

# UC San Diego

## UC San Diego Previously Published Works

### Title

Improving metabarcoding taxonomic assignment: A case study of fishes in a large marine ecosystem

### Permalink

<https://escholarship.org/uc/item/9v70419c>

### Journal

Molecular Ecology Resources, 21(7)

### ISSN

1755-098X

### Authors

Gold, Zachary  
Curd, Emily E  
Goodwin, Kelly D  
[et al.](#)

### Publication Date

2021-10-01

### DOI

10.1111/1755-0998.13450

Peer reviewed

# MOLECULAR ECOLOGY RESOURCES

## Improving Metabarcoding Taxonomic Assignment: A Case Study of Fishes in a Large Marine Ecosystem

Journal:	<i>Molecular Ecology Resources</i>
Manuscript ID	MER-21-0082.R1
Manuscript Type:	Resource Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Gold, Zachary; UCLA, Ecology and Evolutionary Biology Curd, Emily; UCLA, Ecology and Evolutionary Biology Goodwin, Kelly; Atlantic Oceanographic and Meteorological Laboratory, Ocean Chemistry and Ecosystem Division; Southwest Fisheries Science Center Choi, Emma; University of California San Diego Scripps Institution of Oceanography Frale, Benjamin; University of California San Diego Scripps Institution of Oceanography Thompson, Andrew; Southwest Fisheries Science Center Walker Jr, Harold J; University of California San Diego Scripps Institution of Oceanography Burton, Ronald; University of California San Diego Scripps Institution of Oceanography Kacev, Dovi; University of California San Diego Scripps Institution of Oceanography Martz, Lucas Barber, Paul; UCLA, Ecology and Evolutionary Biology
Keywords:	metabarcoding, MiFish primers, California Current Large Marine Ecosystem, eDNA, Environmental DNA, reference database

**Improving Metabarcoding Taxonomic Assignment:  
A Case Study of Fishes in a Large Marine Ecosystem**

12S Taxonomic Assignment Performance

Zachary Gold<sup>1\*</sup>, Emily E. Curd<sup>1</sup>, Kelly D. Goodwin<sup>3</sup>, Emma S. Choi<sup>2</sup>, Benjamin W. Frable<sup>2</sup>,  
Andrew R. Thompson<sup>4</sup>, H. J. Walker, Jr.<sup>2</sup>, Ronald S. Burton<sup>2</sup>, Dovi Kacev<sup>2</sup>, Lucas D. Martz<sup>2</sup>,  
Paul H. Barber<sup>1</sup>

<sup>1</sup> Department of Ecology and Evolutionary Biology, University of California Los Angeles, Los Angeles, California, USA

<sup>2</sup> Scripps Institution of Oceanography, University of California San Diego, La Jolla, California, USA

<sup>3</sup> Ocean Chemistry and Ecosystems Division, Atlantic Oceanographic and Meteorological Laboratory, National Oceanic and Atmospheric Administration, stationed at Southwest Fisheries Science Center, La Jolla, California, USA

<sup>4</sup> Southwest Fisheries Science Center, National Oceanic and Atmospheric Administration, La Jolla, California, USA

\*Corresponding Author; email: [zjgold@g.ucla.edu](mailto:zjgold@g.ucla.edu)

## 1 **ABSTRACT**

2 DNA metabarcoding is an important tool for molecular ecology. However, its effectiveness  
3 hinges on the quality of reference sequence databases and classification parameters employed.  
4 Here we evaluate the performance of MiFish *12S* taxonomic assignments using a case study of  
5 California Current Large Marine Ecosystem fishes to determine best practices for  
6 metabarcoding. Specifically, we use a taxonomy cross-validation by identity framework to  
7 compare classification performance between a global database comprised of all available  
8 sequences and a curated database that only includes sequences of fishes from the California  
9 Current Large Marine Ecosystem. We demonstrate that the curated, regional database provides  
10 higher assignment accuracy than the comprehensive global database. We also document a  
11 tradeoff between accuracy and misclassification across a range of taxonomic cutoff scores,  
12 highlighting the importance of parameter selection for taxonomic classification. Furthermore, we  
13 compared assignment accuracy with and without the inclusion of additionally generated  
14 reference sequences. To this end, we sequenced tissue from 597 species using the MiFish *12S*  
15 primers, adding 252 species to GenBank's existing 550 California Current Large Marine  
16 Ecosystem fish sequences. We then compared species and reads identified from seawater  
17 environmental DNA samples using global databases with and without our generated references,  
18 and the regional database. The addition of new references allowed for the identification of 16  
19 additional native taxa representing 17.0% of total sequence reads from eDNA samples, including  
20 species with vast ecological and economic value. Together these results demonstrate the  
21 importance of comprehensive and curated reference databases for effective metabarcoding and  
22 the need for locus-specific validation efforts.

23 **KEYWORDS:** metabarcoding, MiFish primers, California Current Large Marine Ecosystem,  
24 eDNA, environmental DNA, reference database

25

## 26 **INTRODUCTION**

27 Metabarcoding is a process in which multiple species are identified from bulk DNA (e.g.  
28 homogenized gut contents, settlement tile scrapings, etc.) or environmental samples (Bohmann  
29 et al., 2014; Deiner et al., 2017; Taberlet, Coissac, Pompanon, et al., 2012). Metabarcoding is  
30 increasingly used to study marine ecosystems as the ability to sequence tens to hundreds of  
31 millions of reads in a single sequencing run allows the development of novel research questions,  
32 including species mapping, biomonitoring, gut content analyses, and population genomics, all of  
33 which aid understanding of the ecology of marine ecosystems (Baetscher et al., 2019; Closek et  
34 al., 2019; Goodwin et al., 2017; Guo, 2017; Kelly, Port, Yamahara, Martone, et al., 2014;  
35 Sanders et al., 2015; Thompson et al., 2017; Yamahara et al., 2019). In particular, metabarcoding  
36 of environmental DNA (eDNA), freely associated DNA obtained from environmental samples, is  
37 an increasingly attractive approach for marine ecosystem characterization because it can detect a  
38 broad range of diversity from a single liter of seawater, and has the potential to transform marine  
39 biomonitoring efforts (Kelly, Port, Yamahara, Martone, et al., 2014).

40 Metabarcoding typically employs PCR amplification and sequencing of a target gene  
41 (Goodwin et al., 2017) followed by comparison of these sequences to a database of known  
42 reference sequences to identify species present in the sample (Taberlet, Coissac, Hajibabaei, et  
43 al., 2012). Incomplete databases cannot identify all species present, leading to a lack of  
44 assignment despite the actual detection and capture of the sequences, potentially biasing the

45 interpretation of results (Boyer et al., 2016; Deiner et al., 2017; Machida et al., 2017). Thus,  
46 building complete and accurate reference databases is paramount to the success of molecular  
47 ecology monitoring efforts (Schenekar et al., 2020).

48         One approach for maximizing metabarcoding taxonomic assignment is to compare query  
49 sequences to a global database of archived sequences (Camacho et al., 2009; Edgar, 2018b).  
50 Global databases, such as GenBank, include nearly all publicly available sequences for specific  
51 barcode loci and are thus inherently comprehensive (Benson et al., 2018). However, the  
52 inclusion of reference barcodes from non-target or biologically irrelevant species may potentially  
53 bias taxonomic assignment algorithms (Curd et al., 2019). This issue is particularly problematic  
54 for lowest common ancestor taxonomic assignment methods that make inherent assumptions that  
55 each best sequence alignment is equally valid, irrespective of the geographic distributions and  
56 ecologies of these taxa (Curd et al., 2019; Gao et al., 2017), potentially leading to assignments of  
57 biologically implausible species. This problem can be compounded by the occurrence of mis-  
58 annotated sequences, a well-known problem in global reference databases (Heller et al., 2018;  
59 Leray et al., 2019; Nobre et al., 2016; Wakeling et al., 2019).

60         An alternative approach to using global databases for taxonomic classification is to  
61 employ a curated reference database that includes only appropriately annotated sequences for  
62 taxa that occur in a given region (Macheriotou et al., 2019; Poloczanska et al., 2013; Richardson  
63 et al., 2018). However, the inclusion or exclusion of barcodes from a reference database can  
64 affect metabarcoding taxonomic assignments (Macheriotou et al., 2019; Poloczanska et al., 2013;  
65 Richardson et al., 2018), yet few studies systematically addressed this problem (Bergsten et al.,  
66 2012; Stoeckle et al., 2020). As such, it is currently unclear whether global or regional reference

67 databases produce more accurate taxonomic assignments. Systematically quantifying error and  
68 bias associated with global and curated database is essential to identifying best practices for  
69 metabarcoding taxonomic assignment.

70         Critical to such assessments are methods that validate taxonomy prediction and evaluate  
71 the sensitivity to bioinformatic and database parameters. One key method for comparing the  
72 performance of taxonomic classification across different reference databases or classification  
73 parameters is the taxonomy cross-validation by identity (TAXXI) framework (Edgar, 2018a).  
74 The TAXXI framework is executed by using a reference database with known taxonomic  
75 identities that is split into test and training sets and then assigning taxonomy to the training set  
76 using the test set. The TAXXI framework can then be applied to allow taxonomic assignment  
77 performance to be compared across different metabarcodes, reference databases, and different  
78 assignment parameters.

79         Critically, TAXXI approaches allow for comparing the performance of bioinformatic  
80 pipelines within and across loci, including informing the proper selection of classifier parameters  
81 for a given metabarcoding locus (Boyer et al., 2016; Machida et al., 2017). Taxonomic  
82 assignments made by metabarcoding classifiers are particularly influenced by taxonomic cutoff  
83 scores (e.g., exact alignment match or 97% identity threshold) (Edgar, 2018c, 2018a). Using this  
84 cross-validation approach to evaluate the performance of taxonomic assignments for *16S* and  
85 fungal *ITS* metabarcoding loci across a range of classification parameters revealed that percent  
86 identities below 95% had poor classification performance (Edgar 2018a), and highlighted key  
87 tradeoffs between assignment confidence and taxonomic resolution (Edgar, 2018c, 2018a).  
88 Frequently, attempts to balance confidence–resolution tradeoffs leads to the selection of

89 conservative taxonomic cutoff scores to avoid over-classification errors (Alberdi et al., 2018;  
90 Camacho et al., 2009; Port et al., 2015; Siegwald et al., 2017; Wood & Salzberg, 2014).  
91 However, parameter selection is rarely systematically evaluated across different taxonomic  
92 groups or metabarcoding loci, inadvertently leading to poorer quality taxonomic assignments  
93 (Curd et al., 2019; Edgar, 2018a, 2018c). Importantly, the few studies that explored classification  
94 parameter performance across metabarcoding loci found that a “one size fits all” approach (e.g.,  
95 97% identity threshold) is inappropriate across different metabarcoding loci (Curd et al., 2019;  
96 Edgar, 2018c, 2018a). Thus, evaluating the performance of taxonomic assignments across a  
97 range of cutoff scores for a given metabarcoding target is important for maximizing the accuracy  
98 of metabarcoding efforts (Balakirev et al., 2017; Bokulich et al., 2018; Hassanin et al., 2010).

99       Using the TAXXI framework, Curd *et al.* (2019) compared the performance of reference  
100 databases for taxonomic assignment, demonstrating the utility of custom reference libraries. The  
101 *Creating Reference libraries Using eXisting tools (CRUX)* module of the *Anacapa Toolkit*  
102 constructs custom reference databases by querying public sequence archives based on primer sets  
103 defined by the user. Curd et al. (2019) showed that *CRUX*-generated custom reference databases  
104 were more comprehensive and provided improved taxonomic assignment compared to  
105 previously published *COI* reference databases [Midori (Machida et al., 2017) and CO-Arbitrator  
106 (Heller et al., 2018)], yielding results nearly equal to heavily curated reference databases for *16S*  
107 [SILVA (Quast et al., 2012)] and *12S* [MitoFish (Sato et al., 2018)] metabarcodes. The TAXXI  
108 framework thus provides a critical set of tools to evaluate the performance of taxonomic  
109 assignment across classification parameters and reference databases for any metabarcoding locus  
110 of interest.



111           The MiFish Universal Teleost and MiFish Elasmobranch primer sets (Miya et al., 2015)  
112 target the same portions of the mitochondrial *12S* RNA gene, but differ by a few critical base  
113 pairs on the forward primer. These metabarcodes are vertebrate specific, provide species-level  
114 resolution for many fishes, and are well suited to short read-length next-generation DNA  
115 sequencing, such as Illumina platforms (Collins et al., 2019; Jo et al., 2017; Miya et al., 2015;  
116 Valsecchi et al., 2019). As such, they are becoming the standard barcode locus for marine  
117 vertebrate metabarcoding studies (Bista et al., 2017; Closek et al., 2019; Miya et al., 2015;  
118 Thomsen et al., 2016; Valsecchi et al., 2019; Yamamoto et al., 2017). However, *12S* fish  
119 reference databases are relatively incomplete compared to traditional barcoding loci, such as the  
120 655 bp region of the mitochondrial Cytochrome Oxidase I (*COI*) gene (Ardura et al., 2013; Duke  
121 & Burton, 2020; Hastings & Burton, 2008; Ward et al., 2009). For example, there is an extensive  
122 *COI* barcode database of fishes of the California Current Large Marine Ecosystem (Hastings &  
123 Burton, 2008) that, according to the MitoHelper query of the MitoFish database (accessed April  
124 2021) includes 878 of 1,144 (76.7%) species (Iwasaki et al., 2013; Lim & Thompson, 2021)]  
125 facilitating numerous recent metabarcoding studies (Closek et al., 2019; Djurhuus et al., 2020;  
126 Pitz et al., 2020). However, there are relatively few reference *12S* sequences that overlap with  
127 the MiFish primer sets, limiting the utility of *12S* metabarcoding approaches in this region.

128           The California Current Large Marine Ecosystem is a highly productive coastal ecosystem  
129 that extends approximately 3,000 km across most of the Northeast Pacific from Baja California,  
130 Mexico to British Columbia, Canada (Checkley Jr & Barth, 2009; Coleman, 2008; Ekstrom,  
131 2009; Koslow & Davison, 2016). This large marine ecosystem has enormous regional and global  
132 importance (Ekstrom, 2009; Sherman, 1991; Wells et al., 2020), driving an ocean economy

133 valued at over \$56 billion USD, employing over 675,000 people (Block et al., 2011; Koslow &  
134 Davison, 2016; NMFS, 2017) and supporting food security of the region. The California Current  
135 Large Marine Ecosystem also plays a vital role in the cultures and traditional practices of coastal  
136 North American tribes and First Nations by supporting species such as Pacific salmon  
137 (*Oncorhynchus* spp.), orcas (*Orcinus orca*), eulachon (*Thaleichthys pacificus*), and abalone  
138 (*Haliotis* spp.) (Armstrong, 2017; Braje et al., 2017; Brooks et al., 2012; Lepofsky et al., 2017;  
139 Norgaard, 2019; Wadewitz, 2012).

140       Unfortunately, this ecosystem is increasingly facing numerous threats including  
141 overexploitation (Koslow & Davison, 2016), ocean acidification and hypoxia (Chan et al., 2008;  
142 Crozier et al., 2019; Hofmann et al., 2014; Samhoury et al., 2017), pollution (Good et al., 2020;  
143 Halpern et al., 2009), and climate change induced marine heat waves (Rogers-Bennett & Catton,  
144 2019; Santora et al., 2020). Metabarcoding has the power to address many critical management  
145 questions in this region, ranging from shifting species distributions, effectiveness of marine  
146 protected areas, and seasonal patterns of larval fish recruitment, among others (Duke & Burton,  
147 2020; Kelly, Port, Yamahara, Martone, et al., 2014; Port et al., 2015). However, the ability of  
148 metabarcoding efforts to address these important questions hinges on the availability of  
149 comprehensive reference databases and appropriate methods of bioinformatic analysis.

150       To improve the utility of *12S* metabarcoding of marine fishes for the California Current  
151 Large Marine Ecosystem and to address larger questions regarding the impact of bioinformatic  
152 processes on taxonomic classification, we 1) generated and contributed 741 additional MiFish  
153 *12S* sequences representing 597 fish species to global sequence databases; 2) used these  
154 additional sequences to create a reference database curated specifically for the California Current

155 Large Marine Ecosystem; 3) compared the performance of taxonomic assignments made by this  
156 regional curated reference database to those made by global marine vertebrate reference  
157 databases; and 4) assessed the effect of classifier parameters on phylum through species level  
158 assignments of MiFish *I2S* sequences to identify optimal locus-specific bioinformatic  
159 parameters.

160

## 161 **METHODS**

### 162 **Reference Barcode Generation from Fish Tissue Samples**

163 To generate a more complete *I2S* barcode reference database for California Current Large  
164 Marine Ecosystem fishes, we assembled a list of the 1,144 marine teleost and elasmobranch  
165 species that occur in this system (Allen & Horn, 2006; Froese & Pauly, 2010; Hastings &  
166 Burton, 2008; Love, & Passarelli, 2020) (Table S1). From this list, we acquired 741 ethanol-  
167 preserved voucher specimens representing 597 species (Table S1, Table S2) from the Scripps  
168 Institution of Oceanography Marine Vertebrate Collection at the University of California San  
169 Diego. DNA was extracted from each tissue sample using a Chelex 100 extraction method  
170 (Walsh, Metzger, & Higuchi, 1991), as described in the Supplemental Methods. We amplified all  
171 teleost DNA extracts (n=701) using the MiFish Universal Teleost Primers (Miya et al., 2015),  
172 and all elasmobranchs (n=55) using the MiFish Elasmobranch primers (Miya et al., 2015)  
173 following the thermocycler profile of Curd et al., (2019) (Table S3). We Sanger sequenced  
174 purified amplicons (see Supplemental Methods for details), and aligned and trimmed forward  
175 and reverse sequences in Sequencher version 5.4.6 (Nishimura, 2000). We used *R* package *taxize*  
176 (version 0.9.99) (Chamberlain & Szöcs, 2013) to synonymize taxonomic names of all vouchered

177 specimens and GenBank. We then checked the accuracy of generated reference barcodes by  
178 building a UPGMA phylogenetic tree of all reference sequences and California Current Large  
179 Marine Ecosystem fishes using *phangorn* (2.5.5). In addition, we queried each sequence using  
180 *blastn* (Camacho et al., 2009) and removed any sequence that did not cluster or align to known  
181 taxonomic lineages (data available at <https://doi.org/10.5068/D1H963>). The resulting *I2S*  
182 reference barcodes were deposited into GenBank (SAMN19289093–SAMN19289810; Table  
183 S2).

