

1 **The “Minimum Information about an ENvironmental Sequence” (MIENS)**
2 **specification**

3
4 Pelin Yilmaz^{1,2}, Renzo Kottmann¹, Dawn Field³, Rob Knight^{4,5}, James R. Cole^{6,7}, Linda
5 Amaral-Zettler⁸, Jack A. Gilbert^{9,10,11}, Ilene Karsch-Mizrachi¹², Anjanette Johnston¹²,
6 Guy Cochrane¹³, Robert Vaughan¹³, Christopher Hunter¹³, Joonhong Park¹⁴, Norman
7 Morrison^{15,16}, Phillipe Rocca-Serra^{13,17}, Peter Sterk³, Mani Arumugam¹⁸, Laura
8 Baumgartner¹⁹, Bruce W. Birren²⁰, Martin J. Blaser²¹, Vivien Bonazzi²², Tim Booth³,
9 Peer Bork¹⁸, Frederic D. Bushman²³, Pier Luigi Buttigieg^{1,2}, Patrick S. G. Chain^{7,24,25},
10 Emily Charlson²³, Elizabeth K. Costello⁴, Heather Huot-Creasy²⁶, Peter Dawyndt²⁷, Todd
11 DeSantis²⁸, Noah Fierer²⁹, Jed Fuhrman³⁰, Rachel E. Gallery³¹, Dirk Gevers²⁰, Richard A.
12 Gibbs^{32,33}, Michelle Gwinn Giglio²⁶, Inigo San Gil³⁴, Antonio Gonzalez³⁵, Jeffrey I.
13 Gordon³⁶, Robert Guralnick²⁹, Wolfgang Hankeln^{1,2}, Sarah Highlander^{32,37}, Philip
14 Hugenholtz²⁴, Janet Jansson³⁸, Scott T. Kelley³⁹, Jerry Kennedy⁴, Dan Knights³⁵, Omry
15 Koren⁴⁰, Justin Kuczynski¹⁹, Nikos Kyrpides²⁴, Robert Larsen⁴, Christian L. Lauber⁴¹,
16 Teresa Legg²⁹, Ruth E. Ley⁴⁰, Catherine A. Lozupone⁴, Wolfgang Ludwig⁴², Donna
17 Lyons⁴¹, Eamonn Maguire^{13,17}, Barbara A. Methé⁴³, Folker Meyer¹⁰, Sara Nakielny⁴,
18 Karen E. Nelson⁴³, Diana Nemergut⁴⁴, Lindsay K. Neubold³, Josh D. Neufeld⁴⁵, Anna E.
19 Oliver³, Norman R. Pace¹⁹, Giriprakash Palanisamy⁴⁶, Jörg Peplies⁴⁷, Jane Peterson²²,
20 Joseph Petrosino^{32,37}, Lita Proctor⁴⁸, Elmar Pruesse^{1,2}, Christian Quast¹, Jeroen Raes⁴⁹,
21 Sujeevan Ratnasingham⁵⁰, Jacques Ravel²⁶, David A. Relman^{51,52}, Susanna Assunta-
22 Sansone^{13,17}, Patrick D. Schloss⁵³, Lynn Schriml²⁶, Rohini Sinha²³, Erica Sodergren⁵⁴,
23 Aymé Spor⁴⁰, Jesse Stombaugh⁴, James M. Tiedje⁷, Doyle V. Ward²⁰, George M.

24 Weinstock⁵⁴, Doug Wendel⁴, Owen White²⁶, Andrew Whitely³, Andreas Wilke¹⁰,
25 Jennifer R. Wortman²⁶, Frank Oliver Glöckner^{1,2}

26

27

28 1 Microbial Genomics and Bioinformatics Group, Max Planck Institute for Marine
29 Microbiology, D-28359 Bremen, Germany

30 2 Jacobs University Bremen gGmbH, D-28759 Bremen, Germany

31 3 NERC Centre for Ecology and Hydrology, Maclean Building, Benson Lane,
32 Crowmarsh Gifford, Wallingford, Oxfordshire, OX10 8BB, UK

33 4 Department of Chemistry and Biochemistry, University of Colorado, Boulder, CO
34 80309, USA

35 5 Howard Hughes Medical Institute, USA

36 6 Ribosomal Database Project, Michigan State University, 2225A Biomedical and
37 Physical Sciences Building, East Lansing, Michigan 48824-4320, USA

38 7 Center for Microbial Ecology, Michigan State University, 540 Plant and Soil Sciences
39 Building, East Lansing, Michigan 48824-1325, USA

40 8 The Josephine Bay Paul Center for Comparative Molecular Biology and Evolution,
41 Marine Biological Laboratory, Woods Hole, Massachusetts, USA

42 9 Plymouth Marine Laboratory, Prospect Place, Plymouth, UK

43 10 Mathematics and Computer Science Division, Argonne National Laboratory,
44 Argonne, IL 60439, USA

45 11 Dept of Ecology and Evolution, University of Chicago, Chicago, IL 60637, USA

46 12 National Center for Biotechnology Information (NCBI), National Library of
47 Medicine, National Institutes of Health, 8600 Rockville Pike, Bethesda, Maryland 20894,

48 USA

49 13 European Molecular Biology Laboratory (EMBL) Outstation, European
50 Bioinformatics Institute (EBI), Wellcome Trust Genome Campus, Hinxton, Cambridge
51 CB10 1SD, UK.

52 14 School of Civil and Environmental Engineering, Yonsei University, Seoul, 120-749,
53 Republic of Korea

54 15 NERC Environmental Bioinformatics Centre, Oxford Centre for Ecology and
55 Hydrology, Oxford OX1 3SR, UK

56 16 Department of Computer Science, University of Manchester, Oxford Rd., Manchester,
57 UK

58 17 Oxford e-Research Centre, University of Oxford, Oxford OX1 3QG, UK

59 18 Structural and Computational Biology Unit, European Molecular Biology Laboratory,
60 Meyerhofstr. 1, D-69117 Heidelberg, Germany

61 19 Department of Molecular, Cellular and Developmental Biology, University of
62 Colorado, Boulder, CO 80309, USA

63 20 Broad Institute of Massachusetts Institute of Technology and Harvard University,
64 Cambridge, MA 02142

65 21 Department of Medicine and the Department of Microbiology, New York University
66 Langone Medical Center, New York, New York 10017, USA

67 22 National Human Genome Research Institute, National Institutes of Health, Bethesda,
68 Maryland 20892, USA

69 23 Department of Microbiology, University of Pennsylvania School of Medicine, 426A
70 Johnson Pavilion, 3610 Hamilton Walk, Philadelphia, PA 19104

