

1 **Metagenomics uncovers gaps in amplicon-based detection of microbial diversity**

2
3 Emiley A. Eløe-Fadrosh¹, Natalia N. Ivanova¹, Tanja Woyke¹, Nikos C. Kyrpides^{1*}

4
5 ¹Joint Genome Institute, Walnut Creek, CA 94598, USA.

6
7 *Corresponding author. Mailing address: DOE Joint Genome Institute, Walnut Creek,
8 CA 94598. Phone: 925-296-5718. E-mail: nckyrpides@lbl.gov.

9
10 Our understanding of the microbial world has made significant advances through the
11 application of molecular approaches, particularly PCR-based amplification of the small-
12 subunit ribosomal (SSU rRNA) gene¹. The accumulation of SSU rRNA gene sequences
13 has increased dramatically², yet the pace of discovery for new taxonomic lineages
14 uncovered through PCR-based biodiversity surveys has seemingly slowed³. On the
15 other hand, new candidate phyla have been identified using metagenomic and single-
16 cell genomic techniques⁴⁻⁶. The above raises the question: Have we approached
17 saturation or are there systematic biases in PCR-based surveys that preclude discovery
18 of additional microbial lineages? Arguably, there is a wealth of microbial clades that are
19 systematically under-represented or missed altogether, leaving major taxonomic “blind
20 spots”⁷. PCR amplification biases, including primer mismatches, are well-recognized
21 pitfalls in biodiversity surveys⁸, yet a comprehensive analysis of the prevalence of such
22 “blind spots” has not been undertaken⁹. Here, we systematically surveyed primer fidelity
23 in SSU rRNA gene sequences recovered from over 6,000 assembled metagenomes
24 sampled globally. Our findings show approximately 10% of environmental microbial
25 sequences might be missed from classical PCR-based SSU rRNA gene surveys, mostly
26 members of the Candidate Phyla Radiation (CPR)⁴ and as-yet uncharacterized
27 Archaea. These results underscore the extent of uncharacterized microbial diversity,
28 and provide fruitful avenues for describing additional phylogenetic lineages.

29 We compared SSU rRNA gene sequences ≥ 1000 bp recovered from
30 metagenomes (38,454 SSU rRNA) and isolate genomes (25,439 SSU rRNA) available
31 through the Integrated Microbial Genomes with Microbiome Samples (IMG/M)¹⁰ against
32 commonly used PCR primer sets (515F-806R from the Earth Microbiome Project^{2,11};
33 515F-926R¹²; 341F-785R¹³; and 357F-926R from the Human Microbiome Project¹⁴),
34 and generated a weighted score based on the number of mismatches (**Fig. 1a**;
35 **Supplementary Table 1**; Methods). Unexpectedly, this analysis indicates that a
36 minimum of 9.6% of environmental bacterial and/or archaeal sequences based on
37 metagenomic data might not be recovered using a targeted PCR survey (**Fig. 1a**). The
38 newly modified primers part of the Earth Microbiome Project appear to more fully
39 capture SSU rRNA diversity (9.6% metagenomic sequences would likely be missed),
40 while the relatively poor performance of the Human Microbiome Project primer set may
41 be due to the narrow bacterial target range (22.4% bacterial metagenomic sequences
42 would likely be missed). Notably, combining any two primer pairs does not significantly
43 improve taxonomic coverage (**Supplementary Fig. 1**). Even when the four best-
44 performing primer sets are combined together, 5.5% of sequences remain that would be
45 missed based on metagenomic data, suggesting the use of multiple primer sets might
46 slightly improve recovery but does not fully resolve the issue (**Supplementary Fig. 1**).
47 Overall, these results are most likely an underestimate of diversity missed since our
48 data consisted of assembled metagenomic contigs, primarily representing abundant
49 organisms in a sample and neglecting less abundant phylotypes (e.g. the ‘rare
50 biosphere’). Taken together, we hypothesize that amplification-unbiased exploration of
51 microbial diversity via metagenome and metatranscriptome sequencing will
52 unquestionably improve our current view of the microbial tree of life.

