

1 **SeqCode, a nomenclatural code for prokaryotes described from sequence data**

2
3 Prepared as a Brief Communication for Nature Microbiology

4
5 Brian P. Hedlund¹, Maria Chuvochina², Philip Hugenholtz², Konstantinos T. Konstantinidis³,
6 Alison E. Murray⁴, Marike Palmer¹, Donovan H. Parks², Alexander J. Probst⁵, Anna-Louise
7 Reysenbach⁶, Luis M. Rodriguez-R⁷, Ramon Rossello-Mora⁸, Iain C. Sutcliffe⁹, Stephanus N.
8 Venter¹⁰ and William B. Whitman^{11*}

9
10 ¹ School of Life Sciences, University of Nevada, Las Vegas, NV, USA

11 ² The University of Queensland, School of Chemistry and Molecular Biosciences, Australian
12 Centre for Ecogenomics, Brisbane, Australia

13 ³ School of Civil and Environmental Engineering, Georgia Tech, Atlanta, GA, USA

14 ⁴ Division of Earth and Ecosystem Sciences, Desert Research Institute, Reno, NV, USA

15 ⁵ Department of Chemistry, Environmental Microbiology and Biotechnology (EMB), Group for
16 Aquatic Microbial Ecology and Centre of Water and Environmental Research
17 (ZWU), University of Duisburg-Essen, Essen, Germany.

18 ⁶ Biology Department, Portland State University, Portland, OR, USA

19 ⁷ Department of Microbiology and Digital Science Center (DiSC), University of Innsbruck,
20 Innsbruck, Austria

21 ⁸ Marine Microbiology Group, Department of Animal and Microbial Diversity, Mediterranean
22 Institute of Advanced Studies (CSIC-UIB), Esporles, Illes Balears, Spain

23 ⁹ Faculty of Health & Life Sciences, Northumbria University, Newcastle upon Tyne, UK

24 ¹⁰ Department of Biochemistry, Genetics and Microbiology, University of Pretoria, Pretoria,
25 South Africa

26 ¹¹ Department of Microbiology, University of Georgia, Athens, GA, USA

27
28 For correspondence: William B. Whitman, whitman@uga.edu

29

30 **Abstract**

31 Most prokaryotes are not available as pure cultures and therefore ineligible for naming under the
32 International Code of Nomenclature of Prokaryotes. Here we summarize the development of the
33 SeqCode, a code of nomenclature under which genome sequences can serve as nomenclatural
34 types. The SeqCode operates through self-registration (<https://seqco.de/>), provides a reproducible
35 and objective framework for all prokaryotes, regardless of cultivability, and facilitates
36 communication across microbiological disciplines.

37 **Manuscript body**

38 It is widely recognized that the requirement of the International Code of Nomenclature of
39 Prokaryotes (ICNP) for deposition of axenic and viable cultures as nomenclatural types has
40 hindered the development of a nomenclature for uncultured and fastidious cultured prokaryotes
41 (Archaea and Bacteria) and thus effective communication of microbial diversity (Konstantinidis
42 et al., 2017; Murray et al., 2020). For example, as-yet-uncultivated taxa account for ~85% of the
43 phylogenetic diversity of prokaryotes (Nayfach et al., 2021), and named prokaryotes account for
44 only <0.2% of total species (Sutcliffe et al., 2021). By excluding the uncultured majority, a
45 substantial portion of the tree of life is relegated to poorly ordered, ambiguous, and often
46 synonymous names or alphanumeric codes, the latter of which have limited mnemonic value
47 (Miller 1956).