71 24 DOE Joint Genome Institute, Walnut Creek, CA 94598, USA
72 25 Los Alamos National Laboratory, Bioscience Division, Los Alamos, New Mexico,
73 USA
74 26 Institute for Genome Sciences, University of Maryland School of Medicine,
75 Baltimore, MD 21201, USA
76 27 Department of Applied Mathematics and Computer Science, Ghent University,
77 Krijgslaan 281, 9000 Ghent, Belgium
78 28 Center for Environmental Biotechnology, Lawrence Berkeley National Laboratory,
79 Berkeley, CA, USA
80 29 Department of Ecology and Evolutionary Biology, University of Colorado, Boulder,
81 CO 80309, USA
82 30 Department of Biological Sciences, University of Southern California, Los Angeles,
83 CA, USA
84 31 National Ecological Observatory Network (NEON), Boulder, CO 80301, USA
85 32 Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX
86 33 Department of Molecular and Human Genetics, Baylor College of Medicine, Houston,
87 TX
88 34 Department of Biology, University of New Mexico, LTER Network Office, MSC03
89 2020, Albuquerque, NM 87131, USA
90 35 Department of Computer Science, University of Colorado, Boulder, CO 80309, USA
91 36 Center for Genome Sciences and Systems Biology, Washington University School of
92 Medicine, St. Louis, MO 63108, USA
93 37 Department of Molecular Virology and Microbiology, Baylor College of Medicine,

94 Houston, TX
95 38 Lawrence Berkeley National Laboratory, Earth Science Division, Berkeley, CA, USA
96 39 Department of Biology, San Diego State University, 5500 Campanile Drive, San
97 Diego , CA 92182-4614 USA
98 40 Department of Microbiology, Cornell University, Ithaca NY 14853, USA
99 41 Cooperative Institute for Research in Environmental Sciences, University of Colorado,
100 Boulder, USA
101 42 Lehrstuhl für Mikrobiologie, Technische Universität München, D-853530 Freising,
102 Germany
103 43 J. Craig Venter Institute, Rockville, Maryland, United States of America
104 44 Department of Environmental Sciences, University of Colorado, Boulder, CO 80309,
105 USA
106 45 Department of Biology, University of Waterloo, Ontario, N2L 3G1, Canada
107 46 Environmental Sciences Division, Oak Ridge National Laboratory, Mail Stop 6407
108 Oak Ridge, TN, USA
109 47 Ribocon GmbH, D-28359 Bremen, Germany
110 48 The National Science Foundation, 4201 Wilson Boulevard, Arlington, Virginia 22230,
111 USA
112 49 VIB - Vrije Universiteit Brussel, 1050 Brussels, Belgium
113 50 Canadian Centre for DNA Barcoding, Biodiversity Institute of Ontario, University of
114 Guelph, 50 Stone Road, Guelph, ON, Canada N1G 2W1
115 51 Departments of Microbiology and Immunology and of Medicine, Stanford University
116 School of Medicine, Stanford, CA 94305, USA

- 117 52 Veterans Affairs Palo Alto Health Care System, Palo Alto, CA 94304, USA
- 118 53 Department of Microbiology and Immunology, 5641 Medical Science Bldg. II, 1150
119 West Medical Center Dr., Ann Arbor, Michigan 48109-5620
- 120 54 The Genome Center, Department of Genetics, Washington University in St. Louis
121 School of Medicine, St. Louis, Missouri, USA
- 122
- 123

124 **Summary**

125 **We present the Genomic Standards Consortium’s (GSC) “Minimum Information**
126 **about an ENvironmental Sequence” (MIENS) standard for describing marker**
127 **genes. Adoption of MIENS will enhance our ability to analyze natural genetic**
128 **diversity across the Tree of Life as it is currently being documented by massive**
129 **DNA sequencing efforts from myriad ecosystems in our ever-changing biosphere.**

130 **Acronyms**

- 131 amoA: ammonia monooxygenase-alpha subunit
- 132 BOLI: Barcode of Life Initiative
- 133 CBOL: Consortium for the Barcode of Life
- 134 COI: cytochrome c oxidase I
- 135 DDBJ: DNA DataBank of Japan
- 136 DOE-JGI: Department of Energy Joint Genome Institute
- 137 DOI: Digital Object Identifier
- 138 DRA: DDBJ Sequence Read Archive
- 139 dsrAB: dissimilatory sulfite reductase
- 140 ENA: European Nucleotide Archive
- 141 EnvO: Environment Ontology
- 142 GAZ: Gazetteer
- 143 GCDML: Genomic Contextual Data Markup Language
- 144 GSC: Genomic Standards Consortium
- 145 gyrA: DNA gyrase (type II topoisomerase), subunit A
- 146 HSP70: 70 kilodalton heat shock protein
- 147 ICoMM: International Census of Marine Microbes
- 148 INSDC: International Nucleotide Sequence Database Collaboration
- 149 ISA: Investigation/Study/Assay Infrastructure
- 150 ISO: International Organization for Standardization
- 151 ITS: internal transcribed spacer region
- 152 LSU: large subunit

- 153 MICROBIS: The Microbial Oceanic Biogeographic Information System
- 154 MIENS: Minimum Information about an Environmental Sequence
- 155 MIGS/MIMS: Minimum Information about a Genome/Metagenome Sequence
- 156 MIRADA-LTERS: Microbial Inventory Research Across Diverse Aquatic Long Term
157 Ecological Research Sites
- 158 MLST: multi-locus sequence typing
- 159 NGS: next generation sequencing
- 160 nifH: dinitrogenase reductase
- 161 ntcA: nitrogen regulator gene
- 162 OBO: Open Biological and Biomedical Ontologies
- 163 phnA: phosphonoacetate hydrolase gene
- 164 phnJ: carbon-phosphorous lyase complex subunit
- 165 PMID: Pubmed ID
- 166 RDP: Ribosomal Database Project
- 167 recA: recombinase A subunit
- 168 rpoB: beta subunit of the bacterial RNA polymerase
- 169 rRNA: ribosomal RNA
- 170 SI: International System of Units
- 171 SRA: Sequence Read Archive
- 172 SSU: small subunit
- 173 URL: Uniform Resource Locator
- 174 WGS84: World Geodetic System 84
- 175 XML Schema: Extensible Markup Language Schema

176 **Big Data need Standards**

177 The term Big Data is increasingly being used to describe the vast capacity of high-
178 throughput experimental methodologies, especially next-generation sequencing, to
179 generate data ^{1,2}. Sharing and re-use of such data, and translating such data into
180 knowledge, requires widely-adopted standards that are best developed within the auspices
181 of international working groups ³. Here we describe a new standard, developed by a large
182 and diverse community of researchers, to describe one of the most abundant and useful
183 types of sequence data – that of marker gene data sets.