53 An evaluation of base-specific biases for the commonly used PCR primer sets
54 revealed a subset of bases contributing to the percentage of metagenomic SSU rRNA
55 gene sequences that would likely be missed in PCR-based surveys (**Supplementary**
56 **Fig. 2**). These bases, or “hot spots,” could be candidates for increased degeneracy in

57 the current primers to capture a greater fraction of the microbial diversity or serve as
58 guides in the design of new primer sets. Regardless, modifications of the specific bases
59 contributing to the inefficiency of these primers would need to be experimentally
60 validated prior to proposing improved primers.

61 It has been previously shown experimentally that primer mismatches can
62 significantly affect species evenness⁴. In order to verify this, we evaluated the SSU
63 rRNA gene sequences with primer mismatches from a set of eight matched
64 metagenomic datasets and SSU rRNA surveys from a diversity of environmental locales
65 to determine whether the predicted mismatches would impact recovery in SSU rRNA
66 surveys (**Supplementary Table 2**). In all eight matched metagenomic datasets and
67 SSU rRNA surveys, the computationally predicted missed sequences were not
68 recovered in the SSU rRNA survey. We further evaluated the SSU rRNA gene
69 sequences with primer mismatches from the large-scale Human Microbiome Project
70 metagenomic datasets and compared them with their corresponding SSU rRNA
71 surveys¹⁴. We searched more than 34 million reads from over 4,200 Human Microbiome
72 Project SSU rRNA surveys against the 130 sequences from the HMP metagenomes
73 predicted to be missed and found only 2,060 matches. Although we observe a small
74 number of matches, these represent only 0.006% relative abundance in the SSU rRNA
75 surveys and are in contrast to the two orders of magnitude greater abundance based on
76 shotgun metagenome data. Together, these results are consistent with our
77 computational predictions, and suggest that the primer mismatches would likely
78 significantly reduce or prevent the recovery of taxonomic “blind spots” in PCR-based
79 surveys.

80 Our analysis revealed phylogenetic patterns for those sequences that would
81 presumably be missed with the widely used environmental primer set 515F-806R
82 (overall primer weighted score ≥ 1)^{2,11}. As anticipated, members of the recently
83 described Candidate Phyla Radiation (CPR⁴; including Parcubacteria (OD1),
84 Microgenomates (OP11), WWE3, Berkelbacteria (ACD58), Saccharibacteria (TM7),

85 WS6, Peregrinibacteria (PER), and Kazan phyla) collectively represented 70% of SSU
86 rRNA gene sequences that would likely be missed in PCR-based surveys (**Fig. 1b**;
87 **Supplementary Fig. 3**). The overall length of the SSU rRNA gene further compounds
88 these findings for the CPR since the prevalence of encoded introns may hamper
89 amplification fidelity⁴. Within the domain Archaea, more than half of the taxonomic “blind
90 spots” are phylogenetically positioned outside of the recognized phylum-level lineages,
91 revealing significant untapped archaeal diversity awaiting discovery (**Fig. 1b**;
92 **Supplementary Fig. 3**). Recent efforts to resolve archaeal diversity through genome-
93 resolved metagenomic analyses have yielded substantial progress towards a better
94 understanding of archaeal evolutionary history^{5,15}.

95 Genomic mapping of the tree of life has been accelerated through application of
96 both metagenomic and single-cell genomic sequencing of samples taken directly from
97 the environment without the arduous task of cultivation. Within the last few years,
98 significant advances in high-throughput single-cell genomics have provided some of the
99 first genomic insight for a wealth of candidate phylogenetic lineages previously known
100 only through SSU rRNA gene sequencing⁶. Further, deep sequencing of environmental
101 samples, combined with improved metagenomic assembly and binning methods are
102 yielding complete or near-complete genomes from many novel bacterial and archaeal
103 lineages (see^{4,5,16}). We performed a phylogenetic analysis of all SSU rRNA gene
104 sequences derived from metagenomes (regardless of whether these sequences had
105 primer mismatches) to identify the sequences that could not be placed with known
106 bacterial or archaeal phyla. These data suggest that bacterial diversity has been
107 charted extensively, with minimal SSU rRNA gene orphan sequences not assignable to
108 any phylum (**Table 1**). This is in stark contrast to that of the Archaea, where significant
109 diversity exists beyond the currently described major lineages. Additionally, there is
110 likely to be sizable taxonomic novelty at more refined taxonomic levels, such as class
111 and order.