48 To address this problem, Murray et al. (2020) proposed two paths, which were endorsed by 121
49 authors and signatories from 22 countries and six continents (Murray et al., 2020). ‘Plan A’ was
50 based on proposals by Whitman (2015) that DNA sequences could serve as nomenclatural types
51 and be incorporated into the existing ICNP infrastructure. However, the International Committee
52 on Systematics of Prokaryotes (ICSP) rejected Whitman’s proposal (Sutcliffe et al., 2020), thus
53 triggering “Plan B”, which called for a new code of nomenclature (Murray et al., 2020). To
54 further engage the community in the implementation of “Plan B”, we organized a series of online
55 workshops ([https://www.isme-](https://www.isme-microbes.org/reports-sponsored-events)
56 [microbes.org/reports-sponsored-events](https://www.isme-microbes.org/reports-sponsored-events)) that garnered 848
57 registrants from a broad range of microbiology disciplines and 42 countries. Ninety percent of
58 participants reported that they would use a new code that accepts DNA sequences as types
59 ([https://www.isme-](https://www.isme-microbes.org/sites/default/files/reports/Path_forward_Naming_Uncultivated.pdf)
60 [microbes.org/sites/default/files/reports/Path_forward_Naming_Uncultivated.pdf](https://www.isme-microbes.org/sites/default/files/reports/Path_forward_Naming_Uncultivated.pdf)). Given strong
61 participation and near-unanimous support, we acted on a variety of community recommendations
62 (Table S1) to complete the SeqCode (formally The International Code of Nomenclature of
63 Prokaryotes Described from Sequence Data; see Additional Information) and made progress on
64 systems to implement it.

64 The SeqCode uses genome sequence data as common currency for typification of both cultivated
65 and uncultivated microorganisms and follows the tenets of the ICNP by observing similar rules
66 of priority. In essence, these rules state that the earliest validly published name for a taxon in a

67 particular position is the correct name, observing historical precedent and stabilizing
68 nomenclature. The SeqCode also recognizes the priority of ICNP names provided they do not
69 violate the priority of SeqCode names, thus minimizing divergence between the systems.
70 Taxonomic names will be captured in the SeqCode Registry, a simple self-registration portal
71 through which names and nomenclatural types (e.g., genome sequences for species) are
72 registered, validated, and linked to metadata. In the best-case scenario, data will be entered and
73 reviewed prior to publication, allowing automated checks and curators to guide users through the
74 naming process. Following peer review and publication of the manuscript describing the taxa,
75 the manuscript Digital Object Identifier (DOI) is entered into the Registry, completing the valid
76 publication of the name/s (Figure 1, Path 1). However, the SeqCode also enables registration of
77 previously published names, such as *Candidatus* names that conform to its rules. In that case, the
78 *Candidatus* designation could be dropped, and the names given priority under the SeqCode
79 (Figure 1, Path 2; see Additional Information). While the SeqCode itself is necessarily
80 comprehensive, we have also developed resources to guide the community, including a glossary
81 and examples (see Supplementary Information). Table 1 summarizes recommended minimal
82 standards for sequences and reporting requirements. We endorse high quality standards for use of
83 the SeqCode but expect standards to evolve to keep pace with community feedback and
84 methodological improvements.

85 Potential users may ask: (i) What is the difference between *Candidatus* status and valid
86 publication under the SeqCode? In reply, *Candidatus* is a provisional status lacking priority and
87 standing and is relegated to a non-legislative appendix of the ICNP. *Candidatus* status was
88 developed for organisms for which “more than a mere nucleic acid sequence is available”. Since
89 its inception, visualization of the taxon in a natural sample has been recommended (Murray and
90 Stackebrandt 1995; Parker et al., 2019), but this is rarely implemented. It has been argued that
91 *Candidatus* names should be granted priority under the ICNP (Whitman et al., 2019); however,
92 this proposal was also rejected (Sutcliffe et al., 2020). As a result, many *Candidatus* names may
93 prove to be ephemeral. (ii) What are the consequences for taxonomic names that are published in
94 primary literature but not validly published under the SeqCode? Although the community is free
95 to publish taxonomic names that do not comply with codes of nomenclature, we argue that codes
96 of nomenclature and taxonomic frameworks serve the greater community by promoting
97 objectivity, best practices, communication, and data interoperability. However, the unique
98 restrictions of the ICNP regarding viable and accessible type strains have alienated many
99 microbiologists and engendered a sense of normalcy in publishing names outside of the
100 regulation of the ICNP. The SeqCode addresses this problem by providing an efficient and user-
101 friendly resource that serves the common interests of the wider research community. The
102 SeqCode embraces Findability, Accessibility, Interoperability, and Reusability (FAIR)
103 principles, and the Registry was developed with interoperable data structures to promote sharing
104 SeqCode names across global biodiversity inventories within microbiology and the broader
105 biology research communities (e.g., NCBI (Schoch et al., 2020), GTDB (Parks et al., 2018),

