

1 **Profile hidden Markov model sequence analysis can help remove putative**  
2 **pseudogenes from DNA barcoding and metabarcoding datasets**

3

4 Porter, T. M.<sup>1</sup>, Hajibabaei, M.<sup>1</sup>

5

6 <sup>1</sup>University of Guelph, Centre for Biodiversity Genomics and Department of Integrative  
7 Biology, 50 Stone Road East, Guelph, ON Canada

8

9 Corresponding Author:

10 T.M. Porter, terrimporter@gmail.com

11

12

13 **Abstract**

14 **Background:** Pseudogenes are non-functional copies of protein coding genes that  
15 typically follow a different molecular evolutionary path as compared to functional genes.  
16 The inclusion of pseudogene sequences in DNA barcoding and metabarcoding analysis  
17 can lead to misleading results. None of the most widely used bioinformatic pipelines  
18 used to process marker gene (metabarcoding) high throughput sequencing data  
19 specifically accounts for the presence of pseudogenes in protein-coding marker genes.  
20 The purpose of this study is to develop a method to screen for obvious pseudogenes in  
21 large COI metabarcoding datasets. We do this by: 1) describing gene and pseudogene  
22 characteristics from a simulated DNA barcode dataset, 2) show the impact of two  
23 different pseudogene removal methods on mock metabarcoding datasets with simulated  
24 pseudogenes, and 3) incorporate a pseudogene filtering step in a bioinformatic pipeline  
25 that can be used to process Illumina paired-end COI metabarcoding sequences. Open  
26 reading frame length and sequence bit scores from hidden Markov model (HMM) profile  
27 were used to detect pseudogenes.

28 **Results:** Our simulations showed that it was more difficult to identify pseudogenes from  
29 shorter amplicon sequences such as those typically used in metabarcoding (~300 bp)  
30 compared with full length DNA barcodes that are used in construction of barcode  
31 libraries (~ 650 bp). It was also more difficult to identify pseudogenes in datasets where  
32 there is a high percentage of pseudogene sequences. We show that existing  
33 bioinformatic pipelines used to process metabarcoding sequences already remove some  
34 apparent pseudogenes, especially in the rare sequence removal step, but the addition  
35 of a pseudogene filtering step can remove more.

36 **Conclusions:** The combination of open reading frame length and hidden Markov model  
37 profile analysis can be used to effectively screen out obvious pseudogenes from large  
38 datasets. There is more to learn from COI pseudogenes such as their frequency in  
39 DNA barcode and metabarcoding studies, their taxonomic distribution, and evolution.  
40 Thus, we encourage the submission of verified COI pseudogenes to public databases to  
41 facilitate future studies.

42

### 43 **Key words**

44 Nuclear encoded mitochondrial sequences (nuMT), pseudogene, bioinformatics, COI  
45 mtDNA, DNA barcode, metabarcode, hidden Markov model

46

47

### 48 **Introduction**

49 The mitochondrial cytochrome c oxidase subunit 1 gene, COI, is the official animal  
50 barcode marker and large reference databases are available to help identify COI  
51 metabarcode sequences from soil, water, sediments, or mixed communities such as  
52 those collected from traps [1–3]. Crucially, the COI barcode marker is also a protein  
53 coding gene. This is in contrast with the ribosomal DNA markers typically used for  
54 marker gene studies of prokaryotes or fungi [4–6]. Until recently, the methodology and  
55 bioinformatic pipelines for processing protein coding markers such as COI for animals,  
56 the maturase K gene (matK), or the ribulose biphosphate carboxylase large chain gene  
57 (rbcL) for plants have been treated in very much the same way, even using the same  
58 popular pipelines such as those used to process ribosomal RNA genes.

59 Innovative methods of processing COI sequence data has arisen in recent years.  
60 For example, COI marker analysis need not be limited to operational taxonomic units  
61 (OTUs), but may also include the use of exact sequence variant (ESV) analysis for  
62 improved taxonomic resolution and permit intraspecific phylogeographic analyses [7–  
63 10]. Additionally bioinformatic tools to remove pseudogenes and noise from COI  
64 datasets have become available [11–13]. There are currently few options, however, to  
65 process COI metabarcode reads that specifically handle COI pseudogenes also known  
66 as nuclear encoded mitochondrial sequences (nuMTs). COI pseudogenes have been  
67 discussed in the literature largely with regards to COI barcoding efforts and only  
68 recently have tools appropriate for handling large batches of COI sequences recently  
69 become available [14–17].

70 Pseudogenes are copies of mitochondrial DNA that have been inserted into the  
71 nuclear genome [18]. The mechanism for this is uncertain but may involve the  
72 incorporation of mtDNA during the repair of chromosomal double strand breaks [19].  
73 Some mitochondrial pseudogenes are ‘dead on arrival’ due to the different genetic code  
74 in the nuclear genome [20]. If the pseudogene has only accumulated a few mutations,  
75 the sequence may closely resemble that of a functional COI gene with no frameshift or  
76 internal stop codons and may be referred to as a cryptic pseudogene [21]. More  
77 apparent pseudogenes, on the other hand, may exhibit stark changes in codon usage  
78 bias, transition:transversion ratios, GC content, decreased length, and have unexpected  
79 phylogenetic placement [18]. Since the primers used for PCR will bind to paralogous  
80 regions in pseudogenes, they will amplify nuMTs in addition to or even preferentially to  
81 the target mitochondrial sequence [18, 22]. Including unknown pseudogenes in

82 phylogenetic, biodiversity, or population analyses may introduce noise into analyses,  
83 leading to overestimates of haplotype or species richness, or may lead to misleading  
84 identifications or relationships [14, 16, 23–26].

85 The methods needed to detect different types of pseudogenes will vary depending  
86 on whether or not many changes have accumulated. Cryptic pseudogenes may be  
87 identified by examining raw Sanger chromatograms, similar to looking for evidence for  
88 heteroplasmy, by looking for double peaks. The whole gene region may be examined  
89 looking for the presence of the control region and stop codon. Conserved regions such  
90 as in the inner mitochondrial membrane alpha helices can be examined for changes  
91 [27]. More obvious pseudogenes may accumulate substitutions equally in non-  
92 synonymous and synonymous regions indicating balanced positive and negative  
93 selection at sites across the gene copy or relaxed conservation ( $dN/dS$  ratios  $\sim 1$ ). This  
94 is in contrast with a functional COI gene where substitutions tend to occur in non-  
95 synonymous sites so as to preserve amino acid composition and protein structure and  
96  $dN/dS$  ratios are expected to be  $< 1$ . The result of relaxed purifying selection is the  
97 accumulation of indels, frameshifts, and/or the introduction of premature stop codons.  
98 The objective of this work is to develop methods to remove such apparent pseudogenes  
99 from large COI sequence datasets.

100

## 101 **Bioinformatic Methods**

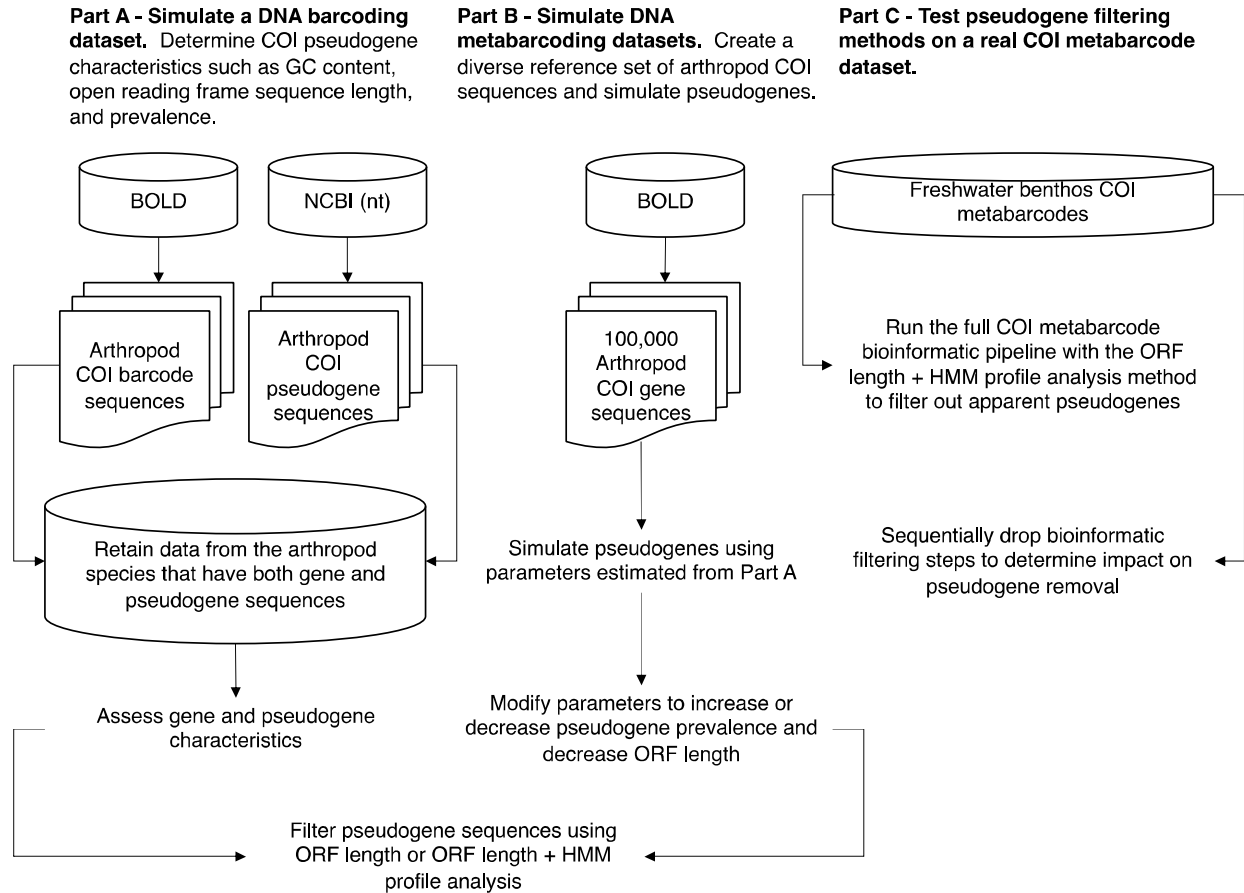
102 We used three approaches in this study: A) We simulated a DNA barcode  
103 dataset by compiling a set of annotated COI genes and pseudogenes from the Barcode  
104 of Life Data System (BOLD) and the National Center for Biotechnology Information