184

## 185 **Reference Database Creation**

186 To test variation in taxonomic assignment among reference databases, we generated three  
187 distinct reference sequence databases: “CRUX-GenBank”, “global”, and “regional” (Table 1 and  
188 Table 2). CRUX-GenBank is a custom *I2S* reference database generated using Creating  
189 Reference libraries Using eXisting tools (*CRUX*) module of the Anacapa Toolkit to query  
190 GenBank for reference barcodes conducted with standard search parameters (Benson et al., 2018;  
191 Curd et al., 2019) and MiFish Universal *I2S* sequences (Table S1) as the user-defined primers.  
192 Briefly, we created this reference database by running *in silico* PCR (Ficetola et al., 2010) on the  
193 European Molecular Biology Laboratory (EMBL) standard nucleotide database (Stoesser et al.,  
194 2002) to generate a seed library of *I2S* references. Next, we used *blastn* (Camacho et al., 2009)  
195 to capture reference barcodes without included primer sequences and to query the seed database  
196 against the NCBI non-redundant nucleotide database (Gold, 2020; Pruitt et al., 2005; sequences  
197 downloaded in October 2019). The resulting *blastn* hits were de-replicated by retaining only the  
198 longest version of each sequence and taxonomy for each accession was retrieved using

199 *Entrez-qiime* (Baker, 2016). The resulting set of reference sequences in the CRUX-GenBank  
200 database included any GenBank reference barcodes that *in silico* amplified to the MiFish *I2S*  
201 primers at the time of this analysis.

202         We created the global database to evaluate whether increasing database completeness  
203 improves taxonomic assignment. To create the global database, we supplemented the CRUX-  
204 GenBank database with 741 additional California Current Large Marine Ecosystem fish *I2S*  
205 barcodes generated for this study (Table S2). Thus, the global database includes all fish *I2S*  
206 reference sequences available at the time of download. From this global database, we created the  
207 regional database, including only *I2S* sequences of fishes known to occur in the California  
208 Current Large Marine Ecosystem. We created this database to specifically test whether databases  
209 curated to specific ecosystems enhance taxonomic assignment performance relative to more  
210 comprehensive databases (“global”). Because of the high degree of similarity between the  
211 MiFish Universal and Elasmobranch loci and the flexibility built into *CRUX*, a single *CRUX*  
212 generated *I2S* reference database performs well for both markers (Curd et al., 2019), so we did  
213 not create separate teleost and elasmobranch databases. Additionally, because the MiFish primer  
214 set amplifies nearly all vertebrate taxa (Miya et al., 2015; Valsecchi et al., 2019), the global  
215 database include teleosts, elasmobranchs, mammals, reptiles, amphibians, birds, etc. All  
216 databases are available at <https://doi.org/10.5068/D1H963>.

217

## 218 **Taxonomy cross-validation by identity comparisons**

219 We implemented the taxonomy cross-validation by identity (TAXXI) framework developed by  
220 (Edgar, 2018a) to 1) compare taxonomic assignment performance metrics for global versus

221 regional reference databases, 2) determine the resolution of taxonomic assignments for all  
222 available MiFish barcodes in the global database, and 3) understand the performance of the  
223 MiFish barcode across taxonomic classifier cutoff scores. Although we use three databases  
224 (global, CRUX-GenBank and regional) on our test dataset below, we did not include the CRUX-  
225 GenBank database in taxonomic cross validation comparisons because the global database  
226 contains all these sequences.

227 The TAXXI analyses were implemented using scripts from Curd et al. (2019) which  
228 adapted TAXXI to the *Anacapa Toolkit* (<https://drive5.com/taxxi/doc/index.html> and  
229 <https://github.com/limey-bean/Anacapa>). We conducted taxonomic assignments using the  
230 *Anacapa Toolkit classifier* which implements the Bayesian Lowest Common Ancestor (BLCA)  
231 classifier (Gao et al., 2017) modified to incorporate sequences from *Bowtie2* (Langmead &  
232 Salzberg, 2012). In brief, amplicon sequence variants (ASVs; exact unique sequences  
233 dereplicated from generated metabarcoding data) are first aligned to reference barcodes using  
234 *Bowtie2* retaining the top 100 alignments. Then the BLCA classifier conducts multiple sequence  
235 alignment for each query ASV to inform a weighted Bayesian posterior probability of taxonomic  
236 assignment. Taxonomy is then ultimately assigned based on the lowest common ancestor of the  
237 total weighted reference database matches; reliability is evaluated through bootstrap confidence  
238 scores which are analogous to percent identity metrics provided by other metabarcoding  
239 classifiers (Gao et al., 2017; See Curd et al. 2019 for full description).

240 We evaluated taxonomic assignment performance by comparing the following metrics: 1)  
241 true positive rate – the number of correct taxonomic assignments divided by the total  
242 opportunities for correct classification, 2) over-classification rate - the number of assignments

243 incorrectly made to additional lower taxonomic ranks divided by the total opportunities to make  
244 an over-classification error, 3) under-classification rate - the number of assignments incorrectly  
245 made to fewer taxonomic ranks divided by the total opportunities to make an under-classification  
246 error, 4) misclassification rate - the number of assignments incorrectly predicted divided by the  
247 opportunities for correct classification, and 5) accuracy - the number of correct assignments  
248 divided by the taxonomic assignment opportunities for which correctness can be determined (R.  
249 C. Edgar, 2018a). The 6) sensitivity was calculated as the true positive rate / (true positive rate +  
250 under-classification rate) as under-classification is analogous to a false negative rate. The 7)  
251 specificity was calculated as 1 - (misclassification rate + over-classification rate) as the  
252 combination of the misclassification rate and over-classification rate is analogous to the false  
253 positive rate.

254

### 255 **Taxonomic Resolution of the MiFish *I2S* primer**

256 To provide insights into which fishes can be resolved to species level using the MiFish *I2S*  
257 primer set, we conducted TAXXI comparisons using the global database as both the test and  
258 training database to assign taxonomy to itself. We then calculated the seven taxonomic  
259 assignment metrics described above. Additionally, we identified families and genera of fishes for  
260 which the MiFish *I2S* locus performed poorly, defined as frequently failing to assign species  
261 level identification. Although all vertebrate sequences in the global database were used in the  
262 taxonomic cross validation, only results for fishes are discussed here.

263

## 264 **Regional vs. global reference databases**

265 To compare the relative ability of regional versus global reference databases to accurately assign  
266 taxonomy, we conducted two additional TAXXI comparisons using the reference databases  
267 created for this study. First, we used the global reference database as a training database to assign  
268 taxonomy to the regional reference database that only contained sequences for fishes known  
269 from the California Current Large Marine Ecosystem. Second, we used the regional reference  
270 database as both the test and training database to assign taxonomy against itself. The taxonomic  
271 assignments made by the global and regional reference databases were compared across the  
272 taxonomic assignment metrics described above.

273

## 274 **Effect of Bootstrap Confidence Scores on Taxonomic Assignment**

275 To understand the performance of the MiFish barcode across a range of taxonomic classifier  
276 cutoff scores, we repeated each of the three TAXXI analyses described above (global-regional,  
277 regional-regional, global-global) using bootstrap confidence cutoff scores of 40, 50, 60, 70, 80,  
278 90, 95, and 100. We then evaluated the effect of bootstrap confidence cutoff scores across the  
279 various taxonomic assignment metrics, as described above.

280

## 281 **eDNA Metabarcoding Case Study**

### 282 **Seawater Sample Collection, DNA Extraction, and Library Generation**

283 To specifically test the impact of *12S* database design on taxonomic assignment in real world  
284 applications, we compared the performance of the three databases in assigning taxonomy to  
285 existing eDNA sequence data as a test case. Briefly, we used MiFish *12S* metabarcoding



286 sequence data generated from three seawater samples collected from 10 m depth from three sites  
287 off eastern Santa Cruz Island, CA in 2017 that were part of a larger ecological study of  
288 biodiversity patterns within rocky reef ecosystems. These sequences were generated using  
289 standard eDNA collection, processing, and sequencing methods, as outlined in Gold et al.,  
290 (2021).

291 We processed this eDNA metabarcoding data three separate times using the *Anacapa*  
292 *Toolkit* (Curd et al., 2019), assigning taxonomy using the CRUX-GenBank, global, and regional  
293 reference databases (Table 2). We used the default *Anacapa Toolkit* parameters and a bootstrap  
294 confidence cutoff score of 60. We then examined the total number of ASVs and taxonomic ranks  
295 identified by each of the three reference databases. We also investigated differences in  
296 taxonomic assignment between single direction ASVs (comprised of forward- and reverse-only  
297 sequence reads) and merged ASVs (merged paired-end sequence reads) to understand the  
298 importance of full length vs. partial length sequences for taxonomic assignment (See  
299 Supplemental Results and Discussion).

300

## 301 **RESULTS**

### 302 **Generation of Novel Barcodes and 3 References Databases**

303 We generated 741 new *I2S* MiFish barcode sequences for 597 California Current Large Marine  
304 Ecosystem fishes (Table S1 and Table S2), 545 teleosts (bony fishes), 49 elasmobranchs  
305 (cartilaginous fishes), and 3 cyclostomatan (jawless fishes) (Table S2). This dataset includes 252  
306 that had no previous *I2S* reference barcodes (Table S1).

307 *CRUX* created a custom *I2S* database comprised of 14,066 taxa and 44,140 sequences  
308 with existing entries in GenBank. Adding the 741 novel sequences, above, resulted in a global  
309 database comprised of 14,321 species and 44,882 sequences. Restricting these sequences to only  
310 fishes from the California Current Large Marine Ecosystem resulted in a curated regional  
311 database that includes 706 out of 1,144 (61.7%) reference *I2S* barcodes from fishes known from  
312 this region. Excluding 382 species missing from the database that are rare in California (n=357)  
313 or not coastal (n=25), resulted in a total coverage of 92.7% of the 763 common coastal fishes in  
314 this region.

315

## 316 **Taxonomy Cross-validation by Identity Comparisons**

### 317 **Regional Versus Global Reference Database Comparisons**

318 The TAXXI quality metrics indicate that the regional reference database yielded more reliable  
319 taxonomy at genus and species ranks relative to the global reference databases across all  
320 bootstrap confidence scores; regional database species level accuracy ranged from 64.2-94.2%  
321 compared to 51.3-90.8% for the global database (Table 2 & Table S4; Figures S1 and S2). This  
322 difference was driven by higher misclassification and under-classification rates for the global  
323 reference databases. In particular, database misclassification rates were higher for the global  
324 compared to the regional reference database across all bootstrap confidence cutoff scores less  
325 than 60 (global reference database misclassification 1.8-4.5%, regional database  
326 misclassification rate 1.3-3.1%) (Table S4). Likewise, global reference database under-  
327 classification rates were higher than regional reference database under-classification rates across

328 all bootstrap confidence cutoff scores (global reference database under-classification 4.8-48.7%,  
329 regional database under-classification rate 2.8-35.8%).

330

### 331 **Taxonomic Resolution of the MiFish 12S primer**

332 Cross validation of the 44,896 sequences within the global database demonstrated that the  
333 MiFish primer set delivered 88.0% sensitivity [true positive rate / (true positive rate + under-  
334 classification rate)] and 98.2% specificity [1 - (misclassification rate + over-classification rate)] at  
335 a bootstrap cutoff score of 60 (Table 2, Table S5), providing species level taxonomic  
336 assignments to 6,762 fish species, genus level resolution to 923 fish species, family level  
337 assignments to 180 fish species, and class level assignments to 2 fish species while  
338 overclassifying 214 fish species (Table S5). While poor taxonomic resolution with the MiFish  
339 primer sets (e.g. assigned taxonomic rank above species) spanned a large number of genera and  
340 families, the genus *Sebastes* and families Cichlidae, Cyprinidae, and Pleuronectidae were  
341 particularly problematic (Figures 4 and 5). Of these, *Sebastes* and Pleuronectidae are highly  
342 prevalent within the California Current Large Marine Ecosystem. A full breakdown of  
343 taxonomic assignment resolution is provided in the Supplemental Results.

344

### 345 **Effect of Bootstrap Confidence Scores on Taxonomic Assignment**

346 Across all TAXXI comparisons, accuracy and true positive rates increased with decreasing  
347 bootstrap confidence cutoff scores (Figure 1, Figures S1 and S2, Table S4). Likewise, the  
348 proportion of species level assignments also increased with decreasing bootstrap confidence  
349 score (Figure 2, Figures S3 and S4). We also found that misclassification rates increased with

350 decreasing bootstrap confidence cutoff score, but at much lower rate (Figure 3, Figures S5 and  
351 S6). These results indicate a clear tradeoff between under-classification and misclassification  
352 across bootstrap confidence cutoff scores.

353

## 354 **eDNA Metabarcoding Example**

### 355 **Unassigned MiFish 12S ASVs**

356 The *Anacapa Toolkit* failed to assign taxonomy to 49.6% (169/341) of ASVs representing 24.5%  
357 (81,002/330,877) of all reads using all three reference databases investigated in this study (Table  
358 S6). Of the 169 unassigned ASVs, 16 were forward-only reads, and 153 were merged reads. To  
359 explore the origins of these unassigned reads, we used BLAST to query all GenBank sequences,  
360 revealing that 94.7% (160/169) of these ASVs aligned to marine prokaryotic and eukaryotic 16S  
361 sequences (Max Alignment Scores 87.9-475). Of these aligned ASVs, 85% (136/160) matched to  
362 uncultured sequences generated from marine metagenomic studies. 80.0% (128/160) of  
363 successfully aligned ASVs matched to bacterial barcodes including those from *Psychromonas*  
364 sp., *Photococcus caeruleum*, *Loktanella* sp., *Leucothrix* sp., and *Gimesia* sp., and cyanobacteria.  
365 A smaller fraction of assigned ASVs (18.8%; 30/160) best aligned to eukaryotic sequences  
366 including those from diatoms (e.g. *Nitzschia alba* and *Eucampia antarctica*) and other marine  
367 microalgae (e.g. Picobiliphytes, *Heterosigma akashiwo*, *Mesopedinella arctica*, and *Phacus*  
368 *warszewiczii*). Given that these 169 unassigned sequences were non-vertebrate, we excluded  
369 these ASVs from all subsequent comparisons. All remaining 172 ASVs were assigned to a class  
370 of vertebrates by at least one of the three reference databases used. Of these vertebrate ASVs, 58  
371 were merged, 107 were forward-only, and 7 were reverse only reads.

372

373 **Comparisons of CRUX-GenBank, Global, and Regional Reference Database Taxonomic**374 **Assignments**

375 The inclusion of additional reference barcodes increased the total number of ASVs and reads  
376 assigned to marine fishes resident in the California Current Large Marine Ecosystem (Tables 2 &  
377 Table S7). Importantly, the inclusion of novel voucher sequences within the global database  
378 resulted in species-level identification for 11 additional California Current Large Marine  
379 Ecosystem fishes including Kelp Bass (*Paralabrax clathratus*), California Moray (*Gymnothorax*  
380 *mordax*), Opaleye (*Girella nigricans*), Giant Kelpfish (*Heterostichus rostratus*), Ocean  
381 Whitefish (*Caulolatilus princeps*), and California Halibut (*Paralichthys californicus*) (Table S8).

382 Use of the regional database largely increased accuracy of taxonomic assignments. The  
383 regional database assigned an ASV to the Black Croaker (*Cheilotrema saturnum*) that was only  
384 assigned to the family Sciaenidae by the global database. Additionally, the regional database  
385 assigned one ASV as Bat Ray (*Myliobatis californica*) and another as Jack Mackerel (*Trachurus*  
386 *symmetricus*), species native to the California Current Large Marine Ecosystem, that the global  
387 database assigned to the non-native species, Common Eagle Ray (*Myliobatis aquila*) and Rough  
388 Scad (*Trachurus lathami*), respectively. However, the regional reference database failed to  
389 resolve the taxonomy of one ASV that the global database assigned to the family of Delphinidae.

390

391 **DISCUSSION**

392 Taxonomic assignment in metabarcoding studies typically employ large public sequence  
393 databases such as GenBank or Barcode of Life (Leray & Knowlton, 2015; Schenekar et al.,

394 2020; Stat et al., 2017), or databases that are curated to specific barcoding markers or taxonomic  
395 groups without consideration of species distributions (e.g. Curd et al. 2019). However,  
396 systematic comparison of these approaches to a curated, region-specific reference database  
397 shows that the region-specific database outperforms the global databases in metabarcoding  
398 taxonomic assignment (Table 1). Accuracy of eDNA metabarcoding only improved by including  
399 GenBank sequences from fishes native to the California Current Large Marine Ecosystem and  
400 supplementing these sequences with additional reference barcodes. Furthermore, examination of  
401 taxonomic assignment over a range of bootstrap cutoff scores revealed key tradeoffs, with lower  
402 bootstrap confidence cutoffs yielding more accurate species assignment, but at the cost of higher  
403 misclassification rates. Combined, these results highlight the importance of reference database  
404 and bootstrap cutoff selection in obtaining the best results from metabarcoding studies.