184

185 ***The wealth of marker gene data sets***

186 The adoption of phylogenetic marker genes as molecular proxies for tracking and
187 cataloguing the diversity of microorganisms has revolutionized the way we view the
188 biological world, and provided us with insights into how life has evolved and how
189 different organisms are genetically related to each other. In the 1970s, studies of small
190 subunit (SSU) ribosomal RNA (rRNA) genes from environmental samples led to the
191 discovery of the domain *Archaea* ⁴ and to the proposal for a three domain classification of
192 life ⁵. Following Darwin's insight that all life is related, SSU rRNA gene surveys allow
193 organisms from any communities, no matter how diverse, to be compared using the same
194 universal phylogenetic tree. This rRNA gene-based molecular approach to characterizing
195 natural communities of organisms provided, for the first time, culture-independent access
196 to the diversity and distribution of microorganisms '*in situ*'. As a result, we are now
197 acutely aware that the vast majority (90-99%) of microorganisms have evaded isolation
198 using existing cultivation methods ⁶⁻⁸.

199 Over the past three decades, the 16S rRNA, 18S rRNA and internal transcribed spacer
200 gene sequences (ITS) from *Bacteria*, *Archaea*, and microbial *Eukaryotes* have provided
201 deep insights into the topology of the tree of life ⁹⁻¹² and the composition of communities
202 of organisms that live in diverse environments, which range from deep sea hydrothermal
203 vents to ice sheets in the Arctic ¹³⁻²⁷.

204 Numerous other phylogenetic marker genes have also proven useful ²⁸: Currently, around
205 40 such phylogenetic marker genes are in wide use, representing well-conserved,
206 housekeeping genes that include initiation factors, for example, RNA polymerase
207 subunits (*rpoB*), DNA gyrases (*gyrB*), DNA recombination and repair proteins (*recA*) and
208 heat shock proteins (*HSP70*) ^{10,29}. Combinations of these genes can also be used in multi-
209 locus sequence typing (MLST) approaches, increasing phylogenetic resolution and
210 differentiating between closely related species of the same genus ^{30,31}. Marker genes can
211 also reveal key metabolic functions rather than phylogeny; examples include nitrogen
212 cycling (*amoA*, *nifH*, *ntcA*) ^{32,33}, sulfate reduction (*dsrAB*) ³⁴ or phosphorus metabolism
213 (*phnA*, *phnI*, *phnJ*) ³⁵⁻³⁷.

214 The molecular approach has been extended beyond microorganisms by its application to
215 phylogeny and systematics of higher *Eukaryotes*. The Barcode of Life Initiative (BOLI)
216 adapted the molecular approach with the standardized use of a specific gene sequence:
217 the 680 base-pair region of mitochondrial cytochrome c oxidase I (COI), as a means of
218 rapid species identification and discrimination ³⁸.

219 In this paper we collectively define all of these different phylogenetic and functional
220 genes (or gene fragments) as ‘marker genes’ as they are used to profile natural genetic
221 diversity across the Tree of Life, and argue that a small amount of additional effort

222 invested in describing them with specific guidelines in our public databases will
223 revolutionize the types of studies that can be performed with these large data resources.
224 This effort is timely, given the need to determine how climate change and various other
225 anthropogenic perturbations of our biosphere are affecting biodiversity, and how marked
226 changes in our cultural traditions and lifestyles are affecting human microbial ecology.

227

228 *The collective value of marker gene sequences*

229 The quality and quantity of marker gene sequence data used to make phylogenetic
230 assignments, to infer metabolic traits, to unravel succession as a function of the
231 environment, and to assess biogeographic distributions continues to increase rapidly due
232 to the availability of next generation sequencing (NGS) technologies. Clearly, specific
233 associations of microbial dynamics with the environment and geography were achieved
234 for cultured microorganisms long before the advent of metagenomics and NGS
235 technologies³⁹⁻⁴². However, with the new powerful technologies at our service, it is
236 possible to unravel the diversity and function of the uncultured majority as well as to
237 study increasingly complex and/or divergent ecosystems. For example, a clear correlation
238 between phylogenetic similarity and similar living conditions was observed using data in
239 available SSU sequence repositories and culture collections⁴³. In addition, two separate
240 global environmental studies established a latitudinal diversity gradient for marine
241 *Bacteria*^{44,45}. Furthermore, it was shown that temporally-driven environmental factors,
242 such as temperature and nutrients, correlate with local seasonal succession of marine
243 microbial communities⁴⁶. In a cross-habitat study, salinity and pH have been suggested
244 to influence bacterial and archaeal community compositions, respectively^{47,48}. In the

245 human body, it has been suggested that the microbial community composition varies
246 systematically across body habitats, individuals and time ⁴⁹. A recent study combined
247 habitat type and 16S rRNA based operational taxonomic units (OTUs) in a graph-
248 theoretic approach to demonstrate that different habitats harbor unique assemblages of
249 co-occurring microorganisms ⁵⁰. For multicellular organisms, modeling approaches to
250 predict global distributions of marine species have been applied in projects such as
251 AquaMaps ⁵¹. Combination of such efforts with the potential of COI to unveil historical
252 processes may successfully be applied in determining factors responsible for the
253 contemporary geographic distributions of these organisms ⁵².

254 Unfortunately, only a few of these large-scale environmental surveys of biodiversity and
255 biogeography have relied on *existing* marker gene sequence data sets found in the public
256 databases ^{43,47,50,53}. Mainly due to the lack of specific guidelines, most marker gene
257 sequences in databases are sparsely annotated with the information that would be
258 required to underpin data integration, comparative studies, and knowledge generation.
259 Even with complex keyword searches, it is currently impossible to reliably retrieve
260 marker gene sequences that have originated from certain environments or particular
261 locations on Earth; for example, all sequences from ‘soil’ or ‘freshwater lakes’ in a
262 certain region of the world.

263 In human health and the study of epidemiology, it would also be desirable to have
264 additional contextual data to help monitor the origins and regional spreading of
265 pandemics ⁵⁴ and study the variation of the human microbiota ⁵⁵⁻⁵⁷. Combining clinical
266 and environmental datasets could provide new insight into where the trillions of bacteria
267 that inhabit our body come from, and could help predict new outbreaks of disease or

268 assist in understanding the normal ecology of occasional pathogens. Already known
269 correlations of some microbial taxa in with different environmental conditions, such as
270 depth in the marine environment^{58,59}, and pH in the soil environment⁶⁰, can be extended
271 further. Careful integration of bacterial, archaeal and eukaryotic SSU and LSU rRNA
272 sequence data with their geographical and environmental context can shed light on new
273 mechanisms by which organisms from these three domains interact.