105 (NCBI) nucleotide (nt) database for the same set of 10 species; B) We created mock  
106 COI metabarcode datasets by mining sequences from BOLD and simulating  
107 pseudogenes, and C) We tested a pseudogene filtering method on a previously  
108 published freshwater benthos COI metabarcode dataset (Figure 1).

109

110 **Figure 1. Overview of methods to determine COI pseudogene characteristics and**  
111 **test methods for pseudogene removal.** Dataflow for our A) simulated DNA barcode  
112 dataset, B) simulated metabarcode datasets, and C) real freshwater COI metabarcode  
113 dataset. Abbreviations: BOLD = Barcode of Life Data System; COI = cytochrome c  
114 oxidase subunit I mtDNA gene; HMM = hidden Markov model; NCBI = National Centre  
115 for Biotechnology Information; nt = nucleotide; ORF = open reading frame.

116



117

118

### 119 *Part A: Simulating a DNA barcoding dataset*

120 To simulate a DNA barcoding dataset where multiple sequences are generated  
121 for the same species, we retrieved high quality sequences from BOLD and known  
122 pseudogenes mined from the NCBI nucleotide database for the same set of species.  
123 Sequences from the BOLD data releases were obtained from  
124 <http://v3.boldsystems.org/index.php/datarelease> . Nucleotide sequences for arthropods  
125 were selected, ensuring that there were no ambiguities in the nucleotide sequences. If  
126 either the nucleotide sequence or amino acid sequence were missing, then the record  
127 was discarded. A FASTA file containing arthropod COI pseudogenes was obtained

128 from the NCBI nucleotide database using an Ebot script with the search term  
 129 “Arthropoda[ORGN] AND pseudogene[TITL] AND (COI[GENE] OR CO1[GENE] OR  
 130 coxI[GENE] OR cox1[GENE]) AND 50:2000[SLEN]”. [28] A few records had to be edited  
 131 by hand to isolate the sequence region associated with the COI pseudogene. We  
 132 retrieved 481 COI nucleotide sequences from BOLD and 112 COI pseudogene  
 133 nucleotide sequences from the NCBI nucleotide database from the same 10 species  
 134 (Table 1). This dataset is further described in Table S1 showing proportion of  
 135 pseudogenes, average length, and average GC content. On average, the length and  
 136 GC content of pseudogenes from these 10 species are slightly shorter and lower,  
 137 respectively, than for COI gene sequences.

138

139 **Table 1: Summary of a simulated DNA barcoding dataset containing known**  
 140 **arthropod COI pseudogenes**

141

Class	Order	Species	Gene sequences (% of total)	Pseudogene sequences (% of total)	Subtotals
Insecta	Coleoptera	<i>Xylosandrus germanus</i>	33	1	34
Insecta	Hemiptera	<i>Bemisia tabaci</i>	252	7	259
Insecta	Hemiptera	<i>Trialeurodes vaporariorum</i>	3	1	4
Insecta	Hemiptera	<i>Triatoma dimidiata</i>	9	1	10
Insecta	Hymenoptera	<i>Ectatomma gibbum</i>	6	1	7
Insecta	Hymenoptera	<i>Halictus rubicundus</i>	29	2	31
Insecta	Hymenoptera	<i>Melissotarsus insularis</i>	135	79	214
Insecta	Orthoptera	<i>Cyphoderris monstrosa</i>	7	14	21
Collembola	Entomobryomorpha	<i>Lepidocyrtus cyaneus</i>	5	1	6
Malacostraca	Decapoda	<i>Goneplax rhomboides</i>	2	5	7
		Subtotals	481 (81)	112 (19)	593

142



143

144 GC content for COI gene and pseudogene sequences were assessed in R using  
145 the 'seqinr' package [29]. We pooled all the sequences together, then proceeded to  
146 filter out just the pseudogene sequences using two different methods:

147 The first method we used to remove pseudogenes involved screening out ESVs  
148 with outlier open reading frame lengths that were very short or very long (SCVUC  
149 v4.1.0). This was done by translating arthropod ESVs using ORFfinder v0.4.3 into  
150 every possible open reading frame on the plus strand, ignoring nested ORFs, minimum  
151 length set to 30. The longest nucleotide (nt) ORFs were retained. Outliers, putative  
152 pseudogenes or genuine sequences with PCR/sequencing errors, were identified as  
153 sequences shorter than the 25<sup>th</sup> percentile ORF length - (1.5 \* interquartile length) and  
154 longer than the 75<sup>th</sup> percentile ORF length + (1.5 \* interquartile length).

155 The second method we used to remove pseudogenes involved profile hidden  
156 Markov model (HMM) analysis (SCVUC v4.3.0). This was done by creating a profile  
157 HMM based on BOLD arthropod barcode sequences using HMMER v3.3 available from  
158 <http://hmmer.org> . From the BOLD data releases iBOL phase 0.50 to 6.50, we retrieved  
159 all arthropod barcodes 600-700 bp in length. We sorted these sequences by  
160 decreasing length using the 'sortbylength' command in VSEARCH. We reduced the  
161 dataset size by clustering by 80% sequence similarity using the 'cluster\_size' command  
162 and retaining the centroids sequences. As described above, arthropod ESVs were  
163 translated and the longest open reading frames were retained for both nucleotide and  
164 amino acid (aa) sequences. The amino acid ORFs were aligned with MAFFT v7.455  
165 using the 'auto' setting [30]. The nucleotide ORFs were also mapped to the amino acid

166 alignment using TRANALIGN (EMBOSS v6.6.0.0) specifying the invertebrate  
167 mitochondrial genetic code [31]. The FASTA file comprised of 6,162 amino acid  
168 sequences was converted to Stockholm format. This reference alignment was turned  
169 into a model that describes the probabilities for travelling a path along the length of the  
170 alignment that moves through match, insert, or deletion states. HMMER was used to  
171 build this nucleotide arthropod COI profile hidden markov model (HMM) using the  
172 'hmmbuild' command. The HMM was indexed using the 'hmmcompress' command.  
173 Individual arthropod amino acid ORFs were then compared with the profile HMM using  
174 the 'hmmsearch' command. One of the hmmsearch outputs is a log odds ratio score (bit  
175 score) that compares the likelihood of the query sequence given the model to the  
176 likelihood of the query sequence given a random sequence model. When a COI gene is  
177 used as the query, we expected a high bit score; whereas when an obvious COI nuMT  
178 is used as the query, we expected a low bit score. In this way, putative pseudogenes or  
179 genuine sequences with PCR/sequencing errors were identified as amino acid ORFs  
180 with short outlier HMMER scores.

181 We also calculated the number of substitutions per non-synonymous and  
182 synonymous sites. Gene sequences and pseudogene sequences were analyzed  
183 separately as follows: Amino acid ORFs were aligned using MAFFT v7.455 using the  
184 'auto' setting. A codon alignment was created using TRANALIGN (EMBOSS) by  
185 mapping the nucleotide ORFs to the amino acid alignment using the invertebrate  
186 mitochondrial genetic code. We used the package 'ggplot2' in Rstudio to create all plots  
187 [32–34]. We used the 'seqinr' function 'kaks' to calculate the number of substitutions for  
188 non-synonymous and synonymous sites [29]. Before calculating dN/dS ratios, we

189 excluded pairwise sequence comparisons where the number of substitutions per  
190 synonymous site was  $< 0.01$  (sequences too similar to yield reliable dN/dS) or  $> 2$  (too  
191 many substitutions, near saturation, to yield a reliable dN/dS).