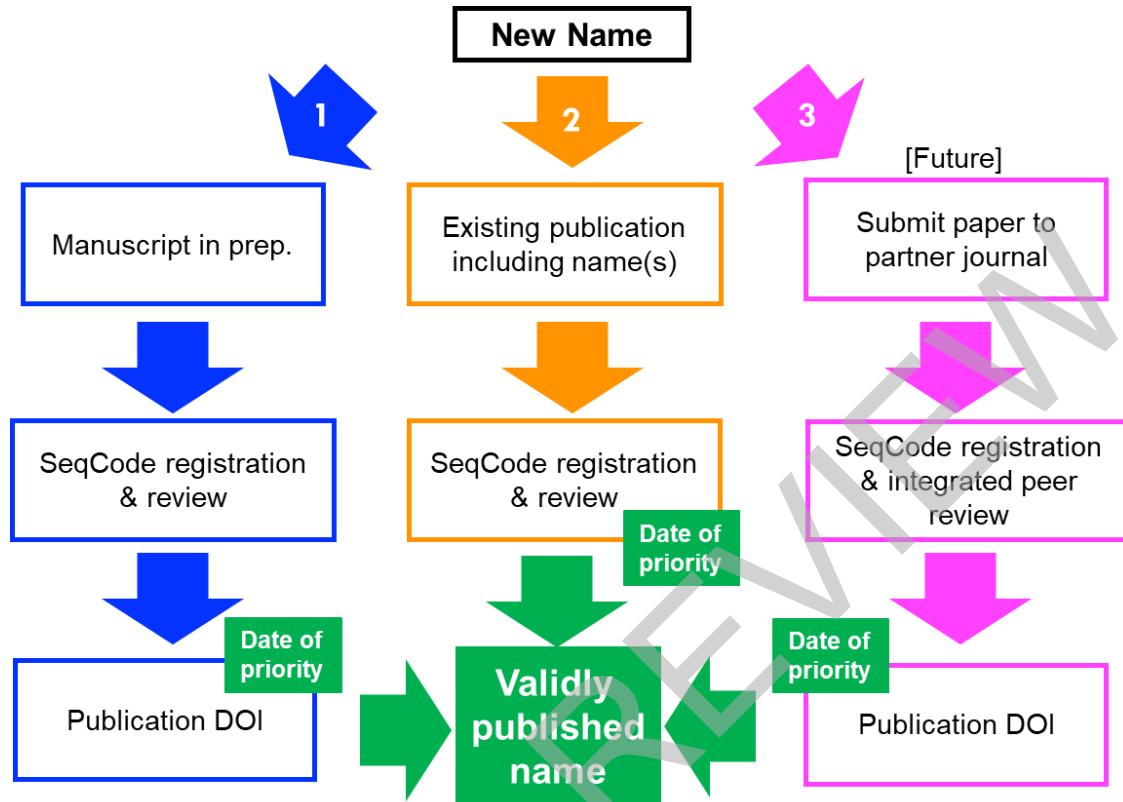
106 MiGA (Rodriguez-R et al., 2016), LPSN (Parte et al., 2020), Catalogue of Life (Roskov et al.,
107 2019), Global Biodiversity Information Facility (GBIF 2020)).