274

275 **The MIENS Specification**

276 Few of the publicly available marker gene datasets contain contextual information about
277 the environment such as geographic location, sampling time, habitat, or about
278 experimental procedures used to obtain the DNA sequences. Such information may or
279 may not be available in associated publications but the ‘costs’ in terms of time and energy
280 to collect this by hand or with semi-automated systems from the literature are prohibitive
281⁶¹. Public databases of the International Nucleotide Sequence Database Collaboration
282 (INSDC; comprised of DDBJ (DNA Data Bank of Japan), ENA (European Nucleotide
283 Archive), and GenBank; <http://www.insdc.org>) depend on information submitted by
284 authors to enrich the value of these sequences. We argue that the only way to change the
285 current practice is to establish a standard of reporting that requires contextual data to be
286 deposited at the time of sequence submission³. The adoption of such a standard would
287 elevate the quality, accessibility, and utility of information that can be collected from
288 INSDC.

289 Here we present a reporting guideline for marker genes, MIENS (Minimum Information
290 about an ENvironmental Sequence), which is based on the “Minimum Information about

291 a (Meta) Genome Sequence” (MIGS/MIMS) specification issued by the Genomic
292 Standards Consortium (GSC)⁶². Since its proposal at the sixth GSC meeting in 2008⁶³,
293 the consortium has been working to build a consensus on an ideal and minimum set of
294 contextual data that should be reported for marker genes retrieved from the environment.
295 The proposed MIENS standard (Table 1) extends the MIGS/MIMS specification for
296 genomes and metagenomes by adding two new report types, a “MIENS-survey” and a
297 “MIENS-culture”, the former being the checklist of choice for uncultured diversity
298 marker gene surveys, the latter designed for marker gene sequences obtained from
299 cultured organisms or any material identifiable via voucher specimens.

300 A specific focus of the extended requirements is the sets of measurements and
301 observations describing particular habitats, termed ‘environmental packages’.

302 The MIENS checklist adopts and incorporates the standards being developed by the
303 Consortium for the Barcode of Life (CBOL) ([http://www.barcoding.si.edu/PDF/
304 DWG_data_standards-Final.pdf](http://www.barcoding.si.edu/PDF/DWG_data_standards-Final.pdf)). Therefore, the specification can be universally applied
305 to any marker gene, from SSU rRNA to COI, to cultured and uncultured organisms, to all
306 taxa and to studies ranging from single individuals to complex communities.

307 The MIENS checklist was developed by collating information from several sources and
308 evaluating it in the framework of the existing MIGS/MIMS specification. These include
309 four independent community-led surveys, examination of the parameters reported in
310 published studies, and examination of compliance with optional features in INSDC
311 documents. The overall goal of these activities was to design the backbone of the MIENS
312 specification that describes the most important aspects of marker gene contextual data,
313 and that would encourage users to deposit this contextual data in a standardized fashion.

314 ***Results of community-led surveys***

315 Community surveys are an excellent way to determine researcher preferences for core
316 descriptors. To date, there have been four online surveys about descriptors for marker
317 genes. In the same manner as the Department of Energy Joint Genome Institute's (DOE-
318 JGI) user survey focusing on general descriptor contextual data for marker genes in 2005,
319 the Ribosomal Database Project (RDP) ^{64,65}, SILVA ⁶⁶ and the Terragenome Consortium
320 ⁶⁷ conducted three more user surveys focusing on prevalent habitats for rRNA gene
321 surveys, general descriptor contextual data for rRNA gene sequences and soil
322 metagenome project contextual data, respectively (supplementary information 1).
323 Additionally, following a special session during the 2005 International Census of Marine
324 Microbes (ICoMM), an extensive set of contextual data items were selected, and were
325 analyzed along with survey results.

326 The results of these user surveys provided valuable insights into community requests for
327 contextual data items to be included in the MIENS specification and the main habitats
328 constituting the environmental packages.

329

330 ***Survey of published parameters***

331 We reviewed published rRNA gene studies, retrieved via SILVA and the ICoMM
332 database MICROBIS (The Microbial Oceanic Biogeographic Information System)
333 (<http://icommmbl.edu/microbis>) to further supplement contextual data items that are
334 included in the respective environmental packages. In total, thirty-nine publications from
335 SILVA; including twenty-three publications with more than 500 sequences, and thirteen
336 others retrieved with habitat-specific study queries; and over 40 ICoMM projects were

337 scanned for contextual data items to constitute the core of the environmental package
338 sub-tables (supplementary information 1).

339

340 *Survey of INSDC source feature qualifiers*

341 As a final analysis step, we surveyed usage statistics of INSDC source feature key
342 qualifier values of rRNA gene sequences contained in SILVA (supplementary
343 information 1). Most striking of these results is that less than 10% of the 1.2 million 16S
344 rRNA gene sequences (SILVA release 100) were associated with even basic information
345 such as latitude/longitude, collection date or PCR primers.

346

347 *The MIENS checklist in full*

348 The MIENS specification provides users with an ‘electronic laboratory notebook’
349 containing core contextual data items required for consistent reporting of marker gene
350 investigations. A number of experts in a wide array of topics, guided by a solid
351 rationalization procedure at each step along the way, led the development of these
352 contextual data items.

353 Project details are hosted in the ‘Investigation’ section of MIENS, facilitating access to
354 the outline of contextual data of a marker gene survey. The ‘Environment’ section
355 provides the geospatial, temporal and environmental context. Fourteen ‘environmental-
356 packages’ were developed, with the assistance from user surveys, publication reviews and
357 expert communities working on their respective environments, and were integrated into
358 the ‘MIMS/MIENS extension’ section. These packages provide a wealth of
359 environmental and epidemiological contextual data fields for a complete description of

360 sampling environments (supplementary information 2). Researchers within The Human
361 Microbiome Project ⁶⁸ contributed the host associated and all human packages. The
362 Terragenome Consortium contributed sediment and soil packages. Finally, ICoMM,
363 Microbial Inventory Research Across Diverse Aquatic Long Term Ecological Research
364 Sites (MIRADA-LTERS), and the Max Planck Institute for Marine Microbiology
365 contributed the water package. The MIENS working group developed the remaining
366 packages (air, microbial mat/biofilm, miscellaneous natural or artificial environment,
367 plant-associated, and wastewater/sludge). The package names describe high-level habitat
368 terms in order to be exhaustive. The miscellaneous natural or artificial environment
369 package contains a generic set of parameters, and is included for any other habitat that
370 does not fall into the other thirteen categories. Whenever needed, multiple packages may
371 be used for the description of the environment.

372 The MIGS/MIMS specifications are applicable to MIENS with respect to the nucleic acid
373 sequence source and sequencing contextual data, but have been complemented with
374 further experimental contextual data such as PCR primers and conditions, or target
375 gene/locus.