192 To assess how pseudogene sequences could be (mis)identified using the top  
193 BLAST hit method, we used the megablast algorithm to find the most similar sequence  
194 in the NCBI nucleotide sequence database [35]. We used this method to verify that the  
195 expected species was a top match (skipping over the top match if it was the same as  
196 the query sequence or if it was an obvious contaminant) and whether or not the top  
197 match was to a gene or pseudogene sequence in the reference database. To further  
198 visualize phylogenetic divergence between gene and pseudogene sequences for each  
199 species, we aligned nucleotide sequences with MAFFT using the 'auto' setting. The  
200 'fdnadist' Phylip method in the EMBOSS package was used to calculate distances using  
201 the Kimura 2-parameter (K2P) model of nucleotide sequence evolution [36, 37]. A  
202 neighbor joining tree was saved in Newick format using the 'fneighbor' Phylip method in  
203 EMBOSS. Statistical support at nodes was calculated by bootstrapping the multiple  
204 sequence alignment 1000 times using the 'fseqboot' Phylip method in the EMBOSS  
205 package then K2P distances and neighbor joining trees were constructed as described  
206 above. A majority rule consensus tree was constructed using the Phylip program  
207 'consense' [37]. Bootstrap values from the consensus tree were mapped to the  
208 phylogram using TreeGraph2 v2.15.0-887 [38]. The tree was mid-point rooted and  
209 nodes rotated or collapsed where necessary to improve readability using FigTree v1.4.4  
210 available from <http://tree.bio.ed.ac.uk/software/figtree/>. Further minor editing to

211 improve readability was performed using Inkscape v1.0.1 available from  
212 <https://inkscape.org/> .

213

#### 214 *Part B: Simulating community sequence data*

215 To test our pseudogene filtering methods on a more taxonomically diverse  
216 community of arthropods, we performed a simulation study. We created an arthropod  
217 COI community based on 100,000 sequences randomly sampled from BOLD. We  
218 manipulated this mock community in different ways described below. In our first mock  
219 community, based on our simulated DNA barcoding results from Part A where ~ 19% of  
220 our dataset represented pseudogenes, we decided to introduce mutations into 19% of  
221 the BOLD sequences. Also based on the results from Part A, we reduced the GC  
222 content in our simulated pseudogenes by 2.5% by replacing G/C bases with an A/T  
223 bases. In our second mock community, we inserted or deleted bases to introduce  
224 frameshift mutations and premature stop codons. To keep the rate of pseudogenization  
225 the same as the first mock community, we introduced indels in 2.5% of the bases in our  
226 simulated pseudogenes. In the third mock community, we split COI barcode sequences  
227 in half to test whether our pseudogene filtering approach would work on shorter barcode  
228 sequences similar in length to those generated in COI metabarcoding studies (~ 300  
229 bp). In a fourth mock community, we doubled the proportion of pseudogenes in the  
230 mock community from 19% to 38%. In the fifth mock community, we halved the  
231 proportion of pseudogenes in the mock community from 19% to 9.5%. Each of these  
232 datasets is further described in Table S1 showing proportion of pseudogenes in the  
233 community, average length, and average GC content.

234

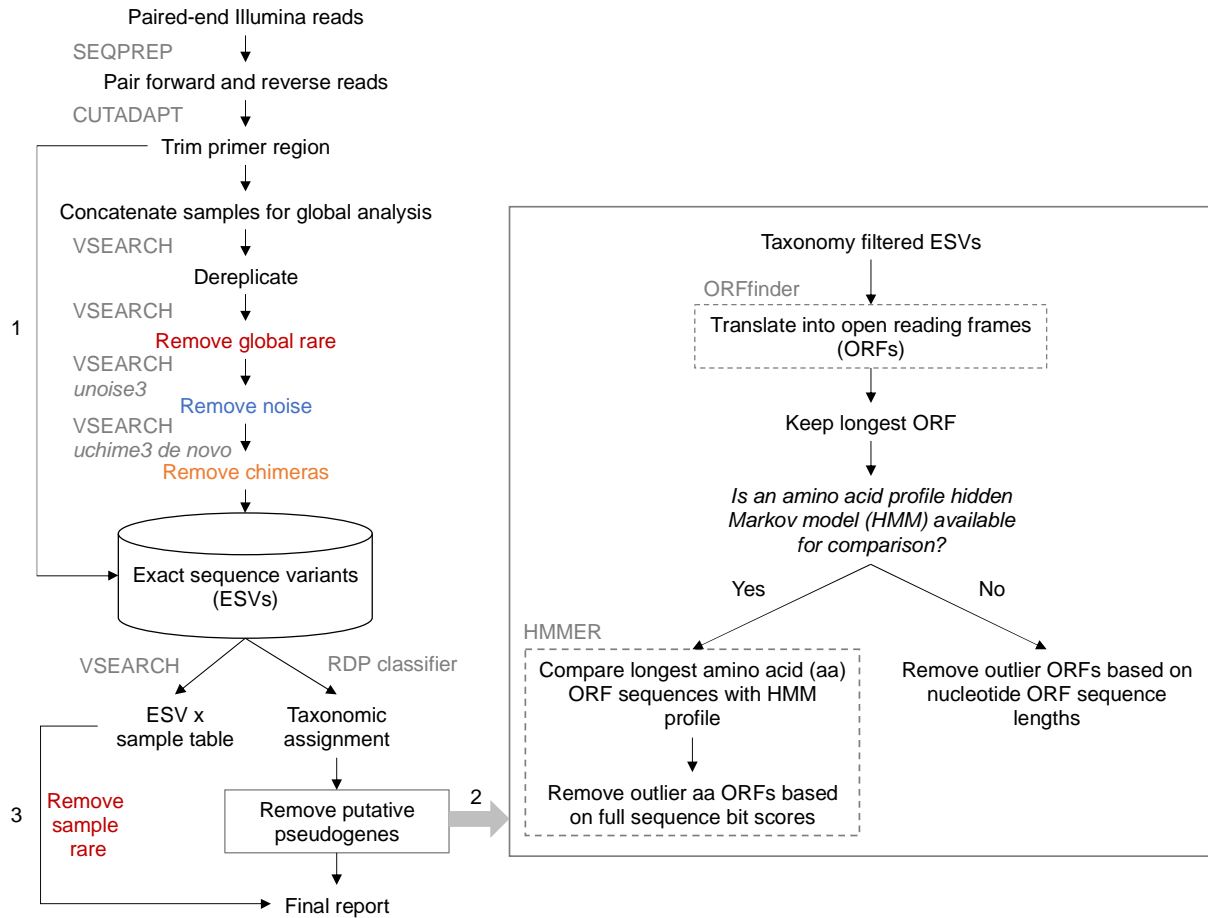
235 *Part C: Test pseudogene filtering methods using a real COI metabarcode dataset*

236 We used a previously published freshwater benthos COI metabarcode dataset to  
237 test our bioinformatic pipeline and two different pseudogene removal strategies [39].  
238 We chose this dataset because it includes results from six different COI amplicons (BR5  
239 [B, ArR5] ~ 310 bp, F230R [LCO1490, 230\_R] ~ 229 bp, ml-jg [mlCOLintF, jgHCO2198]  
240 ~ 313 bp, BF1 [BF1, BR2] ~ 316 bp, BF2 [BF2, BR2] ~ 421 bp, fwh1 [fwhF1, fwhR1] ~  
241 178 bp) currently used in a variety of labs in the freshwater COI metabarcode literature  
242 [40–47]. The primers and their target taxa are listed in Table S2. Each amplicon covers  
243 sites across the COI barcoding region and the mode length ranges from 178 bp (fwh1)  
244 to 421 bp (BF2), averaging ~ 300 bp. The F230R and fwh1 amplicons align to the 5'  
245 end of the barcoding region and the BR5, ml-jg, BF1, and BF2 amplicons align to the 3'  
246 end of the barcode region.

247 A COI metabarcoding bioinformatic pipeline, SCVUC v4.3.0, was used to  
248 process Illumina paired-end reads to output a set of taxonomically assigned ESVs  
249 (available from GitHub at [https://github.com/Hajibabaei-](https://github.com/Hajibabaei-Lab/SCVUC_COI_metabarcode_pipeline)  
250 [Lab/SCVUC\\_COI\\_metabarcode\\_pipeline](https://github.com/Hajibabaei-Lab/SCVUC_COI_metabarcode_pipeline)) (Fig 2). This pipeline runs in a conda  
251 environment using a snakemake pipeline. Conda is an environment and package  
252 manager [48]. It allows most programs and their dependencies to be installed easily  
253 and shared with others. Snakemake is a python-based workflow manager [49]. The  
254 snakefile contains the commands need to run a bioinformatic pipeline. The  
255 configuration file allows users to adjust parameter settings.

256

257 **Fig 2. Overview of metabarcoding bioinformatic pipeline that removes apparent**  
258 **pseudogenes.** The SCVUC pipeline begins with Illumina paired-end reads. Arrow 1  
259 indicates where globally rare sequence clusters are removed and quality trimmed reads  
260 are mapped to denoised exact sequence variants (ESVs) to create a sample x ESV  
261 table that contains read numbers. Arrow 2 indicates where pseudogenes can be  
262 removed using two different approaches. The first method translates ESVs, retains the  
263 longest nucleotide open reading frame (ORF), then removes sequences with very small  
264 or very large outlier lengths. The second method translates ESVs, retains the longest  
265 amino acid open reading frame, does a profile HMM analysis, then removes sequences  
266 with very small outlier full sequence bit scores. Arrow 3 indicates where rare sequence  
267 clusters from each sample are removed and read numbers are mapped to the final  
268 report. The final report contains all ESVs for each sample, read numbers, ORF  
269 sequences, and taxonomic assignments with bootstrap support values.