108 In closing, we emphasize a few important points. First, the SeqCode is not intended to
109 discourage cultivation. Cultivation of mixed or pure cultures enables testing properties predicted
110 from genomes under controlled conditions. Furthermore, investigators are strongly encouraged to
111 deposit strains to culture collections to improve strain stability and availability, enable
112 assessment of reproducibility of phenotypic traits, provide resources for biochemistry and
113 biotechnology, and promote international cooperation. Second, like all other codes of
114 nomenclature, the SeqCode does not provide rules or recommendations on the delineation of
115 taxa. Existing and improving approaches and data structures are available for that purpose (e.g.,
116 Parks et al., 2020; Rodriguez-R et al., 2018), and proposals for novel taxa must be settled
117 through peer review. Finally, this is the first version of the SeqCode and we hope that it will
118 evolve as the community engages in further developing the system. Because of our desire to
119 serve the broad microbiology research community, we will engage the community to gather
120 feedback and develop bylaws for SeqCode administration. This code is driven by bottom-up
121 desires to improve communication across the microbial sciences. Thus, we view this ‘SeqCode
122 v1.0’ as a necessary first step toward a unified system of nomenclature to communicate the full
123 diversity of prokaryotes, and we will cooperate with the community toward building this vision.

124 **Acknowledgements**

125 Large portions of the text of the SeqCode were derived from the ICNP, and the editors gratefully
126 acknowledge the many authors who contributed to that code. Funding was provided by the US
127 National Science Foundation (DEB 1841658, DEB 1557042, and EAR 1516680), the US
128 National Institute of General Medical Sciences (GM103440) from the National Institutes of
129 Health, the Spanish Ministry of Science, Innovation and Universities (PGC2018-096956-B-C41),
130 the Australian Research Council (FL150100038), the Deutsche Forschungsgemeinschaft (DFG,
131 German Research Foundation, SFB 1439/1 2021 – 426547801) also supported with European
132 Regional Development Funds (FEDER), and the International Society for Microbial Ecology
133 (ISME). We also thank all participants in the SeqCode workshops, especially guest speakers who
134 graciously shared their expertise: Jongsik Chun, Nicole Dubilier, Emiley Eloie-Fadrosch, Chris
135 Lane, Juncai Ma, Edward Moore, Aharon Oren, Jörg Overmann, Susanne Renner, Vincent
136 Robert, Conrad Schoch, Scott Tinghe, Linhuan Wu, and Arvind Varsani.

Validation of a name under the SeqCode



137
138 **Figure 1. Validation process of a name under the SeqCode.** Currently, two mechanisms exist,
139 with a third possible in the future. The recommended mechanism (left arrow, Path 1) involves
140 draft registration of the name and metadata into the SeqCode Registry prior to publication.
141 Automated data quality and name synonymy checks in conjunction with curator review will lead
142 to provisional acceptance of proposals that comply with SeqCode rules. Completion of the
143 registration process requires the DOI of the effective publication. Once the proposal is accepted
144 and the DOI entered, the registration is complete, marking the time and date of priority. The
145 second (middle arrow, Path 2) is for names that are already published, such as *Candidatus*
146 names. It requires draft registration of the name and metadata into the SeqCode Registry.
147 SeqCode curators review compliance with the SeqCode rules before accepting the proposal.
148 Acceptance of the proposal completes registration and marks the time and date of priority. At
149 that point, the *Candidatus* designation can be removed. The third mechanism could be developed
150 in partnership with one or more journals in the future (right arrow, Path 3) and would involve
151 simultaneous peer review and SeqCode Registry curator review as an integrated path to the
152 validation of proposed names. Issue of the DOI of the accepted paper marks the time and date of
153 priority.
154

155 **Table 1.** Data quality and reporting requirements and recommendations for an isolate genome,
 156 MAG, or SAG to serve as the nomenclatural type for a species named under the SeqCode.
 157 Requirements will be checked as part of the validation process on the SeqCode Registry.
 158 Recommendations are suggested best practices to guide authors and peer reviewers to ensure
 159 high quality data supporting species to be named. See Supplementary Information for examples.