405 In a test dataset for fish eDNA extracted from seawater collected from three sites on  
406 Santa Cruz Island, the regional database performed the best. The regional database identified 16  
407 additional ASVs to species not identified by the CRUX-GenBank database, and an additional 3  
408 fishes that were misidentified by the global database (Table 2). Higher accuracy with increased  
409 database completeness echoes previous research on the importance of complete reference  
410 databases in metabarcoding (Leray et al., 2012; Machida et al., 2017), and greatly improves the  
411 utility of eDNA for monitoring the California Current Large Marine Ecosystem.

412 Although the *I2S* barcodes and reference databases tested here performed well with  
413 regard to annotating fish species (e.g., 91.3% sensitivity and 98.3% specificity across MiFish  
414 reference barcodes), almost half of the ASVs and a quarter of all reads generated in our eDNA  
415 test datasets were not assigned to any fish reference barcode (Table 2). While other

416 metabarcoding studies report similar levels of unassigned taxa (Leray & Knowlton, 2017) and  
417 others have encountered this issue (Goodwin, personal communication), this issue isn't widely  
418 reported in the literature, particularly considering the popularity of the MiFish primers. Further  
419 investigation showed that the vast majority of unassigned ASVs were uncultured bacteria *16S*  
420 loci (Table S6) derived from marine shotgun sequencing metagenomic studies (Bork et al.,  
421 2015). This result highlights that the MiFish Teleost *12S* primer set, while extremely useful for  
422 targeting vertebrate *12S* loci, can also amplify non-target *16S* genes, raising the possibility that  
423 non-target amplification may at best result in lower returns of target sequences, and at worst  
424 artificially increase estimates of fish diversity.

425

## 426 **Importance of Regional reference databases**

427 Given that increased reference database completeness increases the ability to assign ASV's to  
428 species (Table 2), it is logical to assume that databases with more taxonomic coverage are  
429 universally better (Curd et al., 2019). However, our results suggest an unexpected trade-off  
430 between greater diversity of barcodes and ecologically informed taxonomic assignment. For  
431 example, using only the regional database specific to California Current Large Marine  
432 Ecosystem marine fishes, we identified important native taxa like Black Croaker (*Cheilotrema*  
433 *saturnum*) and Bat Ray (*Myliobatis californica*) in eDNA isolated from seawater samples.  
434 However, while the global database contained the largest total number of barcodes, including all  
435 taxa in the regional database, Black Croaker was not identified and Bat Ray was inconsistently  
436 identified across multiple ASVs. The global database failed to identify Black Croaker due to the  
437 high similarity of *12S* barcode sequences within the Family Sciaenidae, specifically within the

438 clade that includes *Cheilotrema*, a genus native to California, as well as *Equetus* and *Pareques*,  
439 non-native coral reef-associated genera (Table S8). Similarity of barcode sequences also explains  
440 the loss of taxonomic resolution in *Myliobatis*.

441 By excluding highly similar non-native *I2S* barcodes, the database curated for the region  
442 of interest provided more accurate species-level assignments and far fewer under-classifications  
443 and misclassifications, demonstrating that a database comprised of only local taxa is preferred to  
444 maximize identification of local species. Yet, this improvement was not universal. For example,  
445 the regional database failed to classify one ASV belonging to the family Delphinidae that was  
446 identified by both the CRUX-GenBank and global databases. This result stems from the regional  
447 database being specific to California Current Large Marine Ecosystem fishes, and could thus not  
448 identify a marine mammal. This shortcoming easily could be overcome, however, by appending  
449 the regional database with barcodes for other marine-associated vertebrate taxa of regional  
450 management interests (Valsecchi et al., 2019). An alternative and taxon agnostic approach  
451 currently employed by the co-authors is to conduct taxonomic assignments twice. First,  
452 taxonomic assignments are conducted using a regional reference database to get the best  
453 taxonomic assignment for focal taxa of interest, and second using a global reference database to  
454 identify as many remaining unidentified ASVs as possible (Gold et al., 2021). We did not  
455 directly report the results of the two-step taxonomic assignment method here as the only  
456 difference between this approach and the taxonomic assignments made using the regional  
457 database alone is the additional assignment of the single Delphinidae ASV.

458 These results highlight the tradeoff between identifying local species from clades with  
459 little genetic variation and providing taxonomic coverage across a broad range of species. As



460 such, researchers need to identify their research priorities when deciding on which reference  
461 databases to use, with a particular focus on defining the scope of the target taxa. Future work  
462 could alleviate this tradeoff by building bioinformatic pipelines that prioritize assignments to a  
463 reference set of resident species, perhaps by including information on species ranges and sample  
464 locations in the assignment algorithm. However, an advantage of the two-step approach outlined  
465 above is that it allows for eDNA studies to address specific ecological questions without having  
466 a specific target list in mind. This approach is particularly important for eDNA studies which  
467 directly test for the presence of invasive species or range shifts associated with climate change  
468 (Bohmann et al., 2014; Klymus et al., 2017).

469

## 470 **Importance of Taxonomic Cutoff Scores**

471 Taxonomic cutoff scores, or percent identity, strongly influenced taxonomic assignments (Edgar,  
472 2018c). Patterns for the MiFish *12S* locus showed a similar pattern with higher true positive and  
473 misclassification rates and lower under-classification rates at lower bootstrap confidence cutoff  
474 scores (Edgar, 2018a). These results highlight a key tradeoff between under-classification and  
475 misclassification for metabarcoding taxonomic assignment, and demonstrate that the decision of  
476 which taxonomic cutoff score can strongly influence results (Edgar, 2018a). Lower bootstrap  
477 confidence cutoffs ensure a higher overall accuracy in species-level identification but come at  
478 the cost of higher misclassification rates to an incorrect species-level assignment.

479 Our results suggest that a TAXXI bootstrap confidence cutoff score of 60 provides a  
480 balance between maximizing species-level assignment accuracy (89.7%, global reference  
481 database) while minimizing misclassification rates (1.7%, global reference database), matching

482 the general findings of Curd et al. (2019). However, in instances in which metabarcoding results  
483 may influence management or health decisions with substantial legal or economic ramifications  
484 (i.e., detection of an endangered or invasive species or discriminating a putative disease causing  
485 microbe) a misclassification error may be valued as a far less desirable outcome than an under-  
486 classification error (Bohmann et al., 2014; Lodge et al., 2012; Wakeling et al., 2019). In such  
487 cases, results indicate that there isn't one single bootstrap confidence cutoff score that  
488 completely ameliorates these tradeoffs (Figure 1).

489         Given that previous work demonstrates that results may not be consistent across loci  
490 (Curd et al., 2019), we can only generalize our results to the MiFish *I2S* primer set. Determining  
491 confidence–resolution tradeoffs in other widely used primer sets will be fundamental for  
492 effectively interpreting metabarcoding results from those loci. Combining the capabilities of  
493 *CRUX* with the TAXXI framework provides a critical set of tools to both generate and evaluate  
494 the performance of a range of metabarcoding loci and reference databases (Table 1; Curd et al.,  
495 2019; Edgar, 2018a), facilitating such studies. Given the growing number of metabarcoding  
496 applications across a broad range of ecosystems and taxa (Curd et al., 2019; Deiner et al., 2017;  
497 Edgar, 2018a), assessing the performance of barcoding markers in the taxonomic group of  
498 interest is critical.

499

## 500 **Importance of Complete Reference Databases**

501 Previous eDNA metabarcoding efforts in the California Current Large Marine Ecosystem report  
502 poor species-level identification and frequent taxonomic assignment to non-native sister taxa  
503 (Closek et al., 2019; Kelly, Port, Yamahara, & Crowder, 2014; Port et al., 2015). For example,

504 an eDNA metabarcoding study in Southern California (Curd et al., 2019) assigned multiple *I2S*  
505 ASVs to *Girella simplicidens*, the Gulf Opaleye, a fish that does not occur in California Current  
506 Large Marine Ecosystem coastal waters (Froese & Pauly, 2010; Love & Passarelli, 2020). This  
507 incorrect assignment occurred due to the lack of *I2S* reference sequences for the local native  
508 Opaleye, *G. nigricans*. By maximizing the number of local reference barcodes, regional  
509 databases allow the reads to be correctly assigned to ecologically and geographically relevant  
510 species.

511 In our eDNA samples, the regional database improved species-level assignments,  
512 identifying an additional 17.0% of the total vertebrate sequence reads. Much of this improvement  
513 was due to the inclusion of reference barcodes for Kelp Bass (*Paralabrax clathratus*), one of the  
514 most abundant marine species in Southern California kelp forest ecosystems and an important  
515 sport fishery target (Pondella II et al., 2015). By including a reference barcode for this species,  
516 the regional database assigned 20 previously unidentified ASVs to *P. clathratus*, which  
517 accounted for 16.4% of our total sequence reads. Thus, even the inclusion of reference barcodes  
518 for a few key native taxa can dramatically improve metabarcoding efforts.

519

## 520 **Taxonomic Assignment Limitations of MiFish primers**

521 Of the 8,084 fishes represented in the global database, the MiFish primers were unable to  
522 provide species level taxonomic assignments to 1,322 species (See Table S5 for complete list of  
523 putative *in silico* taxonomic assignments). Thus, although the MiFish primer set has broad utility  
524 for fish metabarcoding, this portion of *I2S* cannot resolve many fishes to species (Miya et al.,  
525 2015). These results highlight the tradeoff between breadth and specificity of any metabarcoding

526 primer set, a result consistent with previous investigations of the MiFish primer set and universal  
527 barcodes in general (Deiner et al., 2017; Miya et al., 2015). Critically, these results provide much  
528 needed insights into taxonomic blind spots of the MiFish primers, informing primer selection for  
529 future fish metabarcoding applications both in the California Current Large Marine Ecosystem  
530 and globally (Figures 4 and 5).

531 Another key limitation to metabarcoding taxonomic assignment is the prevalence of  
532 sequence misannotations in public sequence repositories. Misannotations arise predominantly  
533 from subtle incidental issues, such as mislabeling of sequences, and thus are particularly difficult  
534 to address bioinformatically (Heller et al., 2018; Nobre et al., 2016; Wakeling et al., 2019). To  
535 date, the onus of identifying and preventing misannotations are on the user and research  
536 community and there remain few systematic methods for identifying and removing misannotated  
537 sequences although Kozlov et al., 2016 is a notable exception (we also note there is a process to  
538 flag and report such sequences are available through GenBank). One potential solution to the  
539 issue of misannotated sequences is the development and maintenance of global curated datasets  
540 (e.g., MitoFish, Silva, and UNITE) (Nilsson et al., 2018; Quast et al., 2012; Sato et al., 2018).  
541 However, while these approaches may work well for a handful of key loci and taxonomic targets,  
542 these approaches are not scalable with the rapid development of additional metabarcoding loci  
543 and targets of interest (Curd et al., 2019). Thus, further efforts to systemically prevent and  
544 address mis-annotations in public sequence repositories clearly are warranted.

545

## 546 **Limitations of Barcoding Efforts**

547 The regional database did not include barcodes for all California Current Large Marine  
548 Ecosystem fishes (Table S1) due to a combination of limited resources, difficulties amplifying  
549 vouchered tissue samples, the onset of the COVID-19 pandemic (Omary et al., 2020), and a lack  
550 of some vouchered reference material within the Marine Vertebrates Collection of the Scripps  
551 Institution of Oceanography. In total, our regional database did not include 438 of 1,144  
552 (38.3%) California Current Large Marine Ecosystem fishes. However, the vast majority of these  
553 (n=357) are rare in the state of California (the focus of the collection and study), others (n=25)  
554 are common but not coastal species. Discounting these, our barcoding efforts provide coverage  
555 for 92.7% of the 763 marine fishes common in this ecosystem, making it an important tool for  
556 metabarcoding studies, despite a small number (n=53) of common coastal species missing from  
557 the database (Table S9).

558 The one major shortcoming of our barcoding efforts is that 7.3% (n=32) of the missing  
559 taxa are rockfishes in the genus *Sebastes*. Rockfishes are ecologically important (Hyde & Vetter,  
560 2007), form the basis of many commercial and recreational fisheries (Lea et al., 1999; Williams  
561 et al., 2010), and declines in rockfish stocks led to the establishment of the largest marine  
562 protected areas in southern California, the Cowcod Conservation Areas (Thompson et al., 2017).  
563 Unfortunately, this shortcoming cannot be easily overcome through additional *I2S* barcoding  
564 because rockfish are a recent and diverse radiation comprised of 110 species (Ingram & Kai,  
565 2014) and *I2S* fails to resolve most *Sebastes* to species-level (Hyde & Vetter, 2007; Yamamoto  
566 et al., 2017). Thus, effective metabarcoding of *Sebastes* will require designing novel *Sebastes*-  
567 specific metabarcoding primers that target a more rapidly evolving region of the mitochondrial

568 genome (e.g. *CytB*) (Min et al., 2020; Thompson et al., 2017). Importantly, this *Sebastes*  
569 example highlights the importance of comprehensively evaluating the taxonomic performance of  
570 a particular locus (here MiFish *I2S*) for a given taxonomic group and the difficulty of using  
571 metabarcoding methods for delineating species within an adaptive radiation.

572         Despite these limitations, however, the current regional California Current Large Marine  
573 Ecosystem *I2S*-specific reference database includes all but one non-*Sebastes* nearshore species  
574 monitored by the Channel Islands National Kelp Forest Monitoring Program (n=80, Sprague et  
575 al., 2013), as well as by PISCO, the Partnership for Interdisciplinary Studies of Coastal Oceans  
576 (n=76; the only missing species is White Sea Bass *Atractoscion nobilis*; Caselle, Rassweiler,  
577 Hamilton, & Warner, 2015; Pondella II et al., 2015). Further, there is now a *I2S* reference  
578 sequence for 98 of the 100 most abundant ichthyoplankton species collected by the California  
579 Cooperative Oceanic Fisheries Investigation (CalCOFI) from the California Current Large  
580 Marine Ecosystem between 1951-2019 (only missing Showy Bristlemouth *Cyclothone signata*  
581 and White Barracudina, *Arctozenus risso*) (Moser, 1993). Moreover, in real world application,  
582 this reference barcode database assigned taxonomy to over 90% of vertebrate ASVs detecting a  
583 broad range of ecologically and commercially important nearshore rocky reef species (Pondella  
584 II et al., 2019). As such, our barcoding efforts represents an important genetic resource for  
585 coastal California marine metabarcoding monitoring efforts.

586

### 587 **Off Target Limitations of MiFish primers**

588 High numbers of unidentified ASVs are a common feature of barcoding and metabarcoding  
589 studies (e.g. Leray & Knowlton, 2017). These unidentified ASVs are typically attributed to

590 incomplete reference databases (Curd et al., 2019; Ransome et al., 2017; Schenekar et al., 2020)  
591 and/or novel biodiversity (Barber & Boyce, 2006; Boussarie et al., 2018). However, given that  
592 the regional database includes 92.7% of fishes common in this coastal ecosystem, it was  
593 extremely surprising that half of all ASVs and a quarter of all sequences generated in our eDNA  
594 test datasets could not be assigned.

595         In fact, the vast majority of these sequences and ASVs did not belong to vertebrates, but  
596 instead uncultured marine bacteria, specifically matching to *16S*, rather than the target *12S* locus.  
597 Since mitochondria represent the capture of microbial endosymbionts by ancient eukaryotes  
598 (Roger et al., 2017) and that this capture occurred in the sea, it perhaps is not surprising that  
599 primers designed to target vertebrate *12S* might also capture marine prokaryotes. Similarly, the  
600 homology between vertebrate *12S* and prokaryotic and bacterial *16S* genes is well known (Crews  
601 & Attardi, 1980) suggesting capturing microbial *16S* with vertebrate *12S* primers is not  
602 surprising. However, this particular feature of the MiFish primer set previously has not been  
603 widely reported in the scientific literature (Minamoto et al., 2020), potentially impacting the  
604 interpretation of unidentified ASVs in other fish metabarcoding studies.

605         These findings highlight the importance of accurate universal metabarcoding primer  
606 design, especially in outlining both target and non-target sequences. In the design of the MiFish  
607 Teleost *12S* primers, uncultured marine microbe *16S* sequences were not considered as potential  
608 alternative targets for the primer set, resulting in the selection of a metabarcoding locus with a  
609 high degree of non-target amplification (Miya et al., 2015). This finding is important for the  
610 marine vertebrate eDNA community, which has recently converged on the MiFish *12S* primers  
611 as the vertebrate barcode of choice (Closek et al., 2019; Miya et al., 2020; O'Donnell et al.,

612 2017; Valsecchi et al., 2019; Yamahara et al., 2019). At best, this non-target amplification of  
613 microbial DNA will lead to wasted sequencing effort, as every microbial sequence generated  
614 reduces the number of vertebrate sequences captured. Such a situation would be particularly  
615 problematic for relatively rare targets. At worst, it could result in incorrect interpretation of  
616 unidentified ASVs and lead to incorrect biomonitoring assessments (Cordier et al., 2018). This  
617 problem is of particular concern in biodiversity hotspots such as the Coral Triangle where  
618 reference databases are incomplete, as well as in environments with high abundance of bacteria  
619 relative to vertebrate biomass such as in some pelagic midwater and deep-sea habitats where  
620 recent eDNA sample collection efforts have struggled to detect vertebrate sequences (K. Pitz  
621 personal communication).