270

271

272

273

274

275

276

277

278

279

280

Raw paired-end reads are merged using SEQPREP v1.3.2 [50]. This step looks for a minimum Phred quality score of 20 in the overlap region and requires a minimum 25 bp overlap. Primers are trimmed in two steps using CUTADAPT v2.6 requiring a Phred quality score of 20 at the ends to count matches/mismatches, no more than 3 Ns are allowed, and trimmed reads need to be at least 150 bp [51]. Sequence files are combined for a global analysis. Reads are dereplicated using VSEARCH v2.14.1 [52]. Denoised exact sequence variants (ESVs) are also generated using VSEARCH using the unoise3 algorithm [53]. This step clusters reads by 100% sequence identity, removes sequences with predicted errors, and globally rare sequences. Here we define

281 rare sequences as clusters containing only one or two sequences. Putative chimeric  
282 sequences are removed using the uchime3\_denovo algorithm in VSEARCH [54].  
283 Denoised ORFs (ESVs) are taxonomically assigned using a naive Bayesian classifier  
284 trained with a COI reference set comprised of sequences mined from GenBank and the  
285 BOLD data releases [55, 56]. Rare sequences clusters are removed from each sample  
286 before printing the final file.

287 We used the pipeline with the two different pseudogene removal methods  
288 described in Part A. We then modified the pipeline to skip over several steps, one at a  
289 time, to see how this would affect the removal of apparent pseudogenes using the  
290 ORFfinder + profile HMM method: rare sequence removal, noise removal, chimeric  
291 sequence removal.

292

## 293 **Results**

294

295 Our DNA barcode simulation that included 10 species with both gene and  
296 pseudogene sequences allowed us to compare differences in GC content, length, and  
297 dN/dS ratios. In Figure 2, we show that COI pseudogenes tend to have a slightly lower  
298 median GC content, shorter ORF lengths, and shorter full sequence bit score values in  
299 HMM profile analyses. Figure S1 shows how COI genes tend to accumulate  
300 substitutions in synonymous sites where a nucleotide changes does not result in the  
301 change of an amino acid; whereas COI pseudogenes tend to accumulate substitutions  
302 in non-synonymous sites where a nucleotide change results in the change of an amino  
303 acid. After correcting for pairwise comparisons that could yield unreliable dN/dS ratios,



304 where the number of substitutions at synonymous sites is  $< 0.01$  or  $> 2$ , we were only  
305 able to calculate dN/dS for COI gene sequences but not for pseudogene sequences.  
306 Due to the length variation in COI pseudogenes and their resulting ORFs it was difficult  
307 to obtain reliable codon alignments for dN/dS analysis. This method may be more  
308 suitable for detecting cryptic pseudogenes that have open reading frame lengths similar  
309 to functional COI ORFs. Top BLAST hit analysis shows that all pseudogenes had a top  
310 BLAST hit to another sequence from the expected species (92% - 100% identity). In  
311 some cases, the top BLAST match for a known pseudogene was to another COI  
312 sequence annotated as a nuclear copy of a mitochondrial gene. More often, the top  
313 match for a pseudogene was to a COI gene sequence. This indicates that in some  
314 cases, careful analysis of top BLAST hit output could help flag putative pseudogenes.  
315 Figures S2-S11 show COI phylograms for each species. In some cases, pseudogenes  
316 form their own clusters (ex. *Bemisia tabaci*, *Goneplax rhomboides*), often on long  
317 branches (ex. *Bemisia tabaci*, *Xylosandrus germanus*, *Triatoma dimidiata*, *Trialeurodes*  
318 *vaporariorum*, *Goneplax rhomboides*, *Ectatomma gibbum*), but occasionally  
319 pseudogenes are found in clades intermixed with regular genes and little sequence  
320 divergence to distinguish them (ex. *Melissotarsus insularis*, *Lepidocyrtus cyaneus*,  
321 *Halictus rubicundus*, *Cyphoderris monstrosa*).

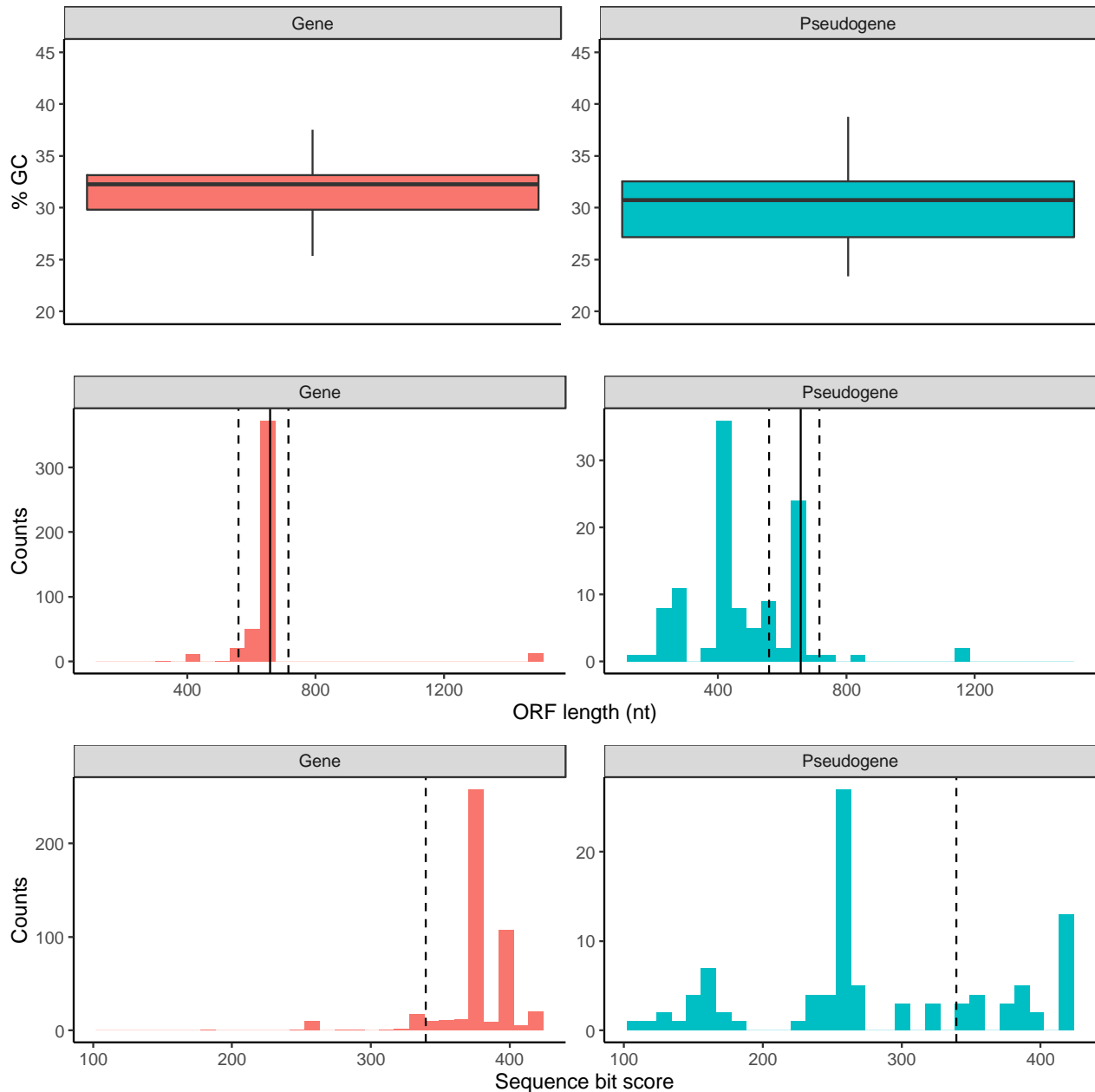
322 Table 2 compares the sensitivity and specificity of two pseudogene removal  
323 methods on this dataset. Figure S12 shows how we calculated sensitivity and  
324 specificity for each pseudogene removal method. Sensitivity refers to the true positive  
325 rate, in this case the number of pseudogenes correctly filtered out of the dataset.  
326 Specificity refers to the true negative rate, in this case, the number of genes correctly

327 retained. For our DNA barcoding simulated dataset including COI gene and  
328 pseudogene sequences from 10 species, sensitivity (73%) is slightly higher for the  
329 ORFfinder + HMM profile analysis pseudogene removal method and the specificity is  
330 the same for each pseudogene removal method (90%).

331

332 **Fig 2. Arthropod COI pseudogenes tend to have lower GC content, shorter open**  
333 **reading frames, and smaller sequence bit scores.** Based on the simulated DNA  
334 barcoding dataset described in Table 1. The top panel shows GC content (%) in gene  
335 and pseudogene sequences. The middle panel shows the sequence length distribution  
336 for the longest retained open reading frame. The solid vertical line indicates the length  
337 of a typical COI barcode at 658 bp. The two vertical dashed lines shows the boundaries  
338 for identifying ORFs with outlier lengths. The bottom panel shows the sequence bit  
339 score distribution after searching our sequences against a COI arthropod nucleotide  
340 profile hidden Markov model. The vertical dashed line shows the boundary for  
341 identifying small outlier scores.