112 Habitat distribution of those sequences unaffiliated with currently recognized
113 phyla based on our phylogenetic analysis provides insight into the suite of
114 environmental locales potentially hosting as-yet uncharacterized microbial life (**Fig. 2**).
115 The habitats where more unclassified sequences are found include “extreme” habitats
116 with unique environmental parameters (e.g. extremes in temperature, pressure, and
117 chemical composition), favoring a distinct composition of microbial communities in these
118 environments. Our data does suggest that these habitats may harbor more divergent
119 phylogenetic groups, specifically within the archaea. On the other hand, environments
120 such as marine cold seeps also comprise a wealth of uncharacterized microbial
121 diversity with sampling challenges traditionally hampering their genomic exploration for
122 the as-yet uncharacterized microbial life. More recent targeted sampling efforts have
123 begun to shed light on these unique environments^{4-6,15}.

124 Taken together, we posit that more recent technologies including metagenomics
125 and single-cell genomics are bringing into focus the taxonomic “blind spots” that have
126 thus far eluded detection. While PCR-based SSU rRNA surveys provided a sizeable
127 snapshot of the phylogenetic structure of the microbial world, we are poised to make
128 substantial strides forward in our efforts to more comprehensively capture microbial
129 diversity. We anticipate new sequencing technologies such as single-molecule
130 sequencing will enable full recovery of genomic content for those bacterial and archaeal
131 phyla that are known only through SSU rRNA, as well as candidate phyla discovered
132 through systematic large-scale data mining efforts. Uncharted microbial diversity will
133 soon find its place in the tree of life.

134

135 **Methods**

136 All SSU rRNA gene sequences ≥ 1000 bp were retrieved from the Integrated Microbial
137 Genomes with Microbiome Samples (IMG/M)¹⁰ database and manually curated to
138 remove mitochondrial and chloroplast sequences. A total of 3,003 mitochondrial and
139 668 chloroplast sequences were identified and culled from the metagenomic dataset.

140 PrimerProspector¹⁷ was used to evaluate primer sets against SSU rRNA gene
141 sequences ≥ 1000 bp recovered from metagenomes (38,454) and isolate genomes
142 (25,439). Primers used in the analysis are listed in **Supplementary Table 1**. Weighted
143 scores were calculated as follows: non-3' mismatches = 0.4; 3' mismatches = 1; non-3'
144 gaps = 1; and 3' gaps = 3. Therefore, a weighted score ≥ 1 could consist of three non-3'
145 mismatches and was considered a missed sequence. Validation of the eight matched
146 metagenomes and corresponding SSU rRNA surveys (**Supplementary Table 2**) was
147 performed using blat¹⁸ with the metagenomic sequences trimmed to the 515F-806R
148 region. Similarly, the Human Microbiome Project metagenomic datasets and the
149 corresponding SSU rRNA surveys¹⁴ were validated using the metagenomic sequences
150 trimmed to the 357F-926R region. Taxonomic affiliation was determined by aligning all
151 SSU rRNA genes using cmalign (--matchonly option)¹⁹ against bacterial and archaeal
152 hmm models and building a phylogenetic tree using RAxML²⁰ with reference sequences
153 to identify coherent clades. Taxonomic affiliation was assigned as follows: Cultured
154 phyla, percentage of metagenomic SSU rRNA gene sequences matching phyla from
155 cultured representatives; Uncultured phyla, percentage of metagenomic SSU rRNA
156 gene sequences matching candidate phyla without cultured representatives; and
157 Unaffiliated with current phyla, percentage of metagenomic SSU rRNA gene sequences
158 that do not match a known, named phylum. Primer sequence logos were generated
159 using WebLogo²¹.