622 Previous applications of MiFish *I2S* primer sets did not identify high rates of non-  
623 homologous sequences (Collins et al., 2019; Miya et al., 2015). Interestingly, these studies used  
624 higher annealing temperatures (60-65°C) and fewer PCR cycles than those used in this study  
625 which may potentially explain why we observed high rates of *I6S* amplification. We note that we  
626 used the touchdown PCR method in order to successfully amplify eDNA from sea water  
627 samples. A white paper from *The eDNA Society* in Japan using the original 65°C annealing  
628 methods highlighted that the application of a size-selection step during library preparation (either  
629 via gel extraction or dual size-selection bead clean up) can be used remove off-target sequences  
630 and help ameliorate this issue (Minamoto et al., 2020; Miya & Sado, 2019). We also confirm  
631 here that non-target ASVs are substantially longer in length than vertebrate *I2S* fragments.  
632 Incorporating these practices to reduce microbial cross-amplification will improve the  
633 application of MiFish *I2S* metabarcoding efforts. Ultimately, understanding the full scope of



634 taxa that can be amplified by a given metabarcoding primer is critical for the successful  
635 application and interpretation of results and concerted efforts to validate markers are clearly  
636 warranted.

637

## 638 **Towards improved metabarcoding efforts**

639 The curated California Current Large Marine Ecosystem *I2S*-specific reference database was  
640 designed to improve effectiveness of metabarcoding of California Current Large Marine  
641 Ecosystem fishes. To further improve and expand the taxonomic coverage of the database, we  
642 generated a website that identifies species needing *I2S* reference barcodes and provides the  
643 research community targets for additional barcoding efforts (Zack Gold, 2020). The ability to  
644 update and expand the regional reference database will be especially important as climate change  
645 leads to range expansions of sub-tropical species that may become resident within the California  
646 Current Large Marine Ecosystem (Gentemann et al., 2017; Harvell et al., 2019; Sanford et al.,  
647 2019; Walker et al., 2020). The importance of expanding the database is highlighted by our  
648 detection of Finescale Triggerfish, *Balistes polylepis*, in the eDNA samples, a species that has  
649 only recently become more common off Santa Cruz Island and La Jolla since the 2014-2016  
650 marine heatwave (B. Frable & S. McMillan, personal communication).

651         Additionally, while the MiFish Teleost and Elasmobranch *I2S* loci are important targets  
652 for current marine metabarcoding studies, future efforts and different applications of marine  
653 metabarcoding will likely rely on additional barcoding targets. Here we used the same primer set  
654 to both generate reference barcodes as well as conduct metabarcoding. Although, this choice  
655 limits the applicability and usefulness of our reference barcode generation beyond

656 metabarcoding efforts, it allowed us to more easily and rapidly sequence and generate barcodes  
657 for our intended purpose. Furthermore, recent efforts have found success multiplexing *COI* and  
658 *16S* loci simultaneously provides more species-level identifications than either marker alone,  
659 demonstrating complimentary genetic loci can improve metabarcoding assignments (Duke &  
660 Burton, 2020). Thus future efforts to develop rapid and affordable multilocus barcoding or whole  
661 mitogenomic tools will provide greater resources for marine metabarcoding and population  
662 genomic efforts (Coissac et al., 2016). As these new barcode loci are developed (e.g., *Sebastes*-  
663 specific barcodes), the California Current Large Marine Ecosystem specific reference database  
664 can be expanded to include these loci. Additionally, resources like the SIO Marine Vertebrate  
665 Collection will continue to provide important voucher specimens for advancing marine  
666 molecular ecology resources as they accession new material.

667

## 668 **ACKNOWLEDGEMENTS**

669 We acknowledge NSF GRFP and GRIP [DEG No. 2015204395], UCLA La Kretz Center for  
670 Conservation Genomics, Howard Hughes Medical Institute, and UCLA Department of Ecology  
671 and Evolutionary Biology for funding this research. In addition, we would like to thank Regina  
672 Wetzer, Dean Pentcheff, Adam Wall, Janessa Wall, and Giacomo Bernardi for providing  
673 additional voucher specimens. We would like to thank Julianne K. Passarelli and Milton S. Love  
674 for sharing their comprehensive list of California fish species from their new edition of *Miller*  
675 *and Lea's Guide to the Coastal Marine Fishes of California* (Love & Passarelli, 2020). We  
676 would also like to thank Miya Masaki for sharing protocols and insights on the performance of  
677 the MiFish primer. Lastly, we would like to thank McKenzie Koch, Beverly Shih, Nikita Sridhar,

678 Lauren Man, Eric Caldera, Candice Cross, and Erick Zerecero for help conducting all laboratory  
 679 work. We would like to thank William Watson, Phillip A. Hastings, Peggy Fong, Richard F.  
 680 Ambrose, and Thomas B. Smith for providing thoughtful feedback and comments during the  
 681 preparation of the manuscript.  
 682

## 683 REFERENCES

- 684 Alberdi, A., Aizpurua, O., Gilbert, M. T. P., & Bohmann, K. (2018). Scrutinizing key steps for  
 685 reliable metabarcoding of environmental samples. *Methods in Ecology and Evolution*, *9*(1),  
 686 134–147. <https://doi.org/doi:10.1111/2041-210X.12849>
- 687 Allen, L. G., & Horn, M. H. (2006). *The ecology of marine fishes: California and adjacent*  
 688 *waters*. Univ of California Press.
- 689 Ardura, A., Planes, S., & Garcia-Vazquez, E. (2013). Applications of DNA barcoding to fish  
 690 landings: authentication and diversity assessment. *ZooKeys*, *365*, 49–65.  
 691 <https://doi.org/10.3897/zookeys.365.6409>
- 692 Armstrong, C. G. D. (2017). *Historical ecology of cultural landscapes in the Pacific Northwest*.  
 693 Environment: Department of Archaeology.
- 694 Baetscher, D. S., Anderson, E. C., Gilbert-Horvath, E. A., Malone, D. P., Saarman, E. T., Carr,  
 695 M. H., & Garza, J. C. (2019). Dispersal of a nearshore marine fish connects marine reserves  
 696 and adjacent fished areas along an open coast. *Molecular Ecology*, *28*(7), 1611–1623.  
 697 <https://doi.org/10.1111/mec.15044>
- 698 Baker, C. (2016). *bakerccm/entrez\_qiime: entrez\_qiime v2.0*.  
 699 <https://doi.org/10.5281/ZENODO.159607>
- 700 Balakirev, E. S., Saveliev, P. A., & Ayala, F. J. (2017). Complete mitochondrial genomes of the  
 701 Cherskii's sculpin *Cottus cherskii* and Siberian taimen *Hucho taimen* reveal GenBank entry  
 702 errors: Incorrect species identification and recombinant mitochondrial genome.  
 703 *Evolutionary Bioinformatics*, *13*, 1176934317726783.
- 704 Barber, P., & Boyce, S. L. (2006). Estimating diversity of Indo-Pacific coral reef stomatopods  
 705 through DNA barcoding of stomatopod larvae. *Proceedings of the Royal Society B:*  
 706 *Biological Sciences*, *273*(1597), 2053–2061.
- 707 Benson, D. A., Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., Ostell, J., Pruitt, K. D., & Sayers,  
 708 E. W. (2018). GenBank. *Nucleic Acids Research*, *46*(D1), D41–D47.
- 709 Bergsten, J., Bilton, D. T., Fujisawa, T., Elliott, M., Monaghan, M. T., Balke, M., Hendrich, L.,  
 710 Geijer, J., Herrmann, J., & Foster, G. N. (2012). The effect of geographical scale of  
 711 sampling on DNA barcoding. *Systematic Biology*, *61*(5), 851–869.
- 712 Bista, I., Carvalho, G. R., Walsh, K., Seymour, M., Hajibabaei, M., Lallias, D., Christmas, M., &  
 713 Creer, S. (2017). Annual time-series analysis of aqueous eDNA reveals ecologically  
 714 relevant dynamics of lake ecosystem biodiversity. *Nature Communications*, *8*, 14087.
- 715 Block, B. A., Jonsen, I. D., Jorgensen, S. J., Winship, A. J., Shaffer, S. A., Bograd, S. J., Hazen,  
 716 E. L., Foley, D. G., Breed, G. A., Harrison, A.-L., Ganong, J. E., Swithenbank, A.,

- 717 Castleton, M., Dewar, H., Mate, B. R., Shillinger, G. L., Schaefer, K. M., Benson, S. R.,  
718 Weise, M. J., ... Costa, D. P. (2011). Tracking apex marine predator movements in a  
719 dynamic ocean. *Nature*, 475(7354), 86–90. <https://doi.org/10.1038/nature10082>
- 720 Bohmann, K., Evans, A., Gilbert, M. T. P., Carvalho, G. R., Creer, S., Knapp, M., Yu, D. W., &  
721 de Bruyn, M. (2014). Environmental DNA for wildlife biology and biodiversity monitoring.  
722 *Trends in Ecology & Evolution*, 29(6), 358–367. <https://doi.org/10.1016/j.tree.2014.04.003>
- 723 Bokulich, N. A., Kaehler, B. D., Rideout, J. R., Dillon, M., Bolyen, E., Knight, R., Huttley, G.  
724 A., & Caporaso, J. G. (2018). Optimizing taxonomic classification of marker-gene amplicon  
725 sequences with QIIME 2's q2-feature-classifier plugin. *Microbiome*, 6(1), 90.
- 726 Bork, P., Bowler, C., De Vargas, C., Gorsky, G., Karsenti, E., & Wincker, P. (2015). *Tara*  
727 *Oceans studies plankton at planetary scale*. American Association for the Advancement of  
728 Science.
- 729 Boussarie, G., Bakker, J., Wangensteen, O. S., Mariani, S., Bonnin, L., Juhel, J.-B., Kiszka, J. J.,  
730 Kulbicki, M., Manel, S., & Robbins, W. D. (2018). Environmental DNA illuminates the  
731 dark diversity of sharks. *Science Advances*, 4(5), eaap9661.
- 732 Boyer, F., Mercier, C., Bonin, A., Le Bras, Y., Taberlet, P., & Coissac, E. (2016). obitools: A  
733 unix-inspired software package for DNA metabarcoding. *Molecular Ecology Resources*,  
734 16(1), 176–182.
- 735 Braje, T. J., Rick, T. C., Szpak, P., Newsome, S. D., McCain, J. M., Smith, E. A. E., Glassow,  
736 M., & Hamilton, S. L. (2017). Historical ecology and the conservation of large,  
737 hermaphroditic fishes in Pacific Coast kelp forest ecosystems. *Science Advances*, 3(2),  
738 e1601759.
- 739 Brooks, J. F., Carothers, C., Colombi, B. J., Diver, S., Kasten, E., Koester, D., Lien, M. E.,  
740 Menzies, C. R., Reedy-Maschner, K., & Sharakhmatova, V. N. (2012). *Keystone nations:*  
741 *indigenous peoples and salmon across the north Pacific*. School for Advanced Research  
742 Press.
- 743 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T.  
744 L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics*, 10(1), 421.
- 745 Caselle, J. E., Rassweiler, A., Hamilton, S. L., & Warner, R. R. (2015). Recovery trajectories of  
746 kelp forest animals are rapid yet spatially variable across a network of temperate marine  
747 protected areas. *Scientific Reports*, 5, 14102. <https://doi.org/10.1038/srep14102>
- 748 Chamberlain, S. A., & Szöcs, E. (2013). taxize: taxonomic search and retrieval in R.  
749 *F1000Research*, 2.
- 750 Chan, F., Barth, J. A., Lubchenco, J., Kirincich, A., Weeks, H., Peterson, W. T., & Menge, B. A.  
751 (2008). Emergence of anoxia in the California current large marine ecosystem. *Science*,  
752 319(5865), 920.
- 753 Checkley Jr, D. M., & Barth, J. A. (2009). Patterns and processes in the California Current  
754 System. *Progress in Oceanography*, 83(1–4), 49–64.
- 755 Closek, C. J., Santora, J. A., Starks, H. A., Schroeder, I. D., Andruszkiewicz, E. A., Sakuma, K.  
756 M., Bograd, S. J., Hazen, E. L., Field, J. C., & Boehm, A. B. (2019). Marine vertebrate  
757 biodiversity and distribution within the central California Current using environmental  
758 DNA (eDNA) metabarcoding and ecosystem surveys. *Frontiers in Marine Science*, 6, 732.
- 759 Coissac, E., Hollingsworth, P. M., Lavergne, S., & Taberlet, P. (2016). From barcodes to  
760 genomes: extending the concept of DNA barcoding. *Molecular Ecology*, 25(7), 1423–1428.

- 761 Coleman, K. (2008). Research review of collaborative ecosystem-based management in the  
762 California current large marine ecosystem. *Coastal Management*, 36(5), 484–494.
- 763 Collins, R. A., Bakker, J., Wangenstein, O. S., Soto, A. Z., Corrigan, L., Sims, D. W., Genner,  
764 M. J., & Mariani, S. (2019). Non-specific amplification compromises environmental DNA  
765 metabarcoding with COI. *Methods in Ecology and Evolution*, 10(11), 1985–2001.  
766 <https://doi.org/10.1111/2041-210X.13276>
- 767 Cordier, T., Forster, D., Dufresne, Y., Martins, C. I. M., Stoeck, T., & Pawlowski, J. (2018).  
768 Supervised machine learning outperforms taxonomy-based environmental DNA  
769 metabarcoding applied to biomonitoring. *Molecular Ecology Resources*, 18(6), 1381–1391.
- 770 Crews, S., & Attardi, G. (1980). The sequences of the small ribosomal RNA gene and the  
771 phenylalanine tRNA gene are joined end to end in human mitochondrial DNA. *Cell*, 19(3),  
772 775–784.
- 773 Crozier, L. G., McClure, M. M., Beechie, T., Bograd, S. J., Boughton, D. A., Carr, M., Cooney,  
774 T. D., Dunham, J. B., Greene, C. M., & Haltuch, M. A. (2019). Climate vulnerability  
775 assessment for Pacific salmon and steelhead in the California Current Large Marine  
776 Ecosystem. *PloS One*, 14(7), e0217711.
- 777 Curd, E. E., Gold, Z., Kandlikar, G. S., Gomer, J., Ogden, M., O'Connell, T., Pipes, L.,  
778 Schweizer, T. M., Rabichow, L., Lin, M., Shi, B., Barber, P. H., Kraft, N., Wayne, R., &  
779 Meyer, R. S. (2019). Anacapa: an environmental DNA toolkit for processing multilocus  
780 metabarcode datasets. *Methods in Ecology and Evolution*, 10, 1469–1475.  
781 <https://doi.org/https://doi.org/10.1111/2041-210X.13214>
- 782 Deiner, K., Bik, H. M., Mächler, E., Seymour, M., Lacoursière-Roussel, A., Altermatt, F., Creer,  
783 S., Bista, I., Lodge, D. M., de Vere, N., Pfrender, M. E., & Bernatchez, L. (2017).  
784 Environmental DNA metabarcoding: Transforming how we survey animal and plant  
785 communities. *Molecular Ecology*, 26(21), 5872–5895. <https://doi.org/10.1111/mec.14350>
- 786 Djurhuus, A., Closek, C. J., Kelly, R. P., Pitz, K. J., Michisaki, R. P., Starks, H. A., Walz, K. R.,  
787 Andruszkiewicz, E. A., Olesin, E., & Hubbard, K. (2020). Environmental DNA reveals  
788 seasonal shifts and potential interactions in a marine community. *Nature Communications*,  
789 11(1), 1–9.
- 790 Duke, E. M., & Burton, R. S. (2020). Efficacy of metabarcoding for identification of fish eggs  
791 evaluated with mock communities. *Ecology and Evolution*, 10(7), 3463–3476.  
792 <https://doi.org/10.1002/ece3.6144>
- 793 Edgar, R. C. (2018a). Accuracy of taxonomy prediction for 16S rRNA and fungal ITS  
794 sequences. *PeerJ*, 6, e4652.
- 795 Edgar, R. C. (2018b). Taxonomy annotation and guide tree errors in 16S rRNA databases. *PeerJ*,  
796 6, e5030.
- 797 Edgar, R. C. (2018c). Updating the 97% identity threshold for 16S ribosomal RNA OTUs.  
798 *Bioinformatics*, 34(14), 2371–2375.
- 799 Ekstrom, J. A. (2009). California Current Large Marine Ecosystem: Publicly available dataset of  
800 state and federal laws and regulations. *Marine Policy*, 33(3), 528–531.
- 801 Ficetola, G. F., Coissac, E., Zundel, S., Riaz, T., Shehzad, W., Bessière, J., Taberlet, P., &  
802 Pompanon, F. (2010). An in silico approach for the evaluation of DNA barcodes. *BMC*  
803 *Genomics*, 11(1), 434.
- 804 Froese, R., & Pauly, D. (2010). *FishBase*. Fisheries Centre, University of British Columbia.