342



343

344

345 **Table 2. Sensitivity and specificity for two pseudogene filtering methods.** We

346 include results from two approaches: Part A) We used a simulated DNA barcoding

347 dataset with COI gene and pseudogene sequences from 10 species, Part B) we

348 simulated pseudogenes from 100,000 BOLD COI sequences. To simulate

349 pseudogenes, we either decreased the %GC content or introduced indels. Sensitivity

350 refers to the true positive rate, our ability to correctly identify known or simulated  
 351 pseudogenes. Specificity refers to the true negative rate, our ability to correctly identify  
 352 real COI sequences (not pseudogenes).

Experiment	Dataset	Type of mutations introduced	Sensitivity (%)		Specificity (%)	
			ORFfinder	ORFfinder + profile HMM analysis	ORFfinder	ORFfinder + profile HMM analysis
Simulated DNA barcoding dataset. COI genes and pseudogenes from 10 species	Full length COI barcode and pseudogene sequences	N/A	70	73	90	90
Simulated metabarcoding dataset	Full length COI barcode and simulated pseudogenes	GC content reduced	31	27	99	~100
Simulated metabarcoding dataset	Full length COI barcode and simulated pseudogenes	Introduced indels	88	94	~100	~100
Simulated metabarcoding dataset	Short COI barcode and simulated pseudogenes	GC content reduced	17** - 50*	6** - 15*	99	~100
Simulated metabarcoding dataset	Short COI barcode and simulated pseudogenes	Introduced indels	42** - 58*	61** - 87*	99	99* - ~100**
Simulated metabarcoding dataset	Full length COI sequences and twice as many pseudogenes	GC content reduced	17	0	99	~100
Simulated metabarcoding dataset	Full length COI sequences and twice as many pseudogenes	Introduced indels	0	0	~100	~100
Simulated metabarcoding dataset	Full length COI sequences and half as many pseudogenes	GC content reduced	39	36	95	96
Simulated metabarcoding dataset	Full length COI sequences and half as many pseudogenes	Introduced indels	95	98	96	99

353 \* 5' fragment

354 \*\* 3' fragment

355

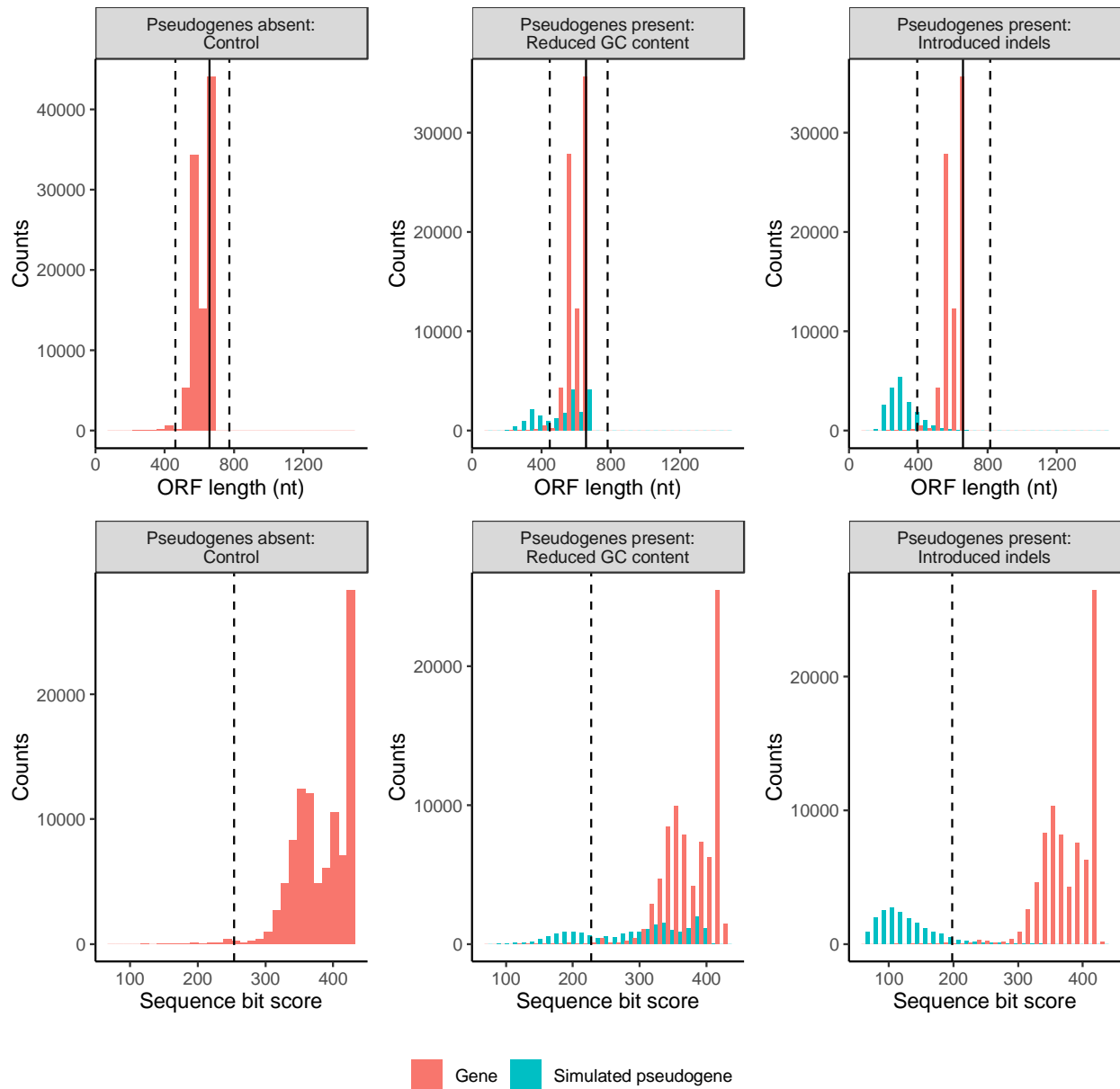
356

357 We used our observations from the simulated DNA barcode dataset with COI  
358 genes and pseudogenes from the same 10 species to guide the creation of a mock  
359 community comprised of 100,000 COI barcode sequences randomly sampled from  
360 BOLD where we could manipulate parameters in different ways. In our simulation study  
361 of full length COI sequences, we found that it was easier to filter out pseudogenes  
362 caused by increased indels (sensitivity 88-94%) rather than reduced GC content  
363 (sensitivity 27-31%) (Fig 3 and Table 2). As shown in Table 2, for full length COI  
364 barcode sequences, each pseudogene removal method performed with similar  
365 specificity (99-100%).

366

367 **Fig 3. In a simulated mock arthropod community, reducing the GC content or**  
368 **introducing indels in COI sequences reduces ORF lengths and sequence bit**  
369 **scores.** Each column shows the results from a particular simulation: a controlled  
370 community with pseudogenes absent, a community with pseudogenes that have a  
371 reduced GC content, and a community with pseudogenes where we have introduced  
372 indels. The top panel shows the length variation of sequences in the longest retained  
373 open reading frame. The solid vertical line indicates the length of a typical COI barcode  
374 at 658 bp. The two vertical dashed lines shows the boundaries for identifying ORFs  
375 with outlier lengths. The bottom panel shows the sequence bit score variation. The

376 vertical dashed line shows the boundary for identifying sequences with low outlier  
377 scores.



378

379

380

381 We also performed additional simulations by adjusting the length of the COI

382 barcodes from full length to half length (~ 329 bp) as this is similar to the length of COI

383 metabarcode sequences. As shown in Fig S13, it is more difficult to filter out short  
384 pseudogenes compared with full length COI barcodes. Table 2 shows that for half-  
385 length COI sequences, pseudogene removal sensitivity is better for pseudogenes  
386 generated by introducing indels (42-87%) rather than with pseudogenes where we  
387 reduced GC content (6-50%). Sensitivity is also generally higher when removing  
388 pseudogenes from the 5' end of the COI barcode region (15-87%) compared with the 3'  
389 end (6-61%). Pseudogene removal specificity is similar across pseudogene types and  
390 removal methods (99-100%).

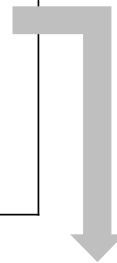
391         Since we don't really know how prevalent pseudogenes are in metabarcode  
392 datasets, we tested the effect of our pseudogene removal methods on a community  
393 where there are many pseudogenes (38% instead of 19% in previous analyses). Figure  
394 S14 shows that doubling the proportion of pseudogenes in the community greatly  
395 reduces the number of simulated pseudogenes removed with either method. As shown  
396 in Table 2, pseudogene removal sensitivity is poor (0-17%) but specificity is high using  
397 either removal method (99-100%). Next, we ran the opposite simulation where there  
398 are few pseudogenes in the community (9.5% instead of 19% in previous analyses).  
399 Figure S15 shows that reducing the number of pseudogenes in the community  
400 increases the number of simulated pseudogenes removed, especially when  
401 pseudogenes are caused by introducing indels. As Table 2 shows, the sensitivity of  
402 pseudogene removal is high when pseudogenes are created by introducing indels (95-  
403 98%), low when pseudogenes are created by reducing GC content (36-39%), and the  
404 specificity is high for any kind of simulated pseudogene or removal method (99-100%).