376 For clarity and ease of use, all items within the MIENS specification are presented with a
377 value syntax description, as well as a clear definition of the item. Whenever terms from a
378 specific ontology are required as the value of an item, these terms can be readily found in
379 the respective ontology browsers, which are linked by URLs in the item definition.

380 Although this version of the MIENS specification does not contain unit specifications, we
381 recommend all units to be chosen from and follow the International System of Units (SI)
382 recommendations. In addition, we strongly urge the community to provide feedback

383 regarding the best unit recommendations for given parameters. To facilitate comparative
384 studies, unit standardization across data sets will be vital in future versions of MIENS.

385

386 *Accessing the MIGS/MIMS/MIENS checklists*

387 The MIGS/MIMS/MIENS checklists are maintained in a relational database system on
388 behalf of the GSC community. This provides a secure and stable mechanism for updating
389 the checklist suite and versioning. An excel version of the checklist is also provided to
390 the community on the GSC web site at: http://genc.org/gc_wiki/index.php/MIENS. The
391 checklist is updated on the GSC web site as development work is carried out on the
392 database end. In the future, we plan to develop programmatic access to this database in
393 order to allow automatic retrieval of the latest version of each checklist for INSDC
394 databases and for GSC community resources. Moreover, the Genomic Contextual Data
395 Markup Language (GCDML) is a reference implementation of the MIGS/MIMS/MIENS
396 checklists by the GSC. It is based on the XML Schema technology and thus serves as an
397 interoperable data exchange format for Web Service based infrastructures⁶⁹.

398

399 **MIENS Adoption by Major Database and Informatics Resources**

400 A variety of efforts are under way to aid sequence submitters in compliance. In the past,
401 the INSDC has issued a reserved 'BARCODE' keyword for the CBOL^{70,71}. Following
402 this model, the INSDC has recently recognized the GSC as an authority for the
403 MIGS/MIMS/MIENS standards and issued it with an official keyword within INSDC
404 nucleotide sequence records⁷². This greatly facilitates automatic validation of the
405 submitted contextual data and provides support for datasets compliant with previous

406 versions by including the checklist version in the keyword.

407 GenBank accepts MIENS metadata in tabular format using the sequin and tbl2asn
408 submission tools, validates MIENS compliance, and reports the MIENS fields in the
409 structured comment block. The ENA Webin submission system provides prepared web
410 forms for the submission of MIENS compliant data; it presents all of the appropriate
411 fields with descriptions, explanations and examples, in addition to validation of the data
412 entered in the forms. An example of a tool that can aid in submission via Sequin or
413 Webin systems is MetaBar ⁷³, a spreadsheet and web-based software tool designed to
414 assist users in the consistent acquisition, electronic storage and submission of contextual
415 data associated with their samples in compliance with the MIGS/MIMS/MIENS
416 specifications.

417 The next-generation Sequence Read Archive (SRA) collects and displays MIENS
418 compliant metadata in the sample and experiment objects. There are several tools that are
419 already available or under development to assist users in SRA submissions. The myRDP
420 SRA PrepKit allows users to prepare and edit their submissions of reads generated from
421 ultra-high-throughput sequencing technologies. A set of suggested attributes in the data
422 forms assist researchers in providing metadata conforming to the MIMS and MIENS
423 specifications. The Investigation/Study/Assay (ISA) Infrastructure is a flexible, freely
424 available software suite that assists in the curation, reporting, and local management of
425 experimental metadata from studies employing one or a combination of technologies,
426 including high-throughput sequencing. Specific ISA configurations (available from
427 http://genc.org/gc_wiki/index.php/Adopters#ISA_infrastructure) have been developed to
428 ensure MIENS compliance by providing templates and validation capability while

429 another tool, ISAconverter, produces SRA.xml documents, thereby facilitating
430 submission to the SRA repository⁷⁴.

431 The SILVA, RDP, Greengenes and the ICoMM resources have participated in the
432 development of MIENS, and are now taking the standardization one step further by
433 establishing tools and resources to aid in compliance.

434 Further detailed guidance for submission processes can be found under the respective
435 wiki pages (http://gensc.org/gc_wiki/index.php/MIENS) of the MIENS standard.

436

437 *Examples of MIENS compliant datasets*

438 Several MIENS compliant reports are included in the supplementary information 3.

439 These include; a 16S rRNA gene survey from samples obtained in the North Atlantic, an

440 18S pyrotag study of anaerobic protists in the permanently anoxic basin of the North Sea,

441 a *pmoA* survey from desert soils of Negev Desert, Israel, a *dsrAB* survey from marine

442 sediments from the Gulf of Mexico, and finally a 16S pyrotag study of bacterial diversity

443 in the Western English Channel (publicly accessible via SRA study accession number

444 SRP001108). Two further MIENS compliant 16S submissions are available in INSDC

445 under the accession numbers GU949561.1 and GU949562.1.

446

447 *MIENS – a ‘living standard’*

448 MIENS, as well as MIGS/MIMS, are ‘living checklists’ and not final specifications.

449 Therefore, further developments, extensions, and enhancements will be recognized, and

450 improved versions of the checklists, if necessitated, will be released annually, while

451 preserving the validity of former versions. A public issue tracking system, which can be

452 reached via <http://mixs.gensc.org/>, is set up to track changes and accomplish feature
453 requests. The final decisions about their implementation will be concluded by the MIENS
454 working group.

455

456 **Conclusions and Call for Action**

457 The GSC is an international working body with a stated mission of working towards
458 richer descriptions of our complete collection of genomes and metagenomes. With the
459 development of the MIENS specification, this mission has been extended to marker gene
460 sequences as well. The GSC is an open initiative that welcomes the participation of the
461 wider community. This includes an open call to contribute to refinements of the MIENS
462 specification or its implementation.

463 The adoption of the MIENS standard by major data providers and organizations as well
464 as the three primary public sequence data repositories (INSDC) with an active poll for
465 MIENS compliant data underlines and seconds the efforts to contextually enrich our
466 marker gene collection, and complements the recent efforts to contextually enrich other
467 (meta) omics data. The MIENS checklist has been developed to the point that it is ready
468 to be used in the publication of sequences. A defined procedure for requesting new
469 features and the stable release cycles will facilitate implementation of the standard across
470 the community. Widespread compliance among authors, adoption by journals and use by
471 informatics resources will vastly improve our collective ability to mine and integrate
472 invaluable sequence data collections for knowledge and application driven research. In
473 particular, the ability to combine microbial community samples collected from any
474 source, using the universal Tree of Life as a yardstick to compare even the most diverse

475 communities, should provide new insights into the dynamic spatial and temporal
476 distribution of microbial life on our planet and even on our own bodies.