- 805 Gao, X., Lin, H., Revanna, K., & Dong, Q. (2017). A Bayesian taxonomic classification method  
806 for 16S rRNA gene sequences with improved species-level accuracy. *BMC Bioinformatics*,  
807 *18*(1), 247.
- 808 Gentemann, C. L., Fewings, M. R., & García-Reyes, M. (2017). Satellite sea surface  
809 temperatures along the West Coast of the United States during the 2014–2016 northeast  
810 Pacific marine heat wave. *Geophysical Research Letters*, *44*(1), 312–319.
- 811 Gold, Zachary, Sprague, J., Kushner, D. J., Zerecero Marin, E., & Barber, P. H. (2021). eDNA  
812 metabarcoding as a biomonitoring tool for marine protected areas. *PLoS Biology*, *In press*.
- 813 Gold, Zack. (2020). *zjgold/FishCARD: California Fish Reference Database Version 0.0*.  
814 Zenodo. <https://doi.org/10.5281/zenodo.4315278>
- 815 Good, T. P., Samhouri, J. F., Feist, B. E., Wilcox, C., & Jahncke, J. (2020). Plastics in the  
816 Pacific: Assessing risk from ocean debris for marine birds in the California Current Large  
817 Marine Ecosystem. *Biological Conservation*, *250*, 108743.
- 818 Goodwin, K. D., Thompson, L. R., Duarte, B., Kahlke, T., Thompson, A. R., Marques, J. C., &  
819 Caçador, I. (2017). DNA sequencing as a tool to monitor marine ecological status. *Frontiers*  
820 *in Marine Science*, *4*, 107.
- 821 Guo, J. (2017). *Metabarcoding Analyses of Gut Microbiome Compositions in Red Abalone*  
822 *(Haliotis Rufescens, Swainson, 1822) Fed Different Macroalgal Diets*.
- 823 Halpern, B. S., Kappel, C. V., Selkoe, K. A., Micheli, F., Ebert, C. M., Kontgis, C., Crain, C. M.,  
824 Martone, R. G., Shearer, C., & Teck, S. J. (2009). Mapping cumulative human impacts to  
825 California Current marine ecosystems. *Conservation Letters*, *2*(3), 138–148.
- 826 Harvell, C. D., Montecino-Latorre, D., Caldwell, J. M., Burt, J. M., Bosley, K., Keller, A.,  
827 Heron, S. F., Salomon, A. K., Lee, L., Pontier, O., Pattengill-Semmens, C., & Gaydos, J. K.  
828 (2019). Disease epidemic and a marine heat wave are associated with the continental-scale  
829 collapse of a pivotal predator (*Pycnopodia helianthoides*). *Science Advances*, *5*(1),  
830 eaau7042. <https://doi.org/10.1126/sciadv.aau7042>
- 831 Hassanin, A., Bonillo, C., Nguyen, B. X., & Cruaud, C. (2010). Comparisons between  
832 mitochondrial genomes of domestic goat (*Capra hircus*) reveal the presence of numts and  
833 multiple sequencing errors. *Mitochondrial DNA*, *21*(3–4), 68–76.
- 834 Hastings, P. A., & Burton, R. S. (2008). *Establishing a DNA sequence database for the marine*  
835 *fish fauna of California*.
- 836 Heller, P., Casaletto, J., Ruiz, G., & Geller, J. (2018). A database of metazoan cytochrome c  
837 oxidase subunit I gene sequences derived from GenBank with CO-ARBitrator. *Scientific*  
838 *Data*, *5*.
- 839 Hofmann, G. E., Evans, T. G., Kelly, M. W., Padilla-Gamiño, J. L., Blanchette, C. A.,  
840 Washburn, L., Chan, F., McManus, M. A., Menge, B. A., & Gaylord, B. (2014). Exploring  
841 local adaptation and the ocean acidification seascape—studies in the California Current  
842 Large Marine Ecosystem. *Biogeosciences*, *11*(4).
- 843 Hyde, J. R., & Vetter, R. D. (2007). The origin, evolution, and diversification of rockfishes of the  
844 genus *Sebastes* (Cuvier). *Molecular Phylogenetics and Evolution*, *44*(2), 790–811.  
845 <https://doi.org/10.1016/j.ympev.2006.12.026>
- 846 Ingram, T., & Kai, Y. (2014). The geography of morphological convergence in the radiations of  
847 Pacific *Sebastes* rockfishes. *The American Naturalist*, *184*(5), E115–31.  
848 <https://doi.org/10.1086/678053>

- 849 Iwasaki, W., Fukunaga, T., Isagozawa, R., Yamada, K., Maeda, Y., Satoh, T. P., Sado, T.,  
850 Mabuchi, K., Takeshima, H., & Miya, M. (2013). MitoFish and MitoAnnotator: a  
851 mitochondrial genome database of fish with an accurate and automatic annotation pipeline.  
852 *Molecular Biology and Evolution*, *30*(11), 2531–2540.
- 853 Jo, T., Murakami, H., Masuda, R., Sakata, M. K., Yamamoto, S., & Minamoto, T. (2017). Rapid  
854 degradation of longer DNA fragments enables the improved estimation of distribution and  
855 biomass using environmental DNA. *Molecular Ecology Resources*, *17*(6), e25–e33.
- 856 Kelly, R. P., Port, J. A., Yamahara, K. M., & Crowder, L. B. (2014). Using environmental DNA  
857 to census marine fishes in a large mesocosm. *PloS One*, *9*(1), e86175.  
858 <https://doi.org/10.1371/journal.pone.0086175>
- 859 Kelly, R. P., Port, J. a., Yamahara, K. M., Martone, R. G., Lowell, N., Thomsen, P. F., Mach, M.  
860 E., Bennett, M., Prahler, E., Caldwell, M. R., & Crowder, L. B. (2014). Harnessing DNA to  
861 improve environmental management. *Science*, *344*(6191).  
862 <https://doi.org/10.1126/science.1251156>
- 863 Klymus, K. E., Marshall, N. T., & Stepien, C. A. (2017). Environmental DNA (eDNA)  
864 metabarcoding assays to detect invasive invertebrate species in the Great Lakes. *PLOS*  
865 *ONE*, *12*(5), e0177643. <https://doi.org/10.1371/journal.pone.0177643>
- 866 Koslow, J. A., & Davison, P. C. (2016). Productivity and biomass of fishes in the California  
867 Current Large Marine Ecosystem: Comparison of fishery-dependent and-independent time  
868 series. *Environmental Development*, *17*, 23–32.
- 869 Kozlov, A. M., Zhang, J., Yilmaz, P., Glöckner, F. O., & Stamatakis, A. (2016). Phylogeny-  
870 aware identification and correction of taxonomically mislabeled sequences. *Nucleic Acids*  
871 *Research*, *44*(11), 5022–5033.
- 872 Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature*  
873 *Methods*, *9*(4), 357.
- 874 Lea, R. N., McAllister, R. D., & VenTresca, D. A. (1999). *Biological aspects of nearshore*  
875 *rockfishes of the genus Sebastes from central California: with notes on ecologically related*  
876 *sport fishes* (Vol. 177). State of California, The Resources Agency, Department of Fish and  
877 Game.
- 878 Lepofsky, D., Armstrong, C. G., Greening, S., Jackley, J., Carpenter, J., Guernsey, B., Mathews,  
879 D., & Turner, N. J. (2017). Historical ecology of cultural keystone places of the Northwest  
880 Coast. *American Anthropologist*, *119*(3), 448–463.
- 881 Leray, M., Boehm, J. T., Mills, S. C., & Meyer, C. P. (2012). Moorea BIOC CODE barcode library  
882 as a tool for understanding predator–prey interactions: insights into the diet of common  
883 predatory coral reef fishes. *Coral Reefs*, *31*(2), 383–388. [https://doi.org/10.1007/s00338-](https://doi.org/10.1007/s00338-011-0845-0)  
884 [011-0845-0](https://doi.org/10.1007/s00338-011-0845-0)
- 885 Leray, M., & Knowlton, N. (2015). DNA barcoding and metabarcoding of standardized samples  
886 reveal patterns of marine benthic diversity. *Proceedings of the National Academy of*  
887 *Sciences*, *112*(7), 201424997. <https://doi.org/10.1073/pnas.1424997112>
- 888 Leray, M., & Knowlton, N. (2017). Random sampling causes the low reproducibility of rare  
889 eukaryotic OTUs in Illumina COI metabarcoding. *PeerJ*, *5*, e3006.
- 890 Leray, M., Knowlton, N., Ho, S.-L., Nguyen, B. N., & Machida, R. J. (2019). GenBank is a  
891 reliable resource for 21st century biodiversity research. *Proceedings of the National*  
892 *Academy of Sciences*, *116*(45), 22651–22656.

- 893 Lim, S. J., & Thompson, L. R. (2021). Mitohelper: A mitochondrial reference sequence analysis  
894 tool for fish eDNA studies. *Environmental DNA*.
- 895 Lodge, D. M., Turner, C. R., Jerde, C. L., Barnes, M. A., Chadderton, L., Egan, S. P., Feder, J.  
896 L., Mahon, A. R., & Pfrender, M. E. (2012). Conservation in a cup of water: estimating  
897 biodiversity and population abundance from environmental DNA. *Molecular Ecology*,  
898 *21*(11), 2555–2558. <https://doi.org/10.1111/j.1365-294X.2012.05600.x>
- 899 Love, M. S., & Passarelli, J. K. (2020). *Miller and Lea's Guide to the Coastal Marine Fishes of*  
900 *California* (2nd.). University of California Agriculture and Natural Resources.
- 901 Macheriotou, L., Guilini, K., Bezerra, T. N., Tytgat, B., Nguyen, D. T., Phuong Nguyen, T. X.,  
902 Noppe, F., Armenteros, M., Boufahja, F., & Rigaux, A. (2019). Metabarcoding free-living  
903 marine nematodes using curated 18S and CO1 reference sequence databases for  
904 species-level taxonomic assignments. *Ecology and Evolution*, *9*(3), 1211–1226.
- 905 Machida, R. J., Leray, M., Ho, S.-L., & Knowlton, N. (2017). Metazoan mitochondrial gene  
906 sequence reference datasets for taxonomic assignment of environmental samples. *Scientific*  
907 *Data*, *4*, 170027.
- 908 Min, M. A., Barber, P. H., & Gold, Z. (2020). MiSebastes: An eDNA metabarcoding primer set  
909 for rockfishes (genus *Sebastes*). *BioRxiv*, 2020.10.29.360859.  
910 <https://doi.org/10.1101/2020.10.29.360859>
- 911 Minamoto, T., Miya, M., Sado, T., Seino, S., Doi, H., Kondoh, M., Nakamura, K., Takahara, T.,  
912 Yamamoto, S., & Yamanaka, H. (2020). An illustrated manual for environmental DNA  
913 research: Water sampling guidelines and experimental protocols. *Environmental DNA*.
- 914 Miya, M., Gotoh, R. O., & Sado, T. (2020). MiFish metabarcoding: a high-throughput approach  
915 for simultaneous detection of multiple fish species from environmental DNA and other  
916 samples. *Fisheries Science*, 1–32.
- 917 Miya, M., & Sado, T. (2019). *Environmental DNA Sampling and Experiment Manual (Version*  
918 *2.1)*.
- 919 Miya, M., Sato, Y., Fukunaga, T., Sado, T., Poulsen, J. Y., Sato, K., Minamoto, T., Yamamoto,  
920 S., Yamanaka, H., Araki, H., Kondoh, M., & Iwasaki, W. (2015). MiFish, a set of universal  
921 PCR primers for metabarcoding environmental DNA from fishes: detection of more than  
922 230 subtropical marine species. *Royal Society Open Science*, *2*(7), 150088.  
923 <https://doi.org/10.1098/rsos.150088>
- 924 Moser, H. G. (1993). *Distributional atlas of fish larvae and eggs in the California Current*  
925 *region: taxa with 1000 or more total larvae, 1951 through 1984* (Issue 31). Marine Life  
926 Research Program, Scripps Institution of Oceanography.
- 927 Nilsson, R. H., Larsson, K.-H., Taylor, A. F. S., Bengtsson-Palme, J., Jeppesen, T. S., Schigel,  
928 D., Kennedy, P., Picard, K., Glöckner, F. O., & Tedersoo, L. (2018). The UNITE database  
929 for molecular identification of fungi: handling dark taxa and parallel taxonomic  
930 classifications. *Nucleic Acids Research*, *47*(D1), D259–D264.
- 931 Nishimura, D. (2000). Sequencher 3.1. 1. *Biotech Software & Internet Report*, *1*(1–2), 24–30.
- 932 NMFS. (2017). *Fisheries economics of the United States, 2015*. NOAA Technical Memorandum.
- 933 Nobre, T., Campos, M. D., Lucic-Mercy, E., & Arnholdt-Schmitt, B. (2016). Misannotation  
934 awareness: a tale of two gene-groups. *Frontiers in Plant Science*, *7*, 868.
- 935 Norgaard, K. M. (2019). *Salmon and acorns feed our people: colonialism, nature, and social*  
936 *action*. Rutgers University Press.



- 937 O'Donnell, J. L., Kelly, R. P., Shelton, A. O., Samhouri, J. F., Lowell, N. C., & Williams, G. D.  
 938 (2017). Spatial distribution of environmental DNA in a nearshore marine habitat. *PeerJ*, *5*,  
 939 e3044. <https://doi.org/10.7717/peerj.3044>
- 940 Omary, M. B., Eswaraka, J., Kimball, S. D., Moghe, P. V., Panettieri, R. A., & Scotto, K. W.  
 941 (2020). The COVID-19 pandemic and research shutdown: staying safe and productive. *The*  
 942 *Journal of Clinical Investigation*, *130*(6).
- 943 Pitz, K. J., Guo, J., Johnson, S. B., Campbell, T. L., Zhang, H., Vrijenhoek, R. C., Chavez, F. P.,  
 944 & Geller, J. (2020). Zooplankton biogeographic boundaries in the California Current  
 945 System as determined from metabarcoding. *Plos One*, *15*(6), e0235159.
- 946 Poloczanska, E. S., Brown, C. J., Sydeman, W. J., Kiessling, W., Schoeman, D. S., Moore, P. J.,  
 947 Brander, K., Bruno, J. F., Buckley, L. B., Burrows, M. T., Duarte, C. M., Halpern, B. S.,  
 948 Holding, J., Kappel, C. V., O'Connor, M. I., Pandolfi, J. M., Parmesan, C., Schwing, F.,  
 949 Thompson, S. A., & Richardson, A. J. (2013). Global imprint of climate change on marine  
 950 life. *Nature Climate Change*, *3*(10), 919–925. <https://doi.org/10.1038/nclimate1958>
- 951 Pondella II, D. J., Caselle, J. E., Claisse, J. T., Williams, J. P., Davis, K., Williams, C. M., &  
 952 Zahn, L. A. (2015). *Baseline Characterization of the Shallow Rocky Reef and Kelp Forest*  
 953 *Ecosystems of the South Coast Study Region*. [https://casegrant.ucsd.edu/news/summaries-](https://casegrant.ucsd.edu/news/summaries-of-projects-selected-for-funding-through-the-south-coast-mpa-baseline-program)  
 954 [of-projects-selected-for-funding-through-the-south-coast-mpa-baseline-program](https://casegrant.ucsd.edu/news/summaries-of-projects-selected-for-funding-through-the-south-coast-mpa-baseline-program)
- 955 Pondella II, D. J., Piacenza, S. E., Claisse, J. T., Williams, C. M., Williams, J. P., Zellmer, A. J.,  
 956 & Caselle, J. E. (2019). Assessing drivers of rocky reef fish biomass density from the  
 957 Southern California Bight. *Marine Ecology Progress Series*, *628*, 125–140.
- 958 Port, J. A., O'Donnell, J. L., Romero-Maraccini, O. C., Leary, P. R., Litvin, S. Y., Nickols, K. J.,  
 959 Yamahara, K. M., & Kelly, R. P. (2015). Assessing vertebrate biodiversity in a kelp forest  
 960 ecosystem using environmental DNA. *Molecular Ecology*.  
 961 <https://doi.org/https://doi.org/10.1111/mec.13481>
- 962 Pruitt, K. D., Tatusova, T., & Maglott, D. R. (2005). NCBI Reference Sequence (RefSeq): a  
 963 curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic*  
 964 *Acids Research*, *33*(suppl\_1), D501–D504.
- 965 Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., Peplies, J., & Glöckner, F.  
 966 O. (2012). The SILVA ribosomal RNA gene database project: improved data processing  
 967 and web-based tools. *Nucleic Acids Research*, *41*(D1), D590–D596.
- 968 Ransome, E., Geller, J. B., Timmers, M., Leray, M., Mahardini, A., Sembiring, A., Collins, A.  
 969 G., & Meyer, C. P. (2017). The importance of standardization for biodiversity comparisons:  
 970 A case study using autonomous reef monitoring structures (ARMS) and metabarcoding to  
 971 measure cryptic diversity on Mo'orea coral reefs, French Polynesia. *PloS One*, *12*(4).
- 972 Richardson, R. T., Bengtsson-Palme, J., Gardiner, M. M., & Johnson, R. M. (2018). A reference  
 973 cytochrome c oxidase subunit I database curated for hierarchical classification of arthropod  
 974 metabarcoding data. *PeerJ*, *6*, e5126.
- 975 Roger, A. J., Muñoz-Gómez, S. A., & Kamikawa, R. (2017). The origin and diversification of  
 976 mitochondria. *Current Biology*, *27*(21), R1177–R1192.
- 977 Rogers-Bennett, L., & Catton, C. A. (2019). Marine heat wave and multiple stressors tip bull  
 978 kelp forest to sea urchin barrens. *Scientific Reports*, *9*(1), 1–9.  
 979 <https://doi.org/10.1038/s41598-019-51114-y>
- 980 Samhouri, J. F., Andrews, K. S., Fay, G., Harvey, C. J., Hazen, E. L., Hennessey, S. M.,