405           Because the ORFfinder + HMM profile analysis method for removing  
406 pseudogenes had the highest sensitivity for short COI metabarcodes when  
407 pseudogenes were simulated by introducing indels, we used this method to test our  
408 ability to remove pseudogenes with a real COI metabarcode dataset. Note that  
409 analyses were limited to only arthropod ESVs because most of the primer sets in the  
410 study were designed to specifically target this group in the original study (Table S2). As  
411 shown in Figure 4, the total number of arthropod ESVs was highest for the F230R  
412 amplicon (1,240) and least for the fwh1 amplicon (320). The greatest number of  
413 pseudogenes was detected and removed from the BR5 amplicon (19) and least for the  
414 ml-jg amplicon (1). Overall, the greatest percentage of pseudogenes out of all ESVs  
415 was detected from the BF2 amplicon (2.8%) and least for the ml-jg amplicon (0.1%).  
416 Because the F230R amplicon detected the greatest ESV richness, we used this  
417 amplicon to determine how existing bioinformatic processing steps affects pseudogene  
418 removal. Using the standard pipeline with ORFfinder + HMM profile analysis  
419 pseudogene removal, three F230R pseudogenes were removed from the dataset.  
420 Omitting the rare sequence removal step from the bioinformatic pipeline resulted in the  
421 largest number of pseudogenes detected, 34. Omitting the denoising step results in 1  
422 pseudogene detected. Omitting the chimera removal step results in 16 pseudogenes  
423 removed. This suggests to us that at least some apparent pseudogenes are probably  
424 already being removed during regular bioinformatic processing, especially during the  
425 rare sequence removal step as we would expect from the literature [53, 54, 57–59].  
426



427 **Fig 4. Removing rare sequences also removes apparent pseudogenes.** The  
 428 number of removed putative pseudogenes was calculated for each of the 5 amplicons  
 429 from a real freshwater COI metabarcode dataset. Note, that we only compared results  
 430 across Arthropoda ESVs. Using the standard bioinformatic pipeline, the F230R amplicon  
 431 recovered the greatest ESV richness (top box) so we used this as a test case for further  
 432 simulations (bottom box). To determine whether current bioinformatic processing steps  
 433 already help to remove apparent pseudogenes, we dropped one step at a time: removal  
 434 of rare sequences, removal of noisy sequences, and removal of chimeric sequences.  
 435

COI primer set	Total Arthropoda ESVs	No. pseudogenes removed	% Pseudogenes
BR5	813	19	2.3
F230R	1,240	3	0.24
ml-jg	1,039	1	0.1
BF1	906	13	1.4
BF2	467	13	2.8
fwh1	320	16	5



F230R Simulation	Remove rare	Remove noise	Remove chimeras	No. pseudogenes removed
Standard pipeline	✓	✓	✓	3
Skip rare removal	✗	✓	✓	34
Skip noise removal	✓	✗	✓	1
Skip chimera removal	✓	✓	✗	16

436

437

## 438 **Discussion**

439

440 Are all the COI sequences filtered out using ORFfinder + HMM profile analysis  
441 nuMTS? This method of sequence removal cannot distinguish between genuine  
442 pseudogenes and technical issues involving PCR or sequencing that causes indels,  
443 frameshifts, or the introduction of premature stop codons. It is possible that even after  
444 bioinformatic processing, artefactual sequences may be missed and subsequently be  
445 removed with these pseudogene removal methods. Although it is possible that genuine  
446 COI sequences could be removed using these methods, the specificity for pseudogenes  
447 is high (96-100%) and the number of COI gene sequences removed is very low in our  
448 simulated DNA barcode and metabarcode datasets.

449 There are also biological reasons why genuine mitochondrial sequences may be  
450 misclassified as pseudogenes. For example, in bivalves, male and female lineages of  
451 mitochondria may lead to fully functional gene copies with divergent sequences [15, 60,  
452 61]. Though this type of sequence could complicate for COI barcoding or phylogenetic  
453 analysis, this would not be filtered out by our methods because as functional COI genes  
454 they are not expected to have frame shifts or shorter length that our method uses to flag  
455 potential pseudogenes. There are also cases in the literature where as a cell ages  
456 oxidative stress damages DNA that is then repaired by enzymes with reduced activity  
457 [15, 62]. Unrepaired mutations including deletions, duplications, and point mutations  
458 can accumulate in aging cells. Since truncated mtDNA can be replicated faster than full  
459 length mtDNA, it is possible for partially deleted mtDNA to accumulate [63]. Similarly,  
460 damaged DNA caused by poor preservation could cause COI sequences with

461 frameshifts or premature stop codons to look like pseudogenes. It is quite likely that  
462 COI sequences with indels that lead to frameshifts and premature stop codons will be  
463 filtered out using the pseudogene removal methods we describe here whether the  
464 changes are technical or biological in nature.

465 How can pseudogenes be avoided? Indicators for the presence of pseudogenes  
466 include extra bands after PCR, sequence ambiguities when comparing both strands,  
467 frameshift mutations, premature stop codons, and unexpected phylogenetic position  
468 [18]. Strategies for avoiding pseudogenes in single specimens may include using  
469 muscle tissue for DNA extraction as it is naturally enriched with mtDNA, purifying  
470 mitochondria before DNA extraction, by amplifying long stretches of mtDNA with PCR,  
471 or targeting RNA using reverse transcription PCR [14, 18]. Even when working with  
472 environmental DNA samples, however, it can be possible to apply some of these  
473 techniques to avoid pseudogenes. For example, mitochondrial enrichment from  
474 homogenized tissues is possible and could be applied to freshwater benthic collections  
475 or insects collected from traps [64]. Additionally, long range PCR targeting  
476 mitochondrial DNA from water samples allowed for the construction of whole  
477 mitogenomes from fish [65]. Environmental RNA has also been used to detect  
478 microbes by targeting ribosomal RNA, this area has just begin to be explored using  
479 messenger RNA to target COI for metabarcoding [66–70]. For large scale studies,  
480 however, introducing additional steps such as mitochondrial purification or reverse  
481 transcription would be costly and time consuming.

482 Our results show that our ability to detect pseudogenes is hindered by short COI  
483 metabarcodes ~ 300 bp in length or if the abundance of sequenced pseudogenes is

484 very high. We show here that in a freshwater benthos COI metabarcode dataset, less  
485 than 3% of arthropod ESVs were removed as putative pseudogenes. It is quite possible  
486 that additional pseudogenes remain in the dataset, undetected by our pipeline. Our  
487 pseudogene removal methods cannot remove all pseudogenes, but remaining  
488 pseudogenes could still be useful for making higher level taxonomic assignments,  
489 though they may inflate richness at the species or haplotype level. Failure to remove  
490 low quality and artefactual sequences can result in inflated richness estimates in  
491 biodiversity studies, as has been shown for grasshoppers and crayfish [14].  
492 Pseudogenes are unlikely to affect community composition or beta diversity analyses if  
493 they are rare in the dataset as these analyses are less likely to be affected by the  
494 presence of rare sequences.

495 The use of phylogenetic based methods is common in COI barcoding studies, but  
496 the presence of pseudogenes could be a complication [14, 24, 26]. For example, a  
497 study of the great apes, showed that nuMTS are commonly sequenced in gorillas and  
498 complicate phylogenetic analyses [71]. It has also been suggested that pseudogenes  
499 are common in *Drosophila melanogaster* and in fish where they were once thought to  
500 be absent [72, 73]. The increasing use of COI metabarcodes for intraspecific analyses  
501 using ESVs could also be impacted by the presence of cryptic pseudogenes. The use  
502 of ORFfinder + HMM profile analysis, screening out hits with low outlier sequence bit  
503 scores, could be used as a first pass method for removing obvious pseudogenes. An  
504 automated method such as what we use in the SCVUC metabarcode pipelines in this  
505 study is more straight-forward to score compared with trying to identify pseudogenes  
506 from phylogenies by eye as branching patterns between genes and pseudogenes are

507 not always clear cut. To detect cryptic pseudogenes careful analysis of species level  
508 sequence alignments should still be carried out to check for sequences with low GC  
509 content, high dN/dS ratios, indels, and codon usage bias.