Information	Requirements	Recommendations
Included in publication proposing new species names under SeqCode^a		
Name	Required for all names	1. Etymologies for all proposed names are recommended. 2. Names with mnemonic cues are recommended.
Interpretation of biological properties	None	Indicate inferred or demonstrated physiological traits and ecological information, such as habitat in the manuscript body and/or protologue.
Designated genome	None	1. Indicate access to genomic assembly (e.g., INSDC accession). 2. Indicate access to raw data (e.g., SRA accession). 3. Demonstrate compliance with GSC standards for isolate genomes (Field et al., 2008) and high-quality SAGs and MAGs (Bowers et al., 2017). 4. Include as much metadata as possible in the publication (see Field et al., 2008).
Evidence of the species, taxonomic rank, and position	None	1. Demonstrate the uniqueness of the species with respect to existing named species and justify the taxonomic rank and position (e.g., Jain et al., 2018, Karthikeyan et al., 2019; Parks et al., 2020; Rodriguez-R et al., 2018). 2. For MAGs and SAGs, compare multiple high-quality genomes representing the species in more than one sample (e.g., Supplemental Information). ^b
Data quality^c and availability necessary for completion of SeqCode Registry		
Type genome assembly quality	1. >90% complete and <5% contaminated; 16S and 23S rRNA genes >75% complete (modified from Bowers et al., 2017). 2. Isolate genome read coverage $\geq 50x$ (Field et al., 2008).	1. >80% of tRNAs present (modified from Bowers et al., 2017). 2. High genome integrity (contig # <100; N50 >25 kb; max. contig >10 kb). 3. MAG/SAG read coverage $\geq 10x$.
INSDC data availability	1. Assembly available in INSDC database. 2. Raw data available in INSDC databases (e.g., Sequence Read Archive) ^d .	1. Data submission using MixS Checklists in INSDC databases (https://gensc.org/mixs/). 2. Include as much metadata as possible in INSDC.
SeqCode Registry	Type genome assembly and raw data INSDC	Provide as much contextual data as possible to facilitate downstream genome comparisons with

accession numbers, respect to provenance.
taxon name, etymology,
rank.

- 160 a. There are purposefully few requirements for the effective publication to accommodate existing and
161 future publications that don't adhere to all recommendations. Critical data will be captured on the
162 SeqCode Registry (Figure 1).
- 163 b. Comparison of multiple high-quality genomic assemblies from multiple samples can support the non-
164 chimeric nature of MAGs and provide confidence of the assembly for both MAGs and SAGs.
- 165 c. Data quality can be assessed by automated pipelines or other approaches. Exceptions for lower data
166 quality should be justified by authors in the effective publication.
- 167 d. Not required for names effectively published before January 1, 2023, to allow for existing published
168 names (e.g., existing *Candidatus* names) and names currently undergoing peer review to be validated
169 under the SeqCode.
- 170