477 **References**

- 478 1 Community cleverness required. *Nature* **455**, 1-1 (2008).
- 479 2 Field, D. *et al.* 'Omics Data Sharing. *Science* **326**, 234-236 (2009).
- 480 3 Taylor, C. F. *et al.* Promoting coherent minimum reporting guidelines for
481 biological and biomedical investigations: the MIBBI project. *Nat Biotechnol* **26**, 889-896
482 (2008).
- 483 4 Woese, C. R. and Fox, E. Phylogenetic structure of the prokaryotic domain: the
484 primary kingdoms. *Proc Nat Acad Sci USA* **74**, 5088-5090 (1977).
- 485 5 Woese, C. R., Kandler, O., and Wheelis, M. L. Towards a natural system of
486 organisms: proposal for the domains *Archaea*, *Bacteria*, and *Eucarya*. *Proc Nat Acad Sci*
487 *USA* **87**, 4576-4579 (1990).
- 488 6 Amann, R. I., Ludwig, W., and Schleifer, K. H. Phylogenetic identification and
489 in-situ detection of individual microbial cells without cultivation. *Microbiol Rev* **59**, 143-
490 169 (1995).
- 491 7 Curtis, T. P., Sloan, W. T., and Scannell, J. W. Estimating prokaryotic diversity
492 and its limits. *Proc Nat Acad Sci USA* **99**, 10494-10499 (2002).
- 493 8 Turrone, F. *et al.* Human gut microbiota and bifidobacteria: from composition to
494 functionality. *Antonie van Leeuwenhoek* **94**, 35-50 (2008).
- 495 9 Ludwig, W. *et al.* Bacterial phylogeny based on comparative sequence analysis.
496 *Electrophoresis* **19**, 554-568 (1998).
- 497 10 Ludwig, W. and Schleifer, K. H. in *Microbial phylogeny and evolution, concepts*
498 *and controversies*, edited by J. Sapp (Oxford university press, New York, 2005), pp. 70-
499 98.

500 11 Ciccarelli, F. D. *et al.* Toward automatic reconstruction of a highly resolved tree
501 of life. *Science* **311**, 1283-1287 (2006).

502 12 Teeling, H. and Glöckner, F. O. RibAlign: a software tool and database for
503 eubacterial phylogeny based on concatenated ribosomal protein subunits. *BMC*
504 *Bioinformatics* **7** (2006).

505 13 Stahl, D. A., Lane, D. J., Olsen, G. J., and Pace, N. R. Analysis of hydrothermal
506 vent associated symbionts by ribosomal RNA sequences. *Science* **224**, 409-411 (1984).

507 14 Pace, N. R., Stahl, D. A., Olsen, G. J., and Lane, D. J. Analyzing natural
508 microbial populations by rRNA sequences. *ASM News* **51**, 4-12 (1985).

509 15 Olsen, G. J. *et al.* Microbial ecology and evolution: a ribosomal RNA approach.
510 *Annu Rev Microbiol* **40**, 337-365 (1986).

511 16 Giovannoni, S. J., Britschgi, T. B., Moyer, C. L., and Field, K. G. Genetic
512 diversity in Sargasso Sea bacterioplankton. *Nature* **345**, 60-63 (1990).

513 17 Ward, D. M., Weller, R., and Bateson, M. M. 16S rRNA sequences reveal
514 numerous uncultured microorganisms in a natural community. *Nature* **345**, 63-65 (1990).

515 18 DeLong, E. F. *Archaea* in coastal marine environments. *Proc Nat Acad Sci USA*
516 **89**, 5685-5689 (1992).

517 19 Fuhrman, J. A., McCallum, K., and Davis, A. A. Novel major archaeobacterial
518 group from marine plankton. *Nature* **356**, 148-149 (1992).

519 20 Pace, N. R. A molecular view of microbial diversity and the biosphere. *Science*
520 **276**, 734-740 (1997).

521 21 Diez, B., Pedros-Alio, C., and Massana, R. Study of Genetic Diversity of
522 Eukaryotic Picoplankton in Different Oceanic Regions by Small-Subunit rRNA Gene

523 Cloning and Sequencing. *Appl Environ Microbiol* **67**, 2932-2941 (2001).

524 22 Hewson, I. and Fuhrman, J. A., Richness and diversity of bacterioplankton species
525 along an estuarine gradient in Moreton Bay, Australia. *Appl Environ Microbiol* **70**, 3425-
526 3433 (2004).

527 23 López-García, P., López-López, A., Moreira, D., and Rodríguez-Valera, F.
528 Diversity of free-living prokaryotes from a deep-sea site at the Antarctic Polar Front.
529 *Fems Microbiol Ecol* **36**, 193-202 (2001).

530 24 Lopez-Garcia, P., Rodriguez-Valera, F., Pedros-Alio, C., and Moreira, D.
531 Unexpected diversity of small eukaryotes in deep-sea Antarctic plankton. *Nature* **409**,
532 603-607 (2001).

533 25 Moon-van der Staay, S. Y., De Wachter, R., and Vaultot, D. Oceanic 18S rDNA
534 sequences from picoplankton reveal unsuspected eukaryotic diversity. *Nature* **409**, 607-
535 610 (2001).

536 26 Huber, J. A., Butterfield, D. A., and Baross, J. A. Temporal changes in archaeal
537 diversity and chemistry in a mid-ocean ridge seafloor habitat. *Appl Environ Microbiol*
538 **68**, 1585-1594 (2002).

539 27 Rappe, M. S. and Giovannoni, S. J. The uncultured microbial majority. *Annu Rev*
540 *Microbiol* **57**, 369-394 (2003).

541 28 Doolittle, W. F. Fun With Genealogy. *Proc Nat Acad Sci USA* **94**, 12751-12753
542 (1997).

543 29 Huynen, M. A. and Bork, P. Measuring genome evolution. *Proc Nat Acad Sci*
544 *USA* **95**, 5849-5856 (1998).

545 30 Ivars-Martinez, E. *et al.* Biogeography of the ubiquitous marine bacterium

546 *Alteromonas macleodii* determined by multilocus sequence analysis. *Mol Ecol* **17**, 4092-
547 4106 (2008).

548 31 Cole, J. R., Konstantinidis, K., Farris, R. J., and Tiedje, J. M. in *Environmental*
549 *Molecular Microbiology*, edited by W.-T. Liu and J.K. Jansson (Caister Academic Press
550 UK, 2010), pp. 1-19.

551 32 Zehr, J. P., Mellon, M. T., and Zani, S. New nitrogen-fixing microorganisms
552 detected in oligotrophic oceans by amplification of nitrogenase (nifH) genes. *Appl*
553 *Environ Microbiol* **64**, 3444-3450 (1998).