510 Hidden Markov model profile analysis is not a commonly used method to process  
511 COI metabarcodes but it is used for many other applications. For example, the ITSx  
512 extractor is a program used to process fungal ITS metabarcodes by identifying and  
513 removing the conserved gene regions adjacent to the internal transcribed spacer  
514 regions (ITS1 and ITS2) [74]. HMMs are already used in the Pfam database of protein  
515 families [75]. HMM analysis is also used to place 16S rRNA gene sequences in a  
516 reference phylogeny in PICRUST2 [76]. The HMM profile analysis approach would be  
517 suitable for identifying gene sequences from protein coding markers such as rbcL and  
518 matK (plants), such that poor hits could be filtered out as putative pseudogenes. A  
519 multi-marker metabarcode pipeline that processes paired-end Illumina reads that  
520 provides a pseudogene filtering step for protein coding markers is the MetaWorks  
521 snakemake pipeline that can be found at <https://github.com/terrimporter/MetaWorks> .  
522 Furthermore, though our current work has focused on arthropod sequences, taxon-  
523 specific HMM profiles could be developed for additional macroinvertebrate groups of  
524 interest for biomonitoring such as tubellaria, gastropoda, bivalvia, polychaeta,  
525 oligochaeta, and hirudinea to permit more refined HMM-profile analyses [46]. It would  
526 also be useful to develop HMM profiles for other commonly used protein coding markers  
527 such as rbcL and matK to facilitate nuMT removal from large plant sequence datasets.

528

529 **Conclusions**

530

531           We have shown that it is possible to screen out obvious pseudogenes using ORF  
532 length filtering alone or combined with HMM profile analysis for greater sensitivity when  
533 pseudogene sequences contain indels. Our pseudogenes removal approach was most  
534 effective on datasets of the full length COI barcode sequence region but is less effective  
535 for shorter sequences (~ 300 bp). This is especially relevant now that newer  
536 sequencing technologies such as LoopSeq (compatible with Illumina sequencing  
537 platforms, but currently only available for RNA genes) or HiFi circular consensus  
538 sequencing (PacBio) could one day be used for COI metabarcoding targeting the full  
539 length of the barcoding region facilitating pseudogene detection [12, 77–79]. It would  
540 also be helpful if COI barcode studies reported and deposited full length verified  
541 pseudogenes into public databases when possible. Having key words such as ‘nuclear  
542 copy of mitochondrial gene’ or ‘pseudogene’ in the description would be essential to  
543 quickly flag hits to such sequences. As the analysis of metabarcode sequences from  
544 protein-coding genes shifts towards the use of exact sequence variants, it is more  
545 important than ever to reduce noise by removing pseudogenes when possible to avoid  
546 inflated richness estimates or misleading phylogenetic results. The incorporation of  
547 pseudogene filtering steps into widely used pipelines such is needed.

548

549

550 **List of abbreviations**

551

552 BLAST - basic local alignment search tool

553 BOLD - Barcode of Life Data System

554 COI - cytochrome c oxidase subunit 1 gene

555 dN/dS - ratio of non-synonymous to synonymous substitutions

556 ESV - exact sequence variant

557 GC content - guanine-cytosine content

558 HMM - Hidden Markov Model

559 ITS - internal transcribed spacer region in the ribosomal RNA operon

560 K2P - Kimura 2-parameter model of nucleotide substitution

561 matK - maturase K gene

562 mtDNA - mitochondrial DNA

563 nuMT - nuclear encoded mitochondrial sequence

564 NCBI - National Center for Biotechnology Information

565 ORF - open reading frame

566 OTU - operational taxonomic unit

567 rbcL - ribulose bisphosphate carboxylate large chain gene

568

569

570

571 **Declarations:**

572 **Ethics approval and consent to participate** - Not applicable

573 **Consent for publication** - Not applicable

574 **Availability of data and materials** - All infiles and scripts used to parse data and

575 generate figures are available from GitHub at xxx. The SCVUC COI metabarcode

576 pipeline used in this study is also available on GitHub from

577 [https://github.com/Hajibabaei-Lab/SCVUC\\_COI\\_metabarcode\\_pipeline](https://github.com/Hajibabaei-Lab/SCVUC_COI_metabarcode_pipeline).

578 **Competing interests** - None

579 **Funding** – This study is funded by the Government of Canada through Genome

580 Canada and Ontario Genomics.

581 **Authors' contributions** – MH and TP conceived of the idea. TP conducted the

582 analyses and wrote the manuscript. MH provided critical input into analysis methods

583 and the manuscript. MH provided funding and computational resources. Both authors

584 edited, read, and approved the final manuscript.

585 **Acknowledgements** We would like to thank the Hajibabaei group for their support and

586 helpful discussions.

587

588

589



## 590 **References**

- 591 1. Hebert PDN, Cywinska A, Ball SL, deWaard JR. Biological identifications through DNA  
592 barcodes. *Proceedings of the Royal Society B: Biological Sciences*. 2003;270:313–21.
- 593 2. Ratnasingham S, Hebert PD. BOLD: The Barcode of Life Data System ([http://www.](http://www.barcodinglife.org)  
594 [barcodinglife.org](http://www.barcodinglife.org)). *Molecular ecology notes*. 2007;7:355–64.
- 595 3. Porter TM, Hajibabaei M. Over 2.5 million COI sequences in GenBank and growing. *PLoS ONE*.  
596 2018;13:e0200177.
- 597 4. Bruns TD, White TJ, Taylor JW. Fungal Molecular Systematics. *Annual Review of Ecology and*  
598 *Systematics*. 1991;22:525–64.
- 599 5. Stackebrandt E, Goebel BM. Taxonomic Note: A Place for DNA-DNA Reassociation and 16S  
600 rRNA Sequence Analysis in the Present Species Definition in Bacteriology. *International Journal*  
601 *of Systematic and Evolutionary Microbiology*. 1994;44:846–9.
- 602 6. Schoch CL, Seifert KA, Huhndorf S, Robert V, Spouge JL, Levesque CA, et al. Nuclear ribosomal  
603 internal transcribed spacer (ITS) region as a universal DNA barcode marker for Fungi.  
604 *Proceedings of the National Academy of Sciences*. 2012;109:6241–6.
- 605 7. Elbrecht V, Vamos EE, Steinke D, Leese F. Estimating intraspecific genetic diversity from  
606 community DNA metabarcoding data. *PeerJ*. 2018;6:e4644.
- 607 8. Porter TM, Hajibabaei M. Putting COI Metabarcoding in Context: The Utility of Exact  
608 Sequence Variants (ESVs) in Biodiversity Analysis. *Front Ecol Evol*. 2020;8:248.
- 609 9. Antich A, Palacin C, Wangenstein OS, Turon X. To denoise or to cluster? That is not the  
610 question. Optimizing pipelines for COI metabarcoding and metaphylogeography. preprint.  
611 *Genetics*; 2021. doi:10.1101/2021.01.08.425760.
- 612 10. Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace operational  
613 taxonomic units in marker-gene data analysis. *The ISME Journal*. 2017;11:2639–43.
- 614 11. Buchner D, Leese F. BOLDigger – a Python package to identify and organise sequences with  
615 the Barcode of Life Data systems. *MBMG*. 2020;4:e53535.
- 616 12. Nugent CM, Elliott TA, Ratnasingham S, Hebert PDN, Adamowicz SJ. debar, a sequence-by-  
617 sequence denoiser for COI-5P DNA barcode data. preprint. *Bioinformatics*; 2021.  
618 doi:10.1101/2021.01.04.425285.
- 619 13. Nugent CM, Elliott TA, Ratnasingham S, Adamowicz SJ. coil: an R package for cytochrome C  
620 oxidase I (COI) DNA barcode data cleaning, translation, and error evaluation. *bioRxiv*. 2019;:35.

- 621 14. Song H, Buhay JE, Whiting MF, Crandall KA. Many species in one: DNA barcoding  
622 overestimates the number of species when nuclear mitochondrial pseudogenes are  
623 coamplified. *PNAS*. 2008;105:13486–91.
- 624 15. Schizas N. Misconceptions regarding nuclear mitochondrial pseudogenes (Numts) may  
625 obscure detection of mitochondrial evolutionary novelties. *Aquatic Biology*. 2012;17:91–6.
- 626 16. Leite LAR. Mitochondrial pseudogenes in insect DNA barcoding: differing points of view on  
627 the same issue. *Biota Neotrop*. 2012;12:301–8.
- 628 17. Andújar C, Creedy TJ, Arribas P, López H, Salces-Castellano A, Pérez-Delgado A, et al. NUMT  
629 dumping: validated removal of nuclear pseudogenes from mitochondrial metabarcoding data.  
630 preprint. *Evolutionary Biology*; 2020. doi:10.1101/2020.06.17.157347.
- 631 18. Bensasson D. Mitochondrial pseudogenes: evolution’s misplaced witnesses. *Trends in  
632 Ecology & Evolution*. 2001;16:314–21.
- 633 19. Hazkani-Covo E, Zeller RM, Martin W. Molecular Poltergeists: Mitochondrial DNA Copies  
634 (numts) in Sequenced Nuclear Genomes. *PLoS Genet*. 2010;6:e1000834.
- 635 20. Adams KL, Palmer JD. Evolution of mitochondrial gene content: gene loss and transfer to the  
636 nucleus. *Molecular Phylogenetics and Evolution*. 2003;29:380–95.
- 637 21. Bertheau C, Schuler H, Krumböck S, Arthofer W, Stauffer C. Hit or miss in phylogeographic  
638 analyses: the case of the cryptic NUMTs. *Molecular Ecology Resources*. 2011;11:1056–9.
- 639 22. Zhang D-X, Hewitt GM. Nuclear integrations: challenges for mitochondrial DNA markers.  
640 *Trends in Ecology & Evolution*. 1996;11:247–51.
- 641 23. Martins J, Solomon SE, Mikheyev AS, Mueller UG, Ortiz A, Bacci M. Nuclear mitochondrial-  
642 like sequences in ants: evidence from *Atta cephalotes* (Formicidae: Attini): Numts in *A.*  
643 *cephalotes* ants. *Insect Molecular Biology*. 2007;16:777–84.
- 644 24. Williams ST, Knowlton N. Mitochondrial Pseudogenes Are Pervasive and Often Insidious in  
645 the Snapping Shrimp Genus *Alpheus*. *Molecular Biology and Evolution*. 2001;18:1484–93.
- 646 25. Moulton MJ, Song H, Whiting MF. Assessing the effects of primer specificity on eliminating  
647 numt coamplification in DNA barcoding: a case study from Orthoptera (Arthropoda: Insecta):  
648 DNA BARCODING. *Molecular Ecology Resources*. 2010;10:615–27.
- 649 26. Buhay JE. “COI-like” Sequences Are Becoming Problematic in Molecular Systematic and DNA  
650 Barcoding Studies. *Journal of Crustacean Biology*. 2009;29:96–110.
- 651 27. Pentinsaari M, Salmela H, Mutanen M, Roslin T. Molecular evolution of a widely-adopted  
652 taxonomic marker (COI) across the animal tree of life. *Scientific Reports*. 2016;6.  
653 doi:10.1038/srep35275.