UNDER REVIEW

171 **References**

- 172
- 173 Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, et al.
174 Minimum information about a single amplified genome (MISAG) and a metagenome-assembled
175 genome (MIMAG) of bacteria and archaea. *Nat Biotechnol.* 2017;35:725-31.
176
- 177 GBIF: The Global Biodiversity Information Facility (2020) What is GBIF? Available from
178 <https://www.gbif.org/what-is-gbif>
179
- 180 Field D, Garrity G, Gray T, Morrison N, Selengut J, Sterk P, et al. The minimum information
181 about a genome sequence (MIGS) specification. *Nat Biotechnol.* 2008 May; 26:541–547.
182
- 183 Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI
184 analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun.* 2018;9:1–
185 8.
186
- 187 Konstantinidis KT, Rosselló-Móra R, Amann R. Uncultivated microbes in need of their own
188 taxonomy. *ISME J.* 2017;11:2399–406.
189
- 190 Miller G. The magical number seven, plus or minus two: some limits on our capacity for
191 processing information. *Psychol Rev* 1956;63(2):81–97.
192
- 193 Murray RGE, Stackebrandt E. Taxonomic note: implementation of the provisional status
194 candidatus for incompletely described prokaryotes. *Int. J. Syst. Bacteriol.* 1995;45:186–7.
195
- 196 Murray AE, Freudenstein J, Gribaldo S, Hatzenpichler R, Hugenholtz P, Kämpfer P, et al.
197 Roadmap for naming uncultivated archaea and bacteria. *Nat Microbiol.* 2020;5:987-94.
198
- 199 Nayfach S, Roux S, Seshadri R, Udwaray D, Varghese N, Schulz F, Wu D, et al. A genomic
200 catalog of Earth's microbiomes. *Nat Biotechnol* 2021;39:499-509.
201
- 202 Parker CT, Tindall BJ, Garrity, GM. International code of nomenclature of prokaryotes. *Int J*
203 *Syst Evol Microbiol*; 2019; 69:S1-S111.
204
- 205 Parks DH, Chuvochina M, Chaumeil P-A, Rinke C, Mussig AJ, Hugenholtz P. A complete
206 domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol.* 2020;38:1079–86.
207
- 208 Parte AC, Sardà Carbasse J, Meier-Kolthoff JP, Reimer LC, Göker M. List of Prokaryotic names
209 with Standing in Nomenclature (LPSN) moves to the DSMZ. *Int J Syst Evol Microbiol*; 2020;70:
210 5607-12.
211
- 212 Rodriguez-R LM, Gunturu S, Harvey WT, Rosselló-Mora R, Tiedje JM, Cole JR, Konstantinidis
213 KT. The Microbial Genomes Atlas (MiGA) webserver: taxonomic and gene diversity analysis of
214 Archaea and Bacteria at the whole genome level. *Nucleic Acids Res.* 2018;46:W282–8.
215

216 Roskov Y, Ower G, Orrell T, Nicolson D, Bailly N, Kirk PM, et al. Species 2000 & ITIS
217 Catalogue of Life, 25th March 2019. Digital resource at www.catalogueoflife.org/col.
218 2021;Species 2000: Naturalis, Leiden, the Netherlands. ISSN 2405-8858.
219
220 Schoch CL, Ciuffo S, Domrachev M, Hotton CL, Kannan S, Khovanskaya R, Leipe D, McVeigh
221 R, O'Neill K, Robbertse B, Sharma S, Soussov V, Sullivan JP, Sun L, Turner S, Karsch-Mizrachi
222 I. NCBI Taxonomy: a comprehensive update on curation, resources and tools. Database. 2020
223 baaa062.
224
225 Sutcliffe IC, Rosselló-Mora R, Trujillo M. Addressing the sublime scale of the microbial world:
226 reconciling an appreciation of microbial diversity with the need to describe species. *New*
227 *Microbes New Infect.* 2021;43:100931.
228
229 Sutcliffe IC, Dijkshoorn L, Whitman WB. Minutes of the International Committee on
230 Systematics of Prokaryotes online discussion on the proposed use of gene sequences as type for
231 naming of prokaryotes, and outcome of vote. *Int. J. Syst. Evol. Microbiol.* 2020;70:4416-7.
232
233
234 Whitman WB. Genome sequences as type material for taxonomic descriptions. *System. Appl.*
235 *Microbiol.* 2015;38:217-222.
236
237 Whitman WB, Sutcliffe I, Rosselló-Mora R. Proposal for changes in the International Code of
238 Nomenclature of Prokaryotes: granting priority to Candidatus names. *Int J Syst Evol Microbiol.*
239 2019;69(7):2174-2175.