- 981 Holsman, K., Hunsicker, M. E., Large, S. I., & Marshall, K. N. (2017). Defining ecosystem  
982 thresholds for human activities and environmental pressures in the California Current.  
983 *Ecosphere*, 8(6), e01860.
- 984 Sanders, J. G., Beichman, A. C., Roman, J., Scott, J. J., Emerson, D., McCarthy, J. J., & Girguis,  
985 P. R. (2015). Baleen whales host a unique gut microbiome with similarities to both  
986 carnivores and herbivores. *Nature Communications*, 6, 8285.
- 987 Sanford, E., Sones, J. L., García-Reyes, M., Goddard, J. H. R. R., & Largier, J. L. (2019).  
988 Widespread shifts in the coastal biota of northern California during the 2014–2016 marine  
989 heatwaves. *Scientific Reports*, 9(1), 4216. <https://doi.org/10.1038/s41598-019-40784-3>
- 990 Santora, J. A., Mantua, N. J., Schroeder, I. D., Field, J. C., Hazen, E. L., Bograd, S. J., Sydeman,  
991 W. J., Wells, B. K., Calambokidis, J., Saez, L., Lawson, D., & Forney, K. A. (2020).  
992 Habitat compression and ecosystem shifts as potential links between marine heatwave and  
993 record whale entanglements. *Nature Communications*, 11(1), 1–12.  
994 <https://doi.org/10.1038/s41467-019-14215-w>
- 995 Sato, Y., Miya, M., Fukunaga, T., Sado, T., & Iwasaki, W. (2018). MitoFish and MiFish  
996 pipeline: a mitochondrial genome database of fish with an analysis pipeline for  
997 environmental DNA metabarcoding. *Molecular Biology and Evolution*, 35(6), 1553–1555.
- 998 Schenekar, T., Schletterer, M., Lecaudey, L. A., & Weiss, S. J. (2020). Reference databases,  
999 primer choice, and assay sensitivity for environmental metabarcoding: Lessons learnt from  
1000 a re-evaluation of an eDNA fish assessment in the Volga headwaters. *River Research and*  
1001 *Applications*.
- 1002 Sherman, K. (1991). The large marine ecosystem concept: research and management strategy for  
1003 living marine resources. *Ecological Applications*, 1(4), 349–360.
- 1004 Siegwald, L., Touzet, H., Lemoine, Y., Hot, D., Audebert, C., & Caboche, S. (2017). Assessment  
1005 of Common and Emerging Bioinformatics Pipelines for Targeted Metagenomics. *PloS One*,  
1006 12(1). <https://doi.org/10.1371/journal.pone.0169563>
- 1007 Stat, M., Huggett, M. J., Bernasconi, R., DiBattista, J. D., Berry, T. E., Newman, S. J., Harvey,  
1008 E. S., & Bunce, M. (2017). Ecosystem biomonitoring with eDNA: metabarcoding across the  
1009 tree of life in a tropical marine environment. *Scientific Reports*, 7(1), 12240.  
1010 <https://doi.org/10.1038/s41598-017-12501-5>
- 1011 Stoeckle, M. Y., Das Mishu, M., & Charlop-Powers, Z. (2020). Improved environmental DNA  
1012 reference library detects overlooked marine fishes in New Jersey, United States. *Frontiers*  
1013 *in Marine Science*, 7, 226.
- 1014 Stoesser, G., Baker, W., van den Broek, A., Camon, E., Garcia-Pastor, M., Kanz, C., Kulikova,  
1015 T., Leinonen, R., Lin, Q., & Lombard, V. (2002). The EMBL nucleotide sequence database.  
1016 *Nucleic Acids Research*, 30(1), 21–26.
- 1017 Taberlet, P., Coissac, E., Hajibabaei, M., & Rieseberg, L. H. (2012). Environmental DNA.  
1018 *Molecular Ecology*, 21(8), 1789–1793. <https://doi.org/10.1111/j.1365-294X.2012.05542.x>
- 1019 Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C., & Willerslev, E. (2012). Towards next-  
1020 generation biodiversity assessment using DNA metabarcoding. *Molecular Ecology*, 21(8),  
1021 2045–2050. <https://doi.org/10.1111/j.1365-294X.2012.05470.x>
- 1022 Thompson, A. R., Chen, D. C., Guo, L. W., Hyde, J. R., & Watson, W. (2017). Larval  
1023 abundances of rockfishes that were historically targeted by fishing increased over 16 years  
1024 in association with a large marine protected area. *Royal Society Open Science*, 4(9), 170639.

- 1025 <https://doi.org/10.1098/rsos.170639>
- 1026 Thomsen, P. F., Møller, P. R., Sigsgaard, E. E., Knudsen, S. W., Jørgensen, O. A., & Willerslev,  
1027 E. (2016). Environmental DNA from seawater samples correlate with trawl catches of  
1028 subarctic, deepwater fishes. *PLoS ONE*, *11*(11).  
1029 <https://doi.org/10.1371/journal.pone.0165252>
- 1030 Valsecchi, E., Bylemans, J., Goodman, S. J., Lombardi, R., Carr, I., Castellano, L., Galimberti,  
1031 A., & Galli, P. (2019). Novel Universal Primers for Metabarcoding eDNA Surveys of  
1032 Marine Mammals and Other Marine Vertebrates. *BioRxiv*, 759746.
- 1033 Wadewitz, L. K. (2012). *The nature of borders: Salmon, boundaries, and bandits on the Salish*  
1034 *Sea*. University of Washington Press.
- 1035 Wakeling, M. N., Laver, T. W., Colclough, K., Parish, A., Ellard, S., & Baple, E. L. (2019).  
1036 Misannotation of multiple-nucleotide variants risks misdiagnosis. *Wellcome Open*  
1037 *Research*, *4*.
- 1038 Walker, H. J., Hastings, P. A., Hyde, J. R., Lea, R. N., Snodgrass, O. E., & Bellquist, L. F.  
1039 (2020). Unusual occurrences of fishes in the Southern California Current System during the  
1040 warm water period of 2014–2018. *Estuarine, Coastal and Shelf Science*, *236*, 106634.  
1041 <https://doi.org/10.1016/j.ecss.2020.106634>
- 1042 Walsh, P. S., Metzger, D. A., & Higuchi, R. (1991). Chelex 100 as a medium for simple  
1043 extraction of DNA for PCR-based typing from forensic material. *Biotechniques*, *10*(4), 506–  
1044 513.
- 1045 Ward, R. D., Hanner, R., & Hebert, P. D. N. (2009). The campaign to DNA barcode all fishes,  
1046 FISH-BOL. *Journal of Fish Biology*, *74*(2), 329–356. [https://doi.org/10.1111/j.1095-](https://doi.org/10.1111/j.1095-8649.2008.02080.x)  
1047 [8649.2008.02080.x](https://doi.org/10.1111/j.1095-8649.2008.02080.x)
- 1048 Wells, R. J. D., Mohan, J. A., Dewar, H., Rooker, J. R., Tanaka, Y., Snodgrass, O. E., Kohin, S.,  
1049 Miller, N. R., & Ohshimo, S. (2020). Natal origin of Pacific bluefin tuna from the California  
1050 Current Large Marine Ecosystem. *Biology Letters*, *16*(2), 20190878.
- 1051 Williams, G. D., Levin, P. S., & Palsson, W. A. (2010). Rockfish in Puget Sound: An ecological  
1052 history of exploitation. *Marine Policy*, *34*(5), 1010–1020.  
1053 <https://doi.org/https://doi.org/10.1016/j.marpol.2010.02.008>
- 1054 Wood, D. E., & Salzberg, S. L. (2014). Kraken: ultrafast metagenomic sequence classification  
1055 using exact alignments. *Genome Biology*, *15*(3), R46. [https://doi.org/10.1186/gb-2014-15-3-](https://doi.org/10.1186/gb-2014-15-3-r46)  
1056 [r46](https://doi.org/10.1186/gb-2014-15-3-r46)
- 1057 Yamahara, K. M., Preston, C. M., Birch, J. M., Walz, K. R., Marin III, R., Jensen, S., Pargett, D.,  
1058 Roman, B., Zhang, Y., & Ryan, J. (2019). In-situ Autonomous Acquisition and Preservation  
1059 of Marine Environmental DNA Using an Autonomous Underwater Vehicle. *Frontiers in*  
1060 *Marine Science*, *6*, 373. <https://doi.org/https://doi.org/10.3389/fmars.2019.00373>
- 1061 Yamamoto, S., Masuda, R., Sato, Y., Sado, T., Araki, H., Kondoh, M., Minamoto, T., & Miya,  
1062 M. (2017). Environmental DNA metabarcoding reveals local fish communities in a species-  
1063 rich coastal sea. *Scientific Reports*, *7*, 40368.
- 1064

## 1065 DATA ACCESSIBILITY

1066 Reference databases and metabarcoding data are publicly available and stored on a Dryad  
1067 repository (<https://doi.org/10.5068/D1H963>). All reference barcode sequences have been  
1068 uploaded to GenBank (BioProject PRJNA731549). Additional supporting information is  
1069 available at <https://github.com/zjgold/FishCARD>.

## 1070 **AUTHOR CONTRIBUTIONS**

- 1071 • Conceptualization ZG, ESC, DK, BF, RSB, KDG, ART, PHB, HJW
- 1072 • Performed Research ZG, ESC, DK, BF, ART
- 1073 • Funding Acquisition ZG, PHB, ART, KDG, DK, RSB
- 1074 • Data Curation ZG, ESC, DK, BF, HJW
- 1075 • Formal Analysis ZG, EEC
- 1076 • Writing – Original Draft Preparation ZG, EEC, ESC, DK, BF, RSB, KDG, ART, PHB,  
1077 HJW

1078

## 1079 **CONFLICTS OF INTEREST**

1080 The authors have no conflicts of interest to report.

1081

1082

1083

1084 **TABLES**

1085

1086 **Table 1. Summary of Cross Validation Results.** Comparison of performance metrics for  
 1087 taxonomic assignments using the global database as a reference to annotate sequences in the  
 1088 global database (global-global)[ (test database-training database)], the regional database (global-  
 1089 regional), and using the regional database as a reference to annotate sequences in itself (regional-  
 1090 regional). Reporting metrics calculated using a taxonomic cutoff score of 60.

<b>Metric</b>	<b>Global-Global</b>	<b>Global-Regional</b>	<b>Regional-Regional</b>
<b>Under-classification Rate</b>	8.6%	11.8%	7.8%
<b>Misclassification Rate</b>	1.7%	1.8%	1.3%
<b>Over-classification Rate</b>	0.0%	0.0%	0.0%
<b>Accuracy</b>	89.7%	86.5%	90.9%
<b>True Positive Rate</b>	89.7%	86.5%	90.9%
<b>Sensitivity</b>	91.3%	88.0%	92.1%
<b>Specificity</b>	98.3%	98.2%	98.7%

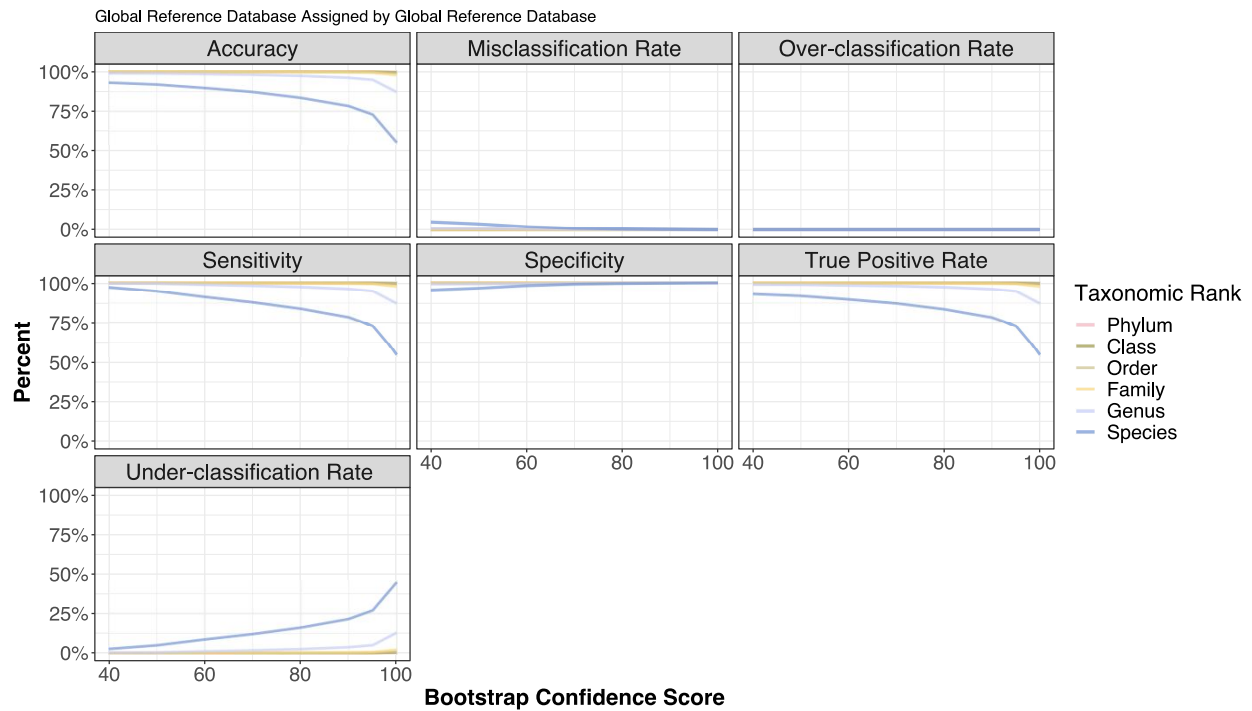
1091

1092 Table 2. Summary of Seawater eDNA Metabarcoding Taxonomic Assignments for Tested  
 1093 Reference Databases.

Metric		Reference Database		
		CRUX-GenBank	Global	Regional
<b>Database</b>	Reference Barcode Origin	GenBank	GenBank + Generated	GenBank + Generated
	Species Included	<b>All</b>	<b>All</b>	<b>California Fishes</b>
<b>Reads</b>	Total Reads	330,877		
	Assigned to NA	81,014	81,002	81,006
	Assigned to Class Level	54,090	-	-
	Assigned to Order Level	727	-	-
	Assigned to Family Level	1,286	1,409	131
	Assigned to Genus Level	952	1,068	1,063
	Assigned to Species Level	192,808	247,398	248,677
<b>ASVs</b>	Total ASVs	341		
	Assigned to NA	172	169	170
	Assigned to Class Level	12	-	-
	Assigned to Order Level	3	-	-
	Assigned to Family Level	5	13	11
	Assigned to Genus Level	4	6	4
	Assigned to Species Level	145	153	156
<b>Taxonomy</b>	Unique Families Identified	31	28	27
	Unique Genera Identified	39	38	39
	Unique Species Identified	38	38	37
	CA Native Species	25	36	37
	Avg. ASVs Per Species	3.8	4.1	4.2

1094

1095 **FIGURES**

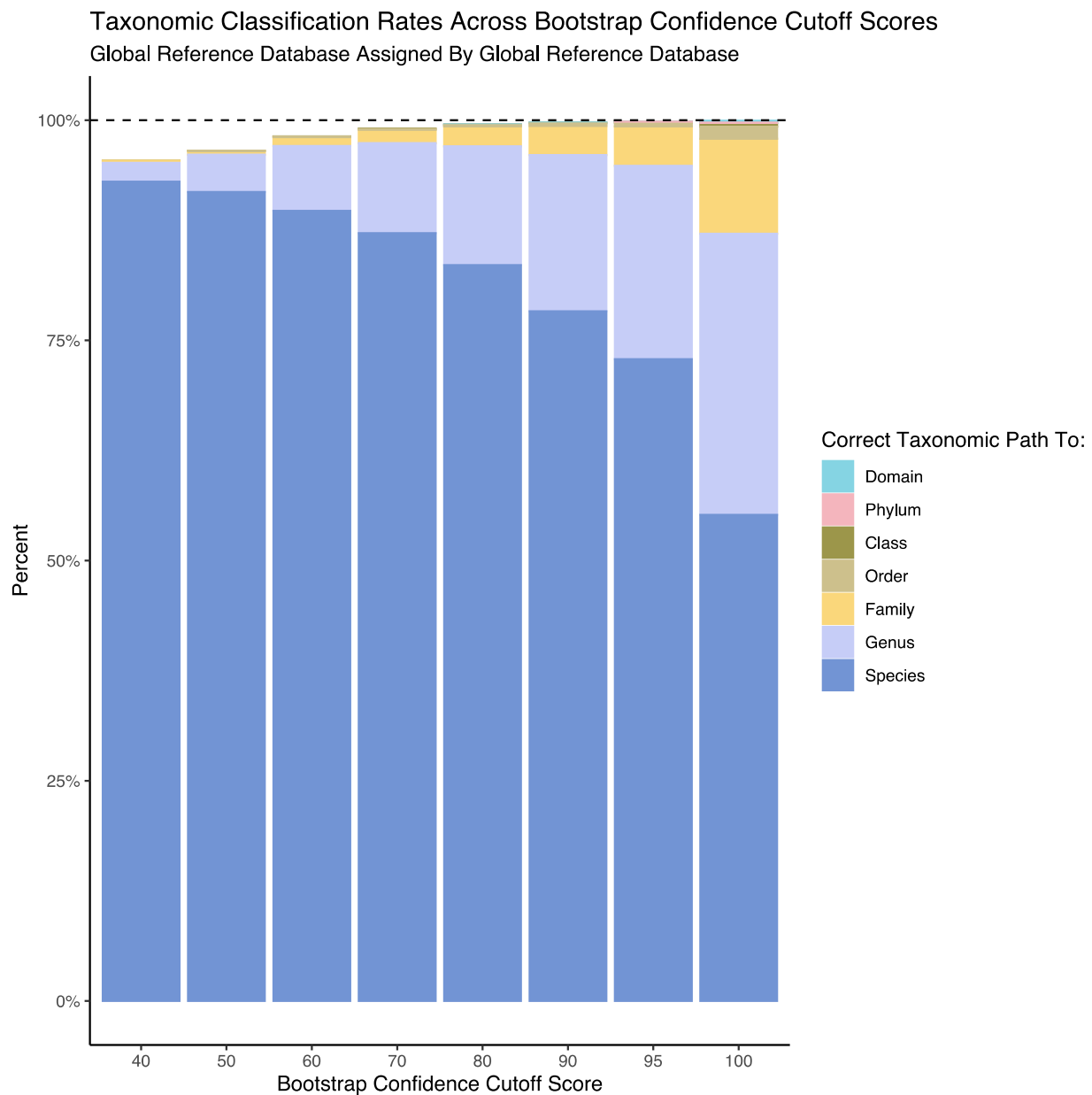


1096

1097 **Figure 1. Effect of TAXXI Bootstrap Confidence Cutoff Scores on Taxonomic Assignment**

1098 **Metrics.** Taxonomy cross-validation by identity (TAXXI) results for taxonomic assignments generated by using  
 1099 the global database as a reference to annotate the sequences in that same database. Accuracy, true positive rate,  
 1100 sensitivity, and misclassification increased with relaxed bootstrap confidence cutoff scores. Under-classification and  
 1101 specificity decreased with relaxed bootstrap confidence cutoff scores. Results for each taxonomic rank are colored.