- 654 28. Sayers EW. Ebot. <http://www.ncbi.nlm.nih.gov/Class/PowerTools/eutils/course.html>.
- 655 29. Charif D, Lobry J. SeqinR 1.0-2: a contributed package to the R project for statistical  
656 computing devoted to biological sequences retrieval and analysis. In: Structural approaches to  
657 sequence evolution: Molecules, networks, populations. New York: Springer Verlag; 2007. p.  
658 207–32.
- 659 30. Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:  
660 Improvements in Performance and Usability. *Molecular Biology and Evolution*. 2013;30:772–80.
- 661 31. Rice P, Longden I, Bleasby A. EMBOSS: The European Molecular Biology Open Software  
662 Suite. *Trends in Genetics*. 2000;16:276–7.
- 663 32. Wickham H. ggplot2: Elegant Graphics for Data Analysis. New York: Springer-Verlag; 2009.  
664 <http://ggplot2.org>.
- 665 33. RStudio Team. RStudio: Integrated Development Environment for R. 2016.  
666 <http://www.rstudio.com/>.
- 667 34. R Core Team. R: A Language and Environment for Statistical Computing. 2017.  
668 <https://www.R-project.org/>.
- 669 35. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and  
670 PSI-BLAST: a new generation of protein database search programs. *Nucleic acids research*.  
671 1997;25:17.
- 672 36. Kimura M. A simple method for estimating evolutionary rates of base substitutions through  
673 comparative studies of nucleotide sequences. *J Mol Evol*. 1980;16:111–20.
- 674 37. Felsenstein J. PHYLIP - Phylogeny Inference Package (Version 3.2). *Cladistics*. 1989;5:164–6.
- 675 38. Stöver BC, Müller KF. TreeGraph 2: Combining and visualizing evidence from different  
676 phylogenetic analyses. *BMC Bioinformatics*. 2010;11:7.
- 677 39. Hajibabaei M, Porter TM, Wright M, Rudar J. COI metabarcoding primer choice affects  
678 richness and recovery of indicator taxa in freshwater systems. *PLoS ONE*. 2019;14:e0220953.
- 679 40. Hajibabaei M, Spall JL, Shokralla S, van Konynenburg S. Assessing biodiversity of a  
680 freshwater benthic macroinvertebrate community through non-destructive environmental  
681 barcoding of DNA from preservative ethanol. *BMC Ecology*. 2012;12:28.
- 682 41. Gibson J, Shokralla S, Porter TM, King I, Konynenburg S van, Janzen DH, et al. Simultaneous  
683 assessment of the macrobiome and microbiome in a bulk sample of tropical arthropods  
684 through DNA metasytematics. *PNAS*. 2014;111:8007–12.

- 685 42. Folmer O, Black M, Hoeh W, Lutz R, Vrijenhoek R. DNA primers for amplification of  
686 mitochondrial cytochrome c oxidase subunit I from diverse metazoan invertebrates. *Molecular*  
687 *marine biology and biotechnology*. 1994;3:294–9.
- 688 43. Gibson J, Shokralla S, Curry C, Baird DJ, Monk WA, King I, et al. Large-Scale Biomonitoring of  
689 Remote and Threatened Ecosystems via High-Throughput Sequencing. *PLOS ONE*.  
690 2015;10:e0138432.
- 691 44. Leray M, Yang JY, Meyer CP, Mills SC, Agudelo N, Ranwez V, et al. A new versatile primer set  
692 targeting a short fragment of the mitochondrial COI region for metabarcoding metazoan  
693 diversity: application for characterizing coral reef fish gut contents. *Frontiers in Zoology*.  
694 2013;10:34.
- 695 45. Geller J, Meyer C, Parker M, Hawk H. Redesign of PCR primers for mitochondrial cytochrome  
696 c oxidase subunit I for marine invertebrates and application in all-taxa biotic surveys. *Mol Ecol*  
697 *Resour*. 2013;13:851–61.
- 698 46. Elbrecht V, Leese F. Validation and Development of COI Metabarcoding Primers for  
699 Freshwater Macroinvertebrate Bioassessment. *Frontiers in Environmental Science*. 2017;5:11.
- 700 47. Vamos E, Elbrecht V, Leese F. Short COI markers for freshwater macroinvertebrate  
701 metabarcoding. *Metabarcoding and Metagenomics*. 2017;1:e14625.
- 702 48. Anaconda. Anaconda Software Distribution. 2016. <https://anaconda.com>.
- 703 49. Koster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine.  
704 *Bioinformatics*. 2012;28:2520–2.
- 705 50. St. John J. SeqPrep. 2016. <https://github.com/jstjohn/SeqPrep/releases>.
- 706 51. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
707 *EMBnet journal*. 2011;17:pp-10.
- 708 52. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool for  
709 metagenomics. *PeerJ*. 2016;4:e2584.
- 710 53. Edgar RC. UNOISE2: improved error-correction for Illumina 16S and ITS amplicon  
711 sequencing. *bioRxiv*. 2016. doi:10.1101/081257.
- 712 54. Edgar R. UCHIME2: improved chimera prediction for amplicon sequencing. *bioRxiv*.  
713 2016;:074252.
- 714 55. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier for Rapid Assignment of  
715 rRNA Sequences into the New Bacterial Taxonomy. *Applied and Environmental Microbiology*.  
716 2007;73:5261–7.

- 717 56. Porter TM, Hajibabaei M. Automated high throughput animal CO1 metabarcode  
718 classification. *Scientific Reports*. 2018;8:4226.
- 719 57. Reeder J, Knight R. The 'rare biosphere': a reality check. *nature methods*. 2009;6:636–7.
- 720 58. Tedersoo L, Nilsson RH, Abarenkov K, Jairus T, Sadam A, Saar I, et al. 454 Pyrosequencing  
721 and Sanger sequencing of tropical mycorrhizal fungi provide similar results but reveal  
722 substantial methodological biases. *New Phytologist*. 2010;188:291–301.
- 723 59. Leray M, Knowlton N. Random sampling causes the low reproducibility of rare eukaryotic  
724 OTUs in Illumina COI metabarcoding. *PeerJ*. 2017;5:e3006.
- 725 60. Zouros E, Oberhauser Ball A, Saavedra C, Freeman KR. An unusual type of mitochondrial  
726 DNA inheritance in the blue mussel *Mytilus*. *Proceedings of the National Academy of Sciences*.  
727 1994;91:7463–7.
- 728 61. Stewart DT, Saavedra C, Stanwood RR, Ball AO, Zouros E. Male and female mitochondrial  
729 DNA lineages in the blue mussel (*Mytilus edulis*) species group. *Molecular Biology and*  
730 *Evolution*. 1995;12:735–47.
- 731 62. Druzhyna NM, Wilson GL, LeDoux SP. Mitochondrial DNA repair in aging and disease.  
732 *Mechanisms of Ageing and Development*. 2008;129:383–90.
- 733 63. Diaz F, Bayona-Bafaluy MP, Rana M, Mora M, Hao H, Moraes CT. Human mitochondrial DNA  
734 with large deletions repopulates organelles faster than full-length genomes under relaxed copy  
735 number control. *Nucleic Acids Research*. 2002;30:4626–33.
- 736 64. Zhou X, Li Y, Liu S, Yang Q, Su X, Zhou L, et al. Ultra-deep sequencing enables high-fidelity  
737 recovery of biodiversity for bulk arthropod samples without PCR amplification. *GigaSci*.  
738 2013;2:4.
- 739 65. Deiner K, Bik HM, Mächler E, Seymour M, Lacoursière-Roussel A, Altermatt F, et al.  
740 Environmental DNA metabarcoding: transforming how we survey animal and plant  
741 communities. *Molecular Ecology*. 2017;26:5872–95.
- 742 66. Tsuru K, Ikeda S, Hirohara T, Shimada Y, Minamoto T, Yamanaka H. Messenger RNA typing of  
743 environmental RNA