are distantly  
221 related based on published phylogenies or those that do not occur in the geographic area of  
222 interest may be exempted from the comparison if sequences are not available. See section 4.5.2 in  
223 ASB 019 and section 3.5 in ASB 029.

224 NOTE: Peer-reviewed literature or internal validation for the species/marker of interest  
225 provides the foundation for evaluating whether hits are appropriate and comprehensive  
226 enough to provide accurate interpretation for reporting.

227 **4.4.2** Whether the interspecific distance for the taxonomic group of interest at the surveyed locus  
228 is greater than intraspecific distance.

229 NOTE: If inter- and intraspecific distances are similar, one should consider using a different  
230 locus or limiting identification to a higher taxonomic level.

231 **4.5** Reporting from BLAST results

232 **4.5.1** It is appropriate to report to the species level when all of these criteria are met:

233 a) The evidentiary sequence(s) has been prepared as outlined in 4.1,

- 234 b) The hit(s) on which conclusions are to be based:
- 235 i. meets the quality criteria as defined in 4.2;
- 236 ii. meets at least two strong support criteria (as defined in 4.3.1), or at least one  
237 strong and one moderate (as defined in 4.3.2) support criteria;
- 238 iii. has been evaluated against the criteria defined in 4.4;
- 239 iv. and when aligned to the evidentiary query sequence, shows 99–100% identity  
240 (inclusive).

241 NOTE: 99% is a conservative threshold, to be applied in instances where no other  
242 information is available for the target taxon. For most species, intraspecific distance will be  
243 greater than 1%; in cases where additional information (*e.g.*, other loci, taxonomies based on  
244 morphological features) indicates species are well-separated, identities lower than 99% may still  
245 warrant a species level identification.

246 NOTE: By default, BLAST results are sorted by E-value, which preferentially weights matches  
247 with higher query coverage, and max-score, based on sequence similarities. This can result in  
248 shorter sequences with higher percent identity being displayed after longer sequences with  
249 lower percent identity. The list may be sorted by the identity value to reveal the highest-  
250 similarity matches. It is critical to consider both the percent identity and the length of the match  
251 when evaluating BLAST results.

252 **4.5.2** It is appropriate to report to a higher taxonomic level when all of these criteria are met:

- 253 a) The evidentiary sequence(s) has been prepared as outlined in 4.1,
- 254 b) The hit(s) meets the quality criteria as defined in 4.2,
- 255 c) The hit(s) has been evaluated against the criteria defined in 4.4,
- 256 d) The hit(s) does not meet the support criteria given in 4.5.1(b)ii, but is from a  
257 peer-reviewed publication and:
- 258 i. The most similar sequences returned by a query are <99% identical and  
259 there is little definitive information on interspecific distance.

260 OR

- 261 ii. All top hits represent a single taxonomic level (*i.e.*, genus, family, order),  
262 but there is a discrepancy at a lower taxonomic level (*e.g.*, hits represent  
263 different species, but they all belong to a single genus).  
264

265

## Annex A (informative)

266 This is not meant to be an all-inclusive list as the group recognizes other publications on this subject  
267 may exist. At the time this standard was drafted, these were the publications available for reference.  
268 Additionally, any mention of a particular software tool or vendor as part of this bibliography is  
269 purely incidental, and any inclusion does not imply endorsement.

270

## Bibliography

- 271 1] Altschul SF. (2014). "BLAST Algorithm." In: *eLS, John Wiley & Sons, Ltd (Ed.)*. doi:  
272 10.1002/9780470015902.a0005253.pub2.
- 273 2] ANSI/ASB Standard 019, Wildlife Forensics General Standards, First Edition, 2019.
- 274 3] ANSI/ASB Standard 029, Report Writing in Wildlife Forensics: Morphology and Genetics,  
275 First Edition, 2019.
- 276 4] ANSI/ASB Standard 048, Wildlife Forensic DNA Standard Procedures, First Edition, 2019.
- 277 5] Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. (2013).  
278 "GenBank." *Nucleic Acids Research* 41(D1):D36-42. Available from:  
279 <https://www.ncbi.nlm.nih.gov/genbank/>.
- 280 6] BLAST® Command Line Applications User Manual [Internet]. Bethesda (MD): National  
281 Center for Biotechnology Information (US); 2008-. Available from:  
282 <https://www.ncbi.nlm.nih.gov/books/NBK279690/>.
- 283 7] Brown TA. *Genomes*. 2nd edition. Oxford: Wiley-Liss; 2002. Chapter 16, Molecular  
284 Phylogenetics. Available from: <https://www.ncbi.nlm.nih.gov/books/NBK21122/>.
- 285 8] International Organization for Standardization. (2017). "ISO/IEC 17025:2005 General  
286 Requirements for the Competence of Testing and Calibration Laboratories." 28 pp.
- 287 9] Lee TRC, Anderson SJ, Tran-Nguyen LTT, Sallam N, Le Ru BP, Conlong D, Powell K, Ward A,  
288 Mitchell A. 2019. Towards a global DNA barcode reference library for quarantine  
289 identifications of lepidopteran stemborers, with an emphasis on sugarcane pests. *Scientific*  
290 *Reports* 9: 7039. Doi: <https://doi.org/10.1038/s41598-019-42995-0>.
- 291 10] Lorenz JG, Jackson WE, Beck JC, Hanner R. (2005). "The problems and promise of DNA  
292 barcodes for species diagnosis of primate biomaterials." *Philosophical Transactions of the*  
293 *Royal Society B* 360, 1869–1877.
- 294 11] Madden T. (2013). "The BLAST Sequence Analysis Tool." In: *The NCBI Handbook, 2nd ed.*  
295 Bethesda, MD. Available from <https://www.ncbi.nlm.nih.gov/books/NBK153387/>.