160

161 **References**

- 162 1. Pace, N. R. *Microbiol. Mol. Biol. Rev.* **73**, 565-576 (2009).
- 163 2. Caporaso, J. G. *et al. ISME J.* **6**, 1621-1624 (2012).
- 164 3. Yarza, P. *et al. Nature Rev. Microbiol.* **12**, 635-645 (2014).
- 165 4. Brown, C. T. *et al. Nature* **523**, 208-211 (2015).
- 166 5. Castelle, C. J. *et al. Current Biol.* **25**, 690-701 (2015).
- 167 6. Rinke, C. *et al. Nature* **499**, 431-437 (2013).
- 168 7. Lynch, M. D. J. & Neufeld, J. D. *Nature Rev. Microbiol.* **13**, 217-229 (2015).
- 169 8. Engelbrektson, A. *et al. ISME J.* **4**, 642-647 (2010).
- 170 9. Woyke, T. & Rubin, E. M. *Science* **346**, 698-699 (2014).

- 171 10. Markowitz, V. M. *et al. Nucleic Acids Res.* **42**, D568-D573 (2014).
172 11. Gilbert, J. A. *et al. Stand. Genomic Sci.* **3**, 249-253 (2010).
173 12. Parada, A., Needham, D. M. & Fuhrman, J. A. *Environ. Microbiol.*, doi:
174 10.1111/1462-2920.13023 (2015).
175 13. Klindworth, A. *et al. Nucleic Acids Res.* **41**, e1 (2013).
176 14. The Human Microbiome Consortium. *Nature* **486**, 215-221 (2012).
177 15. Spang, A. *et al. Nature* **521**, 173-179 (2015).
178 16. Albertsen, M. *et al. Nat. Biotechnol.* **31**, 533-538 (2013).
179 17. Walters, W. A. *et al. Bioinformatics* **27**, 1159-1161 (2011).
180 18. Kent, W. *Genome Res.* **12**, 656-664 (2002).
181 19. Nawrocki, E. P. & Eddy, S. R. *Bioinformatics* **29**, 2933-2935 (2013).
182 20. Stamatakis, A. *Bioinformatics* **22**, 2688-2690 (2006).
183 21. Crooks, G. E. *et al. Genome Res.* **14**, 1188-1190 (2004).
184

185 Correspondence and requests for materials should be addressed to N.C.K.

186

187 **Acknowledgments**

188 This work was conducted by the U.S. Department of Energy Joint Genome Institute, a
189 DOE Office of Science User Facility, under Contract No. DE-AC02-05CH11231 and
190 used resources of the National Energy Research Scientific Computing Center, which is
191 supported by the Office of Science of the U.S. Department of Energy under Contract
192 No. DE-AC02-05CH11231.

193

194 **Author Contributions**

195 E.A.E-F., N.N.I., T. W., and N.C.K. designed the project, analyzed the data, and wrote
196 the manuscript. All authors discussed the results and commented on the manuscript.