554 33 Francis, C. A., Beman, J. M., and Kuypers, M. M. M. New processes and players
555 in the nitrogen cycle: the microbial ecology of anaerobic and archaeal ammonia
556 oxidation. *Isme J* **1**, 19-27 (2007).

557 34 Minz, D. *et al.* Diversity of sulfate-reducing bacteria in oxic and anoxic regions of
558 a microbial mat characterized by comparative analysis of dissimilatory sulfite reductase
559 genes. *Appl Environ Microbiol* **65**, 4666-4671 (1999).

560 35 Gilbert, J. A. *et al.* Potential for phosphonoacetate utilization by marine bacteria
561 in temperate coastal waters. *Environ Microbiol* **11**, 111-125 (2009).

562 36 Martinez, A., W. Tyson, G., and DeLong, E., F. Widespread known and novel
563 phosphonate utilization pathways in marine bacteria revealed by functional screening and
564 metagenomic analyses. *Environ Microbiol* **9999** (2009).

565 37 Thomas, S. *et al.* Evidence for phosphonate usage in the coral holobiont. *Isme J*
566 **4**, 459-461 (2010).

567 38 Hebert, P. D. N., Cywinska, A., Ball, S. L., and Dewaard, J. R. Biological
568 identifications through DNA barcodes. *Proc R Soc Lond B Biol Sci* **270**, 313-321 (2003).

569 39 ZoBell, C. E. and Johnson, F. H. The influence of hydrostatic pressure on the
570 growth and viability of terrestrial and marine bacteria. *J Bacteriol* **57**, 179 (1949).

571 40 Brock, T. D. and Brock, M. L. Relationship between Environmental Temperature
572 and Optimum Temperature of Bacteria along a Hot Spring Thermal Gradient. *J Appl*
573 *Microbiol* **31**, 54-58 (1968).

574 41 Cho, J.-C. and Tiedje, J. M. Biogeography and Degree of Endemicity of
575 Fluorescent Pseudomonas Strains in Soil. *Appl Environ Microbiol* **66**, 5448-5456 (2000).

576 42 Pomeroy, L. R. and Wiebe, W. J. Temperature and substrates as interactive
577 limiting factors for marine heterotrophic bacteria. *Aquat Microb Ecol* **23**, 187-204 (2001).

578 43 von Mering, C. *et al.* Quantitative phylogenetic assessment of microbial
579 communities in diverse environments. *Science* **315**, 1126-1130 (2007).

580 44 Pommier, T. *et al.* Global patterns of diversity and community structure in marine
581 bacterioplankton. *Mol Ecol* **16**, 867-880 (2007).

582 45 Fuhrman, J. A. *et al.* A latitudinal diversity gradient in planktonic marine bacteria.
583 *Proc Nat Acad Sci USA* **105**, 7774-7778 (2008).

584 46 Gilbert, J., A. *et al.* The seasonal structure of microbial communities in the
585 Western English Channel. *Environ Microbiol* **11**, 3132-3139 (2009).

586 47 Lozupone, C. A. and Knight, R. Global patterns in bacterial diversity. *Proc Nat*
587 *Acad Sci USA* **104**, 11436-11440 (2007).

588 48 Auguet, J.-C., Barberan, A., and Casamayor, E. O. Global ecological patterns in
589 uncultured Archaea. *Isme J* **4**, 182-190 (2010).

590 49 Costello, E. K. *et al.* Bacterial community variation in human body habitats across
591 space and time. *Science* **326**, 1694-1697 (2009).

592 50 Chaffron, S., Rehrauer, H., Pernthaler, J., and von Mering, C. A global network of
593 coexisting microbes from environmental and whole-genome sequence data. *Genome Res*
594 **20**, 947-959 (2010).

595 51 Kaschner, K. *et al.* AquaMaps: Predicted range maps for aquatic species,
596 Available at <http://www.aquamaps.org/>, (2008).

597 52 Workshops Report and Recommendations DNA Barcoding of Marine
598 Biodiversity (MarBOL) presented at the MarBOL Workshops, 2009 (unpublished).

599 53 Tamames, J. *et al.* Environmental distribution of prokaryotic taxa. *BMC*
600 *Microbiology* **10**, 85.

601 54 Schriml, L. M. *et al.* GeMInA, Genomic Metadata for Infectious Agents, a
602 geospatial surveillance pathogen database. *Nucl Acids Res* **38**, D754-D764 (2010).

603 55 Palmer, C. *et al.* Development of the Human Infant Intestinal Microbiota. *PLoS*
604 *Biol* **5**, e177 (2007).

605 56 Ravel, J. *et al.* Vaginal microbiome of reproductive-age women. *Proc Nat Acad*
606 *Sci USA e-pub ahead of print* (2010).

607 57 Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic
608 sequencing. *Nature* **464**, 59-65 (2010).

609 58 DeLong, E. F. *et al.* Community genomics among stratified microbial
610 assemblages in the ocean's interior. *Science* **311**, 496-503 (2006).

611 59 Moreira, D. Rodriguez-Valera, F., and Lopez-Garcia, P., Metagenomic analysis of
612 mesopelagic Antarctic plankton reveals a novel deltaproteobacterial group. *Microbiology*
613 **152**, 505-517 (2006).

614 60 Lauber, C. L., Hamady, M., Knight, R., and Fierer, N. Soil pH as a predictor of

615 soil bacterial community structure at the continental scale: a pyrosequencing-based
616 assessment. *Appl Environ Microbiol* **75**, 5111-5120 (2009).

617 61 Hirschman, L. *et al.* Habitat-Lite: a GSC case study based on free text terms for
618 environmental metadata. *OMICS* **12**, 129-136 (2008).

619 62 Field, D. *et al.* The minimum information about a genome sequence (MIGS)
620 specification. *Nat Biotechnol* **26**, 541-547 (2008).

621 63 Field, D. *et al.* Meeting reports from the Genomic Standards Consortium (GSC)
622 Workshops 6 and 7. *SIGS* **1**, 68-71 (2009).

623 64 Cole, J. R. *et al.* The ribosomal database project (RDP-II): introducing myRDP
624 space and quality controlled public data. *Nucl Acids Res* **35**, D169-172 (2007).

625 65 Cole, J. R. *et al.* The Ribosomal Database Project: improved alignments and new
626 tools for rRNA analysis. *Nucl Acids Res* **37**, D141-145 (2009).

627 66 Pruesse, E. *et al.* SILVA: a comprehensive online resource for quality checked
628 and aligned ribosomal RNA sequence data compatible with ARB. *Nucl Acids Res* **35**,
629 7188-7196 (2007).

630 67 Vogel, T. M. *et al.* TerraGenome: a consortium for the sequencing of a soil
631 metagenome. *Nat Rev Micro* **7**, 252-252 (2009).