1102



1103

1104 **Figure 2. Taxonomic Classification Rates Across Bootstrap Confidence Cutoff Scores.**

1105 Results from taxonomy cross-validation by identity (TAXXI) using the global database as the reference to assign

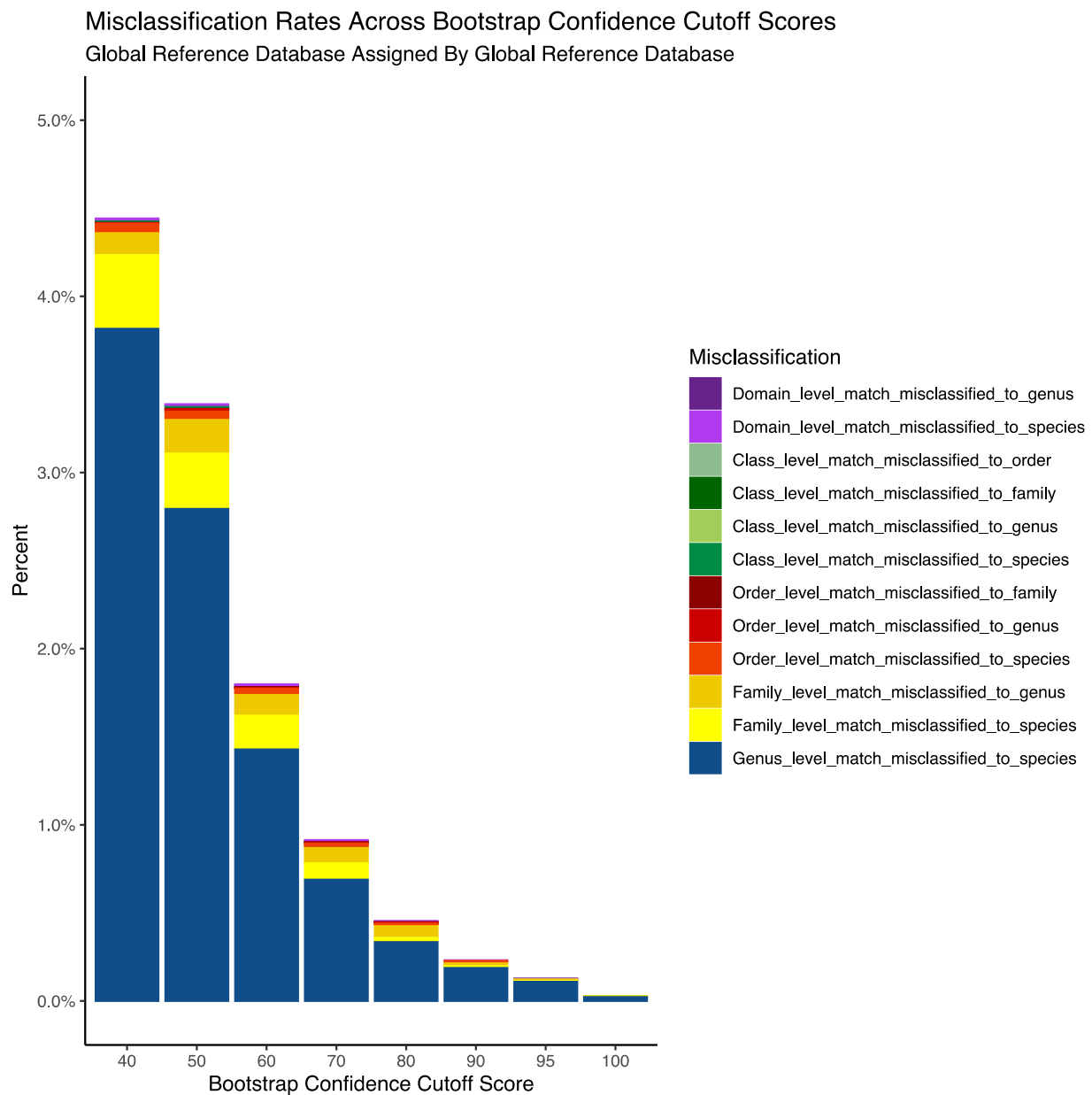
1106 taxonomy to all sequences in that database. Correct species level matches increase with more relaxed bootstrap

1107 confidence cutoff scores. Correct taxonomic level matches are colored by the lowest common ancestor match.

1108 Dotted line indicates 100% and all mismatches were excluded.

1109





1110

1111 **Figure 3. Misclassification Rates Across Bootstrap Confidence Cutoff Scores.** Results from

1112 taxonomy cross-validation by identity (TAXXI) using the global reference database to assign taxonomy to all

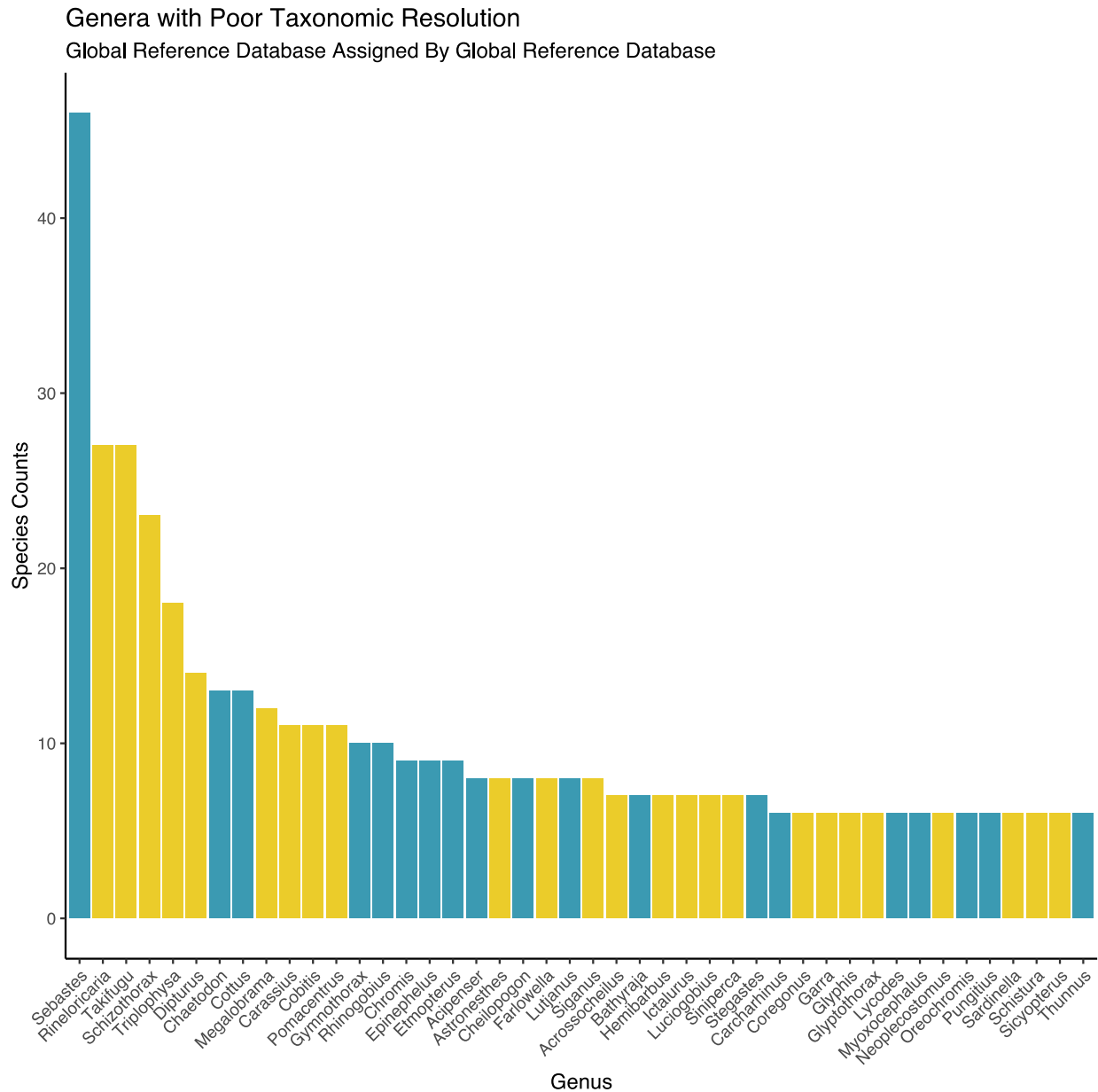
1113 sequences in that database. Misclassification increased with relaxed bootstrap confidence cutoff scores.

1114 Misclassification types are colored.

1115

1116

1117



1118

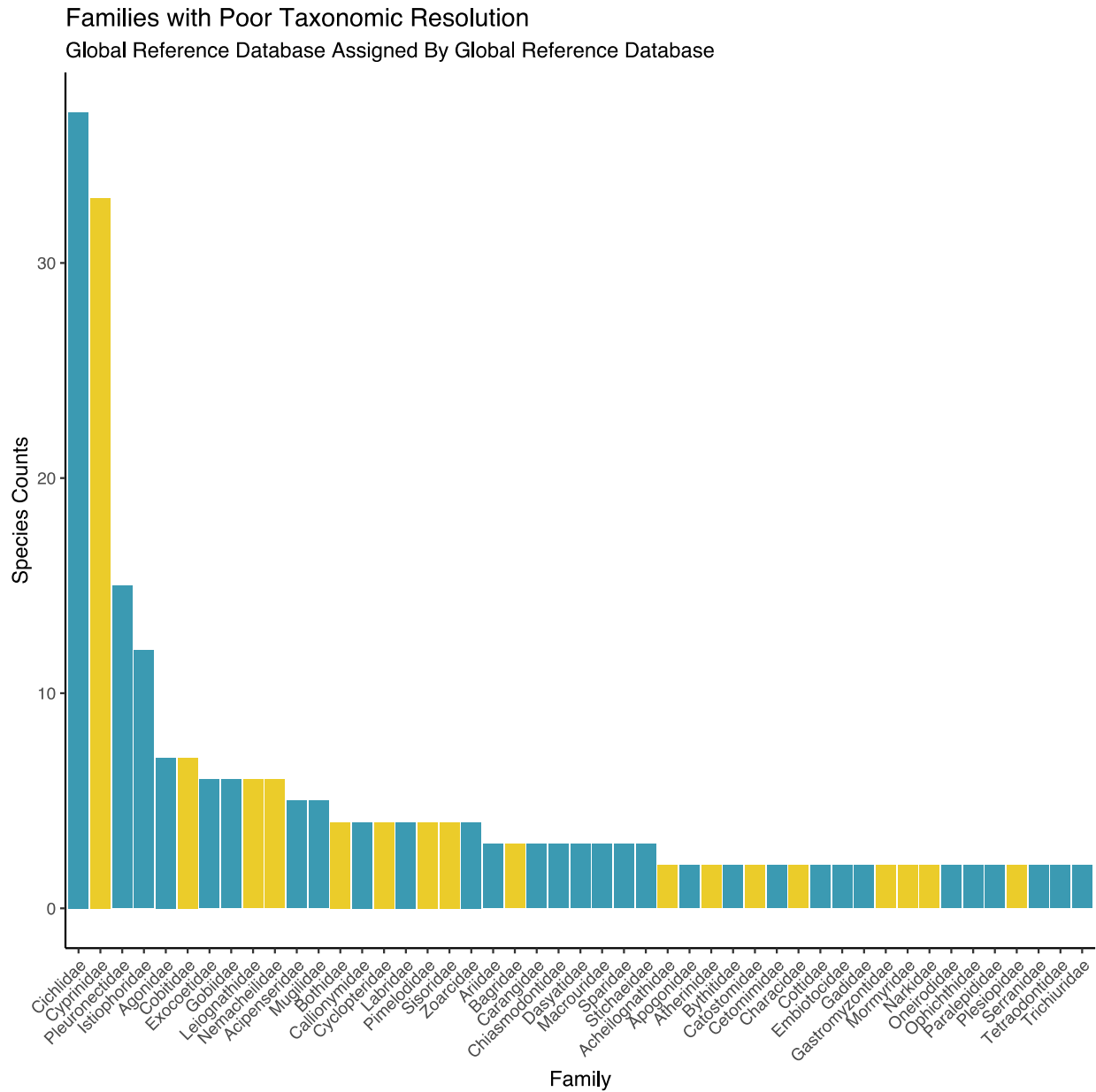
1119 **Figure 4. Genera with Poor Taxonomic Resolution.** Genera poorly resolved to the species level by the

1120 MiFish *I2S* barcode based on results from taxonomy cross-validation by identity (TAXXI) using the global

1121 reference database to assign taxonomy to all sequences in that database using a bootstrap confidence cutoff of 60.

1122 Genera in blue occur in the California Current Large Marine Ecosystem.

1123



1124

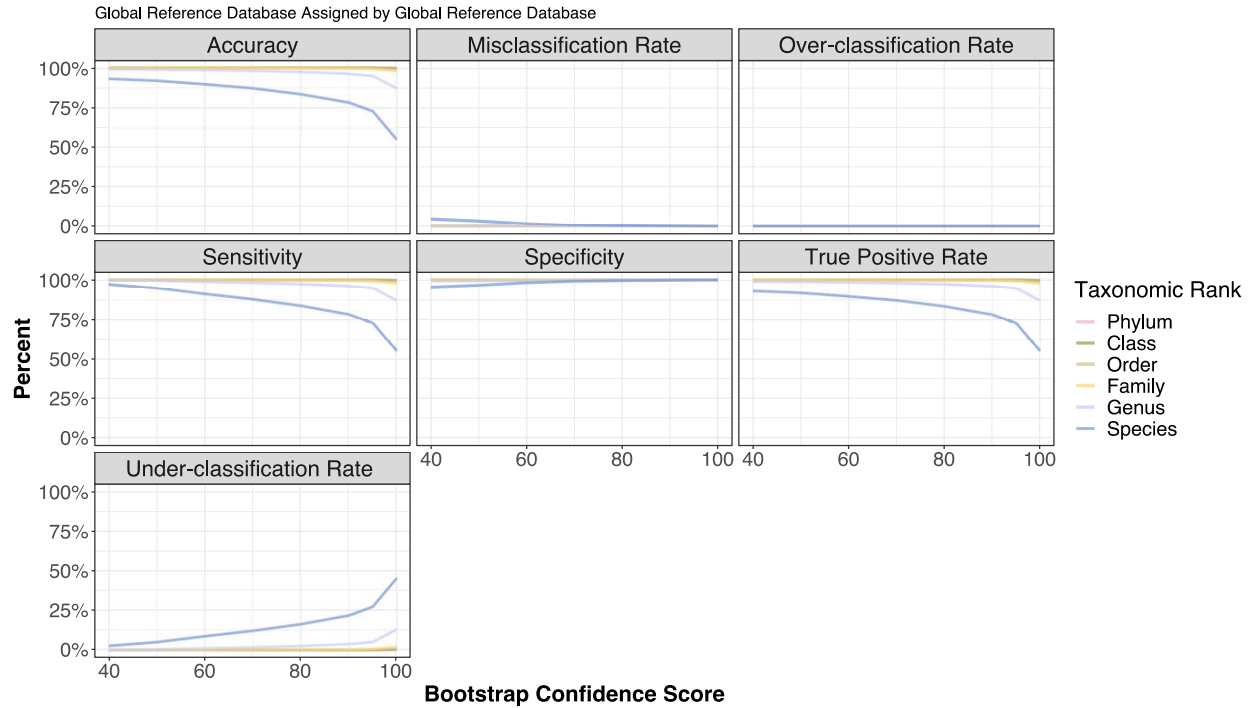
1125

1126 **Figure 5. Families with Poor Taxonomic Resolution.** Families poorly resolved to the species level by

1127 the MiFish 12S barcode based on results from taxonomy cross-validation by identity (TAXXI) using the global

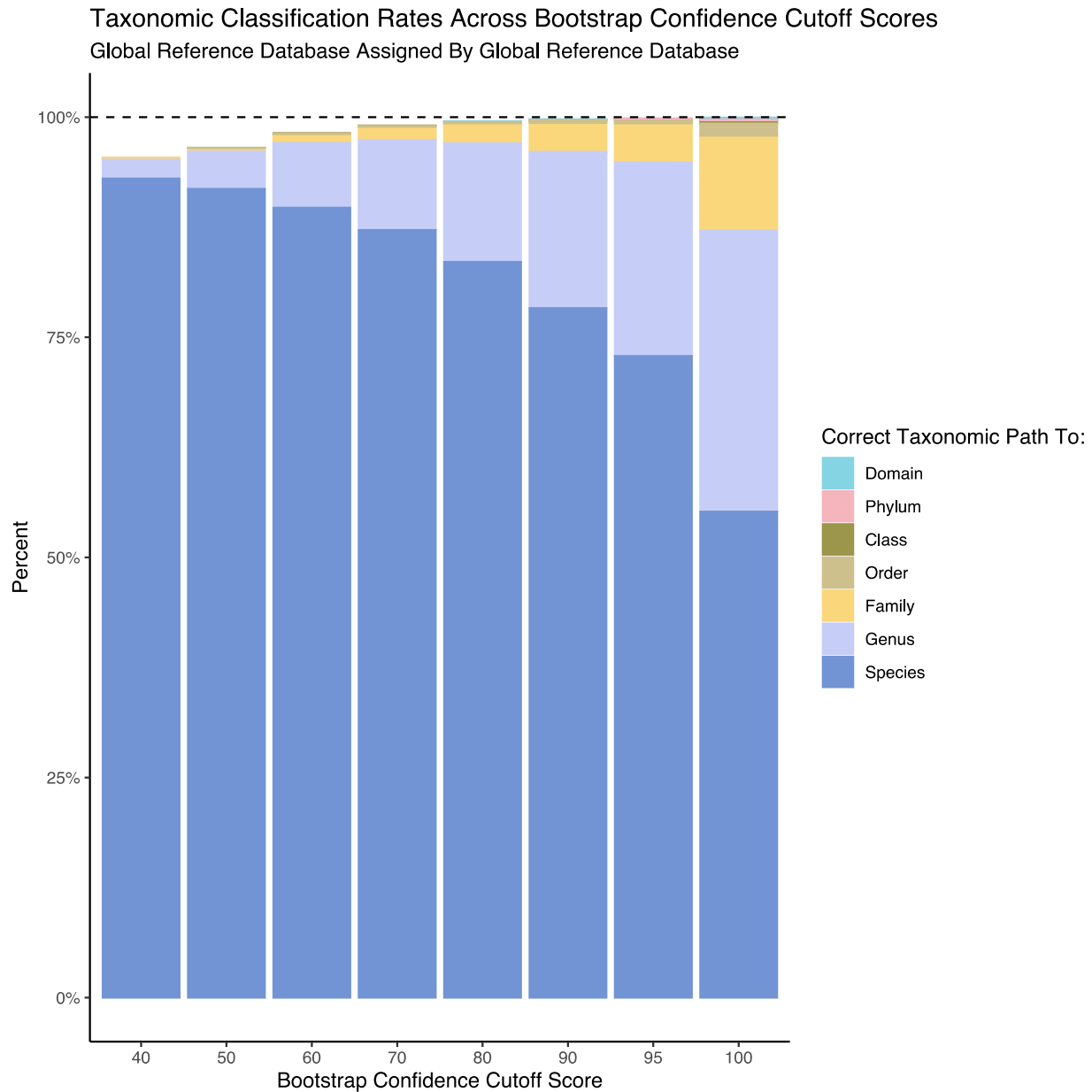
- 1128 reference database to assign taxonomy to all sequences in that database using a bootstrap confidence cutoff (BCC)
- 1129 of 60. Families in blue that occur in the California Current Large Marine Ecosystem.