197

198 **Competing interests**

199 The authors declare no competing interests.

200

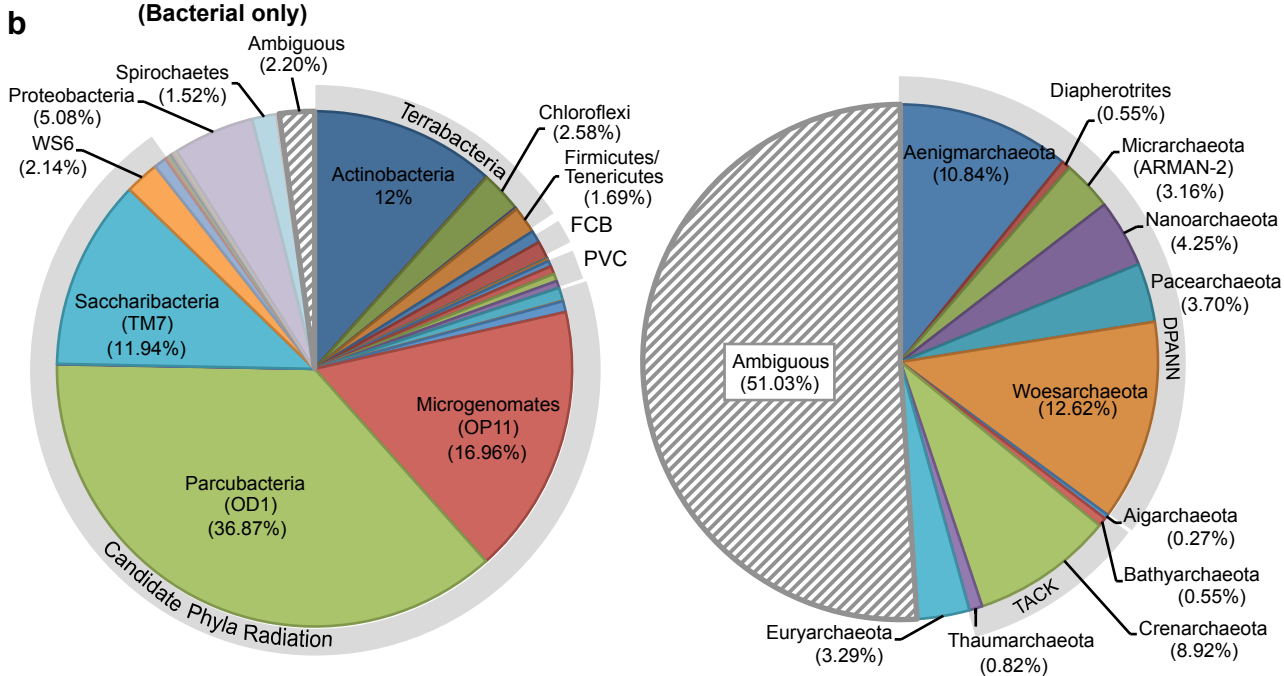
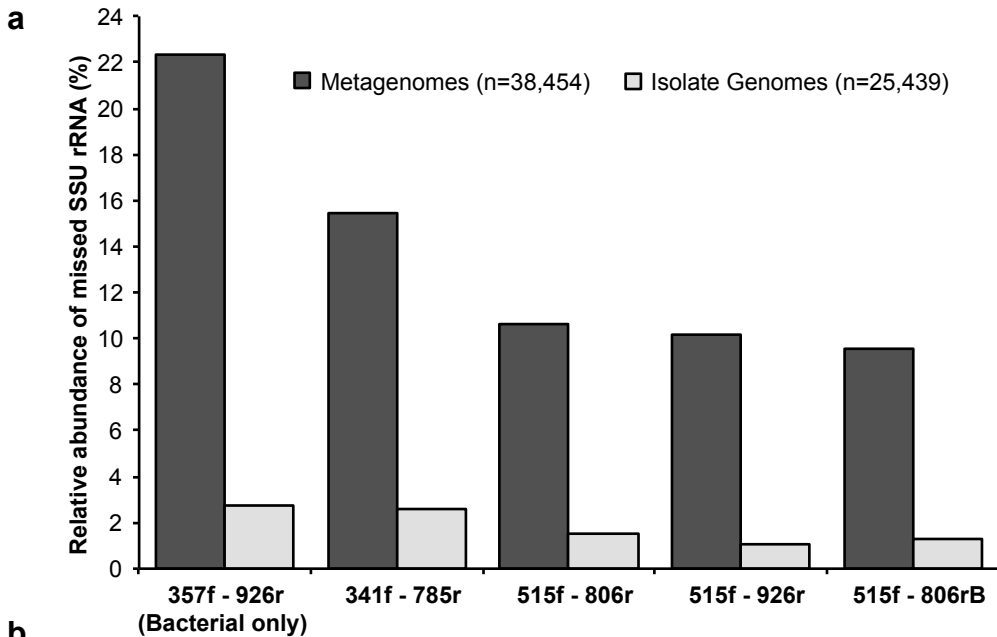
201 **Fig. 1. Primer fidelity for commonly used environmental PCR primers.** Primer sets
202 were evaluated against 38,454 metagenomic and 25,439 isolate genome SSU rRNA
203 gene sequences longer than 1000 bp. (a) Bar chart represents estimated percentage of