632 68 Turnbaugh, P. J. *et al.* The Human Microbiome Project. *Nature* **449**, 804-810
633 (2007).

634 69 Kottmann, R. *et al.* A standard MIGS/MIMS compliant XML schema: Toward
635 the development of the Genomic Contextual Data Markup Language (GCDML). *OMICS*
636 **12**, 115-121 (2008).

637 70 Benson, D. A. *et al.* GenBank. *Nucl. Acids Res.* **35**, D21-25 (2007).

638 71 Benson, D. A. *et al.* GenBank. *Nucl. Acids Res.* **36**, D25-30 (2008).

639 72 Hirschman, L. *et al.* Meeting report: Metagenomics, Metadata and Meta-analysis”
640 (M3) Workshop at the Pacific Symposium on Biocomputing 2010. *SIGS* **2**, 357-360
641 (2010).

642 73 Hankeln, W. *et al.* MetaBar - a tool for consistent contextual data acquisition and
643 standards compliant submission. *BMC Bioinformatics* **11**, 358 (2010).

644 74 Rocca-Serra, P. *et al.* ISA infrastructure: supporting standards-compliant
645 experimental reporting and enabling curation at the community level. *Bioinformatics* **26**,
646 2354-2356 (2010).

647

648

		Report type	
		MIENS survey	MIENS culture
Investigation			
Submitted to INSDC ^[boolean]	Depending on the study (large-scale e.g. done with next generation sequencing technology, or small-scale) sequences have to be submitted to SRA (Sequence Read Archives), DRA (DDBJ Sequence Read Archive) or via the classical Webin/Sequin systems to Genbank, ENA and DDBJ	M	M
Investigation type ^[survey or culture]	Nucleic Acid Sequence Report is the root element of all MIENS compliant reports as standardized by Genomic Standards Consortium (GSC). This field is either MIENS survey or MIENS culture	M	M
Project name	Name of the project within which the sequencing was organized	M	M
Environment			
Geographic location (latitude and longitude ^[float, point, transect and region])	The geographical origin of the sample as defined by latitude and longitude. The values should be reported in decimal degrees and in WGS84 system	M	M
Geographic location (depth ^[integer, point, interval, unit])	Please refer to the definitions of depth in the environmental packages	E	E
Geographic location (elevation of site ^[integer, unit] ; altitude of sample ^[integer, unit])	Please refer to the definitions of either altitude or elevation in the environmental packages	E	E
Geographic location (country and/or sea ^[INSDC or GAZ] ; region ^[GAZ])	The geographical origin of the sample as defined by the country or sea name. Country, sea, or region names should be chosen from the INSDC list (http://insdc.org/country.html), or the GAZ (Gazetteer, v1.446) ontology (http://bioportal.bioontology.org/visualize/40651)	M	M
Collection date ^[ISO8601]	The time of sampling, either as an instance (single point in time) or interval. In case no exact time is available, the date/time can be right truncated i.e. all of these are valid times: 2008-01-23T19:23:10+00:00; 2008-01-23T19:23:10; 2008-01-23; 2008-01; 2008; Except: 2008-01; 2008 all are ISO6801 compliant	M	M
Environment (biome ^[EnvO])	In environmental biome level are the major classes of ecologically similar communities of plants, animals, and other organisms. Biomes are defined based on factors such as plant structures, leaf types, plant spacing, and other factors like climate. Examples include: desert, taiga, deciduous woodland, or coral reef. Environment Ontology (EnvO) (v1.53) terms listed under environmental biome can be found from the link: http://bioportal.bioontology.org/visualize/44405/?conceptid=ENVO%3A00000428	M	M
Environment (feature ^[EnvO])	Environmental feature level includes geographic environmental features. Examples include: harbor, cliff, or lake. EnvO (v1.53) terms listed under environmental feature can be found from the link: http://bioportal.bioontology.org/visualize/44405/?conceptid=ENVO%3A00002297	M	M

Environment (material ^[EnvO])	The environmental material level refers to the matter that was displaced by the sample, prior to the sampling event. Environmental matter terms are generally mass nouns. Examples include: air, soil, or water. EnvO (v1.53) terms listed under environmental matter can be found from the link: http://bioportal.bioontology.org/visualize/44405/?conceptid=ENVO%3A00010483	M	M
MIGS/MIMS/MIENS Extension			
Environmental package ^[air, host-associated, human-associated, human-skin, human-oral, human-gut, human-vaginal, microbial mat/biofilm, miscellaneous natural or artificial environment, plant-associated, sediment, soil, wastewater/sludge, water]	MIGS/MIMS/MIENS extension for reporting of measurements and observations obtained from one or more of the environments where the sample was obtained. All environmental packages listed here are further defined in separate subtables. By giving the name of the environmental package, a selection of fields can be made from the subtables and can be reported	M	M
Nucleic acid sequence source			
Isolation and growth conditions ^[PMID, DOI, or URL]	Publication reference in the form of pubmed ID (PMID), digital object identifier (DOI), or URL for Isolation and growth condition specifications of the organism/material	-	M
Sequencing			
Target gene or locus (e.g. 16S rRNA, 18S rRNA, nif, amoA, rpo, V6, ITS)	Targeted gene, locus or gene region name for marker gene study	M	M
Sequencing method (e.g. dideoxysequencing, pyrosequencing, polony)	Sequencing method used; e.g. Sanger, pyrosequencing, ABI-solid.	M	M

Table 1. Items for the MIENS specification and their mandatory (M), conditionally mandatory (C) (the item is mandatory only when applicable to the study) or recommended (X) status for both MIENS-survey and MIENS-culture checklists. MIENS-survey is applicable to contextual data for marker gene sequences, obtained directly from the environment, without culturing or identification of the organisms. MIENS-culture, on the other hand, applies to the contextual data for marker gene sequences from cultured or voucher-identifiable specimens. Both MIENS-survey and culture checklists can be used for any type of marker gene sequence data, ranging from 16S, 18S, 23S, 28S rRNA to COI, hence the checklists are universal for all three domains of life. ‘-’ denotes that an item is not applicable for a given checklist.

'E' denotes that a field has environment-specific requirements. For example, while 'depth' is mandatory for environments water, sediment or soil; it is optional for human-associated environments. Item names are followed by a short description of the value of the item in parentheses and/or value type in brackets as a superscript. Whenever applicable, value types are chosen from a controlled vocabulary (CV), or an ontology from the Open Biological and Biomedical Ontologies (OBO) foundry (<http://www.obofoundry.org>). This table only presents the very core of MIENS checklists, i.e. only mandatory items for each checklist. Supplementary information 2 in spreadsheet format contains all MIENS items, the tables for environmental packages in the MIMS/MIENS extension, and GenBank structured comment name that should be used for submitting MIENS data to GenBank.