For Review Only



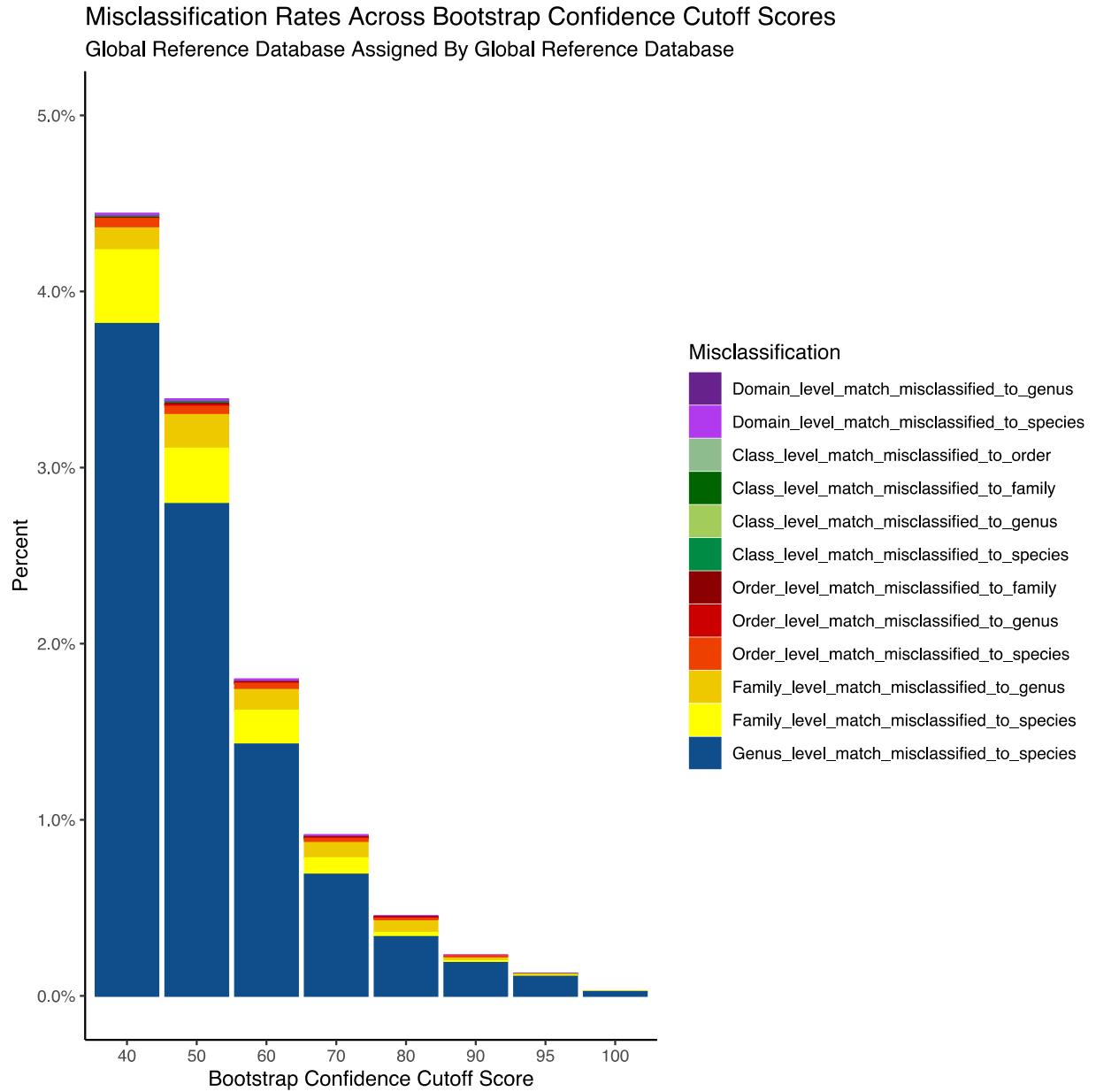
**Figure 1. Effect of TAXXI Bootstrap Confidence Cutoff Scores on Taxonomic Assignment**

**Metrics.** Taxonomy cross-validation by identity (TAXXI) results for taxonomic assignments generated by using the global database as a reference to annotate the sequences in that same database. Accuracy, true positive rate, sensitivity, and misclassification increased with relaxed bootstrap confidence cutoff scores. Under-classification and specificity decreased with relaxed bootstrap confidence cutoff scores. Results for each taxonomic rank are colored.

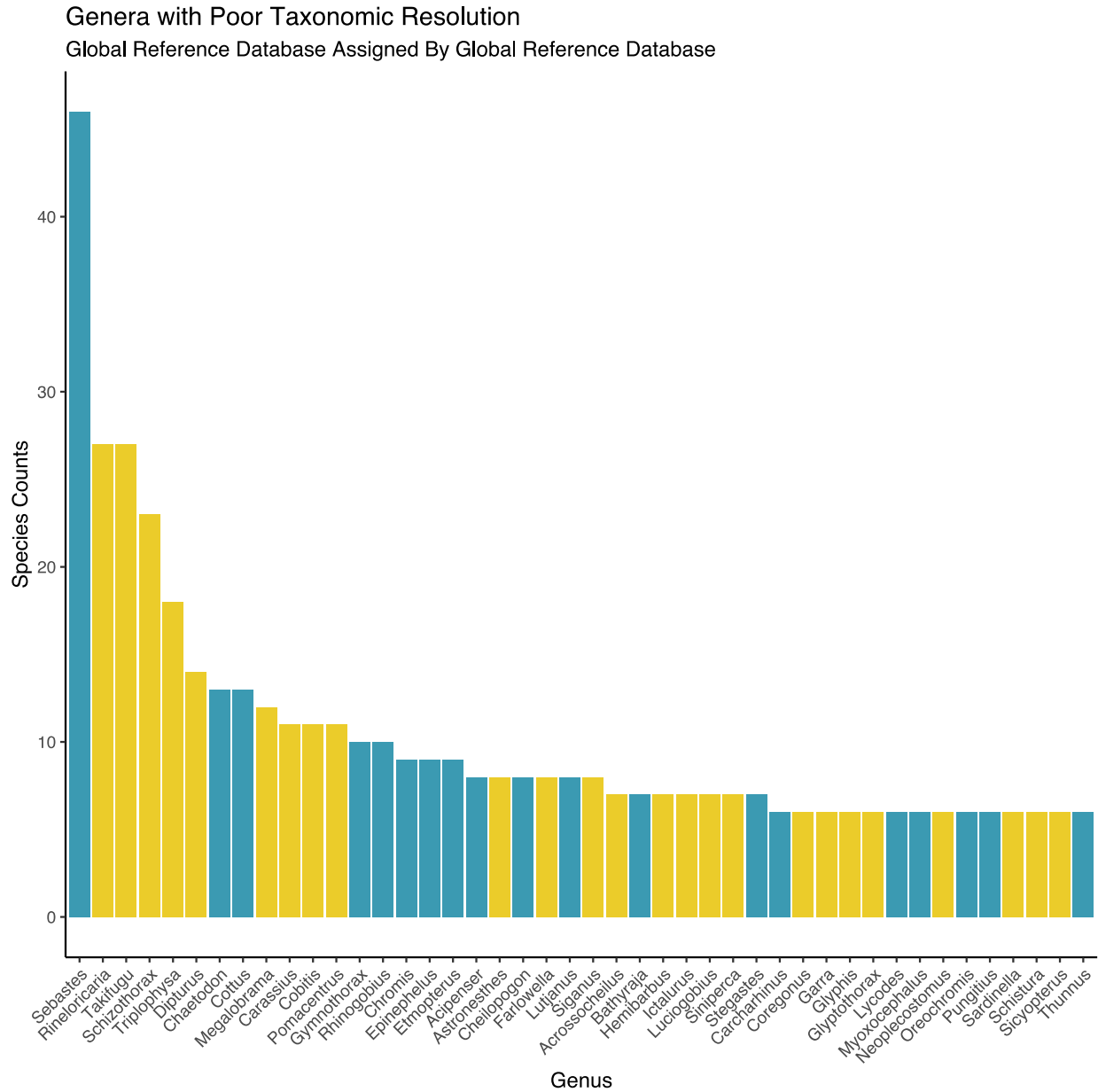


**Figure 2. Taxonomic Classification Rates Across Bootstrap Confidence Cutoff Scores.**

Results from taxonomy cross-validation by identity (TAXXI) using the global database as the reference to assign taxonomy to all sequences in that database. Correct species level matches increase with more relaxed bootstrap confidence cutoff scores. Correct taxonomic level matches are colored by the lowest common ancestor match. Dotted line indicates 100% and all mismatches were excluded.

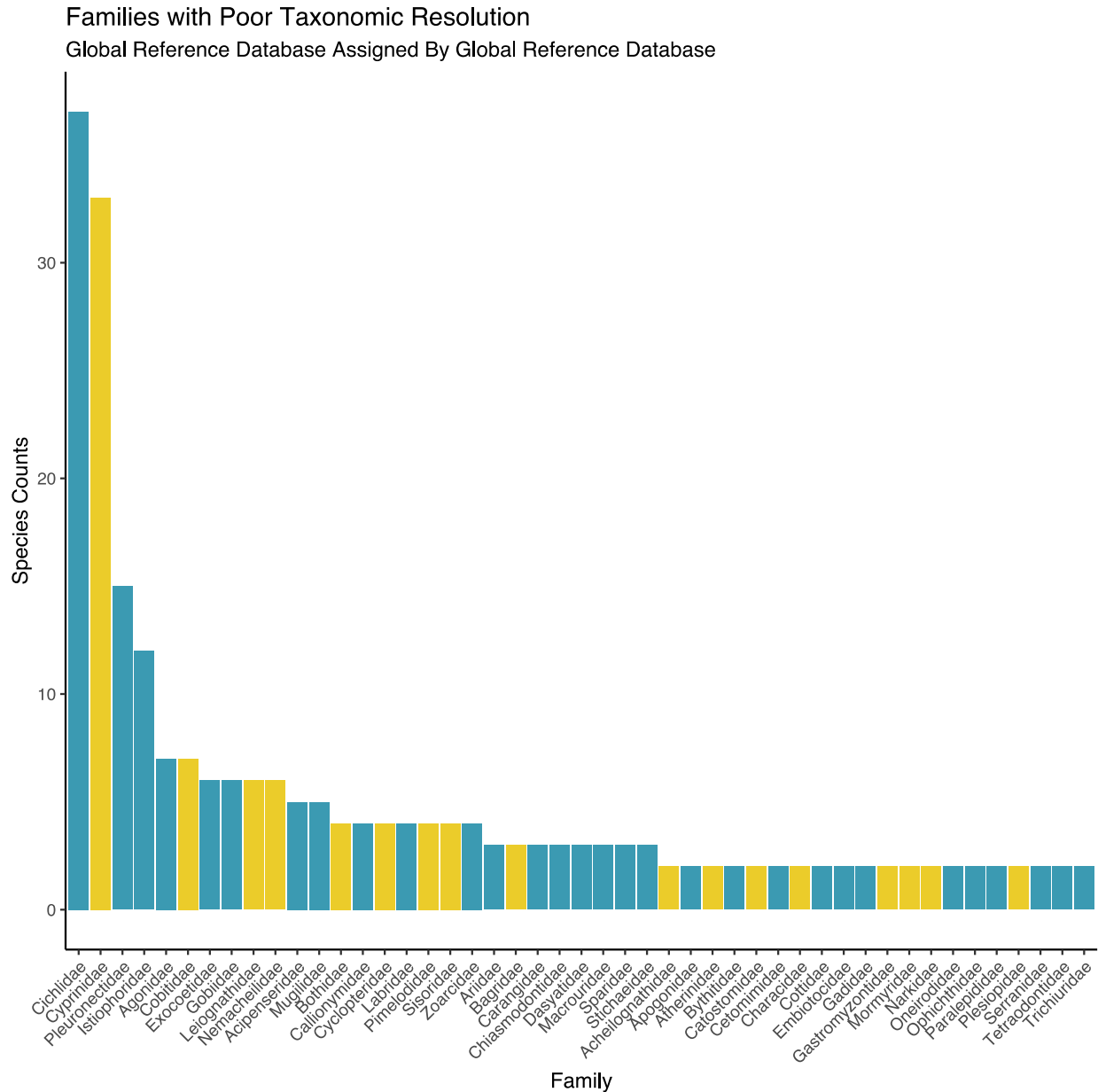


**Figure 3. Misclassification Rates Across Bootstrap Confidence Cutoff Scores.** Results from taxonomy cross-validation by identity (TAXXI) using the global reference database to assign taxonomy to all sequences in that database. Misclassification increased with relaxed bootstrap confidence cutoff scores. Misclassification types are colored.



**Figure 4. Genera with Poor Taxonomic Resolution.** Genera poorly resolved to the species level by the MiFish *12S* barcode based on results from taxonomy cross-validation by identity (TAXXI) using the global reference database to assign taxonomy to all sequences in that database using a bootstrap confidence cutoff of 60. Genera in blue occur in the California Current Large Marine Ecosystem.





1

2

3 **Figure 5. Families with Poor Taxonomic Resolution.** Families poorly resolved to the species level by

4 the MiFish 12S barcode based on results from taxonomy cross-validation by identity (TAXXI) using the global

5 reference database to assign taxonomy to all sequences in that database using a bootstrap confidence cutoff (BCC)

6 of 60. Families in blue that occur in the California Current Large Marine Ecosystem.

7

**Table 1. Summary of Cross Validation Results.** Comparison of performance metrics for taxonomic assignments using the global database as a reference to annotate sequences in the global database (global-global) [ (test database-training database)], the regional database (global-regional), and using the regional database as a reference to annotate sequences in itself (regional-regional). Reporting metrics calculated using a taxonomic cutoff score of 60.

<b>Metric</b>	<b>Global-Global</b>	<b>Global-Regional</b>	<b>Regional-Regional</b>
<b>Under-classification Rate</b>	8.6%	11.8%	7.8%
<b>Misclassification Rate</b>	1.7%	1.8%	1.3%
<b>Over-classification Rate</b>	0.0%	0.0%	0.0%
<b>Accuracy</b>	89.7%	86.5%	90.9%
<b>True Positive Rate</b>	89.7%	86.5%	90.9%
<b>Sensitivity</b>	91.3%	88.0%	92.1%
<b>Specificity</b>	98.3%	98.2%	98.7%

Table 2. Summary of Seawater eDNA Metabarcoding Taxonomic Assignments for Tested Reference Databases.

Metric		Reference Database		
		CRUX-GenBank	Global	Regional
<b>Database</b>	Reference Barcode Origin	GenBank	GenBank + Generated	GenBank + Generated
	Species Included	<b>All</b>	<b>All</b>	<b>California Fishes</b>
<b>Reads</b>	Total Reads	330,877		
	Assigned to NA	81,014	81,002	81,006
	Assigned to Class Level	54,090	-	-
	Assigned to Order Level	727	-	-
	Assigned to Family Level	1,286	1,409	131
	Assigned to Genus Level	952	1,068	1,063
	Assigned to Species Level	192,808	247,398	248,677
<b>ASVs</b>	Total ASVs	341		
	Assigned to NA	172	169	170
	Assigned to Class Level	12	-	-
	Assigned to Order Level	3	-	-
	Assigned to Family Level	5	13	11
	Assigned to Genus Level	4	6	4
	Assigned to Species Level	145	153	156
<b>Taxonomy</b>	Unique Families Identified	31	28	27
	Unique Genera Identified	39	38	39
	Unique Species Identified	38	38	37
	CA Native Species	25	36	37
	Avg. ASVs Per Species	3.8	4.1	4.2