204 missed SSU rRNA gene sequences derived from metagenomes and isolate genomes.
205 Details available in Methods. For HMP primer 357F-926R, only bacterial sequences
206 were tested since they were not designed to include archaeal phylotypes. **(b)** Pie charts
207 represents SSU rRNA gene sequences from metagenomes that would likely be missed
208 with the widely used environmental primer set 515F-806R (overall primer weighted
209 score ≥ 1). The phylum-level affiliation for bacteria (left; 3,366 total missed sequences)
210 and archaea (right; 729 total missed sequences) are shown. Light gray bars along
211 outside of pie charts represent superphyla as follows: Terrabacteria (Actinobacteria,
212 Armatimonadetes (OP10), Chloroflexi, Cyanobacteria, Deinococcus-Thermus,
213 Firmicutes and Tenericutes); FCB (Bacteroidetes, Caldithrix, Chlorobi, Cloacimonetes
214 (WWE1), Fibrobacteres, Gemmatimonadetes, Kryptonia, Latescibacteria (WS3),
215 Marinimicrobia (SAR406), Zixibacteria); PVC (Planctomycetes, Verrucomicrobia,
216 Chlamydiae, Lentisphaerae, Omnitrophica (OP3)); Candidate Phyla Radiation
217 (Berkelbacteria (ACD58), CPR2, CPR3, Gracilibacteria (GN02/BD1-5), Kazan,
218 Microgenomates (OP11), Parcubacteria (OD1), Peregrinibacteria (PER), TM7, WS6,
219 WWE3); DPANN (Diapherotrites, Aenigmarchaeota, Nanoarchaeota,
220 Nanohaloarchaeota, Micrarchaeota (ARMAN-2), Parvarchaeota, Pacearchaeota,
221 Woesearchaeota); and TACK (Thaumarchaeota, Aigarchaeota, Crenarchaeota,
222 Korarchaeota, Bathyarchaeota).

223

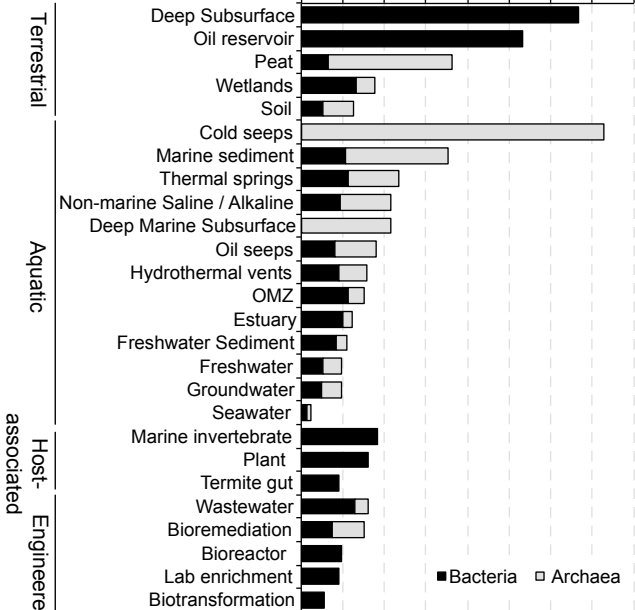
224 **Fig. 2. Habitat distribution of taxonomic novelty.** Distribution of bacterial (black) and
225 archaeal (gray) sequences recovered from metagenomes lacking definitive phylum-level
226 affiliation. A total of 654 bacterial sequences and 816 archaeal sequences were
227 analyzed based on environmental habitat metadata. Data was normalized based on the
228 total number of SSU rRNA gene sequences ≥ 1000 bp recovered from metagenomes
229 for each habitat.

230



Relative abundance of
ambiguous SSU rRNA (%)

0 5 10 15 20 25 30 35 40



■ Bacteria □ Archaea

Table 1. An estimate of taxonomic novelty derived from grouping metagenome-derived bacterial and archaeal SSU rRNA gene sequences into recognized phylum-level lineages.

	Relative abundance of metagenome SSU rRNAs (%)	
	Bacteria (n=34,596)	Archaea (n=3,858)
Cultured Phyla	83.53	46.73
Uncultured Phyla	14.58	32.12
Unaffiliated with Current Phyla	1.89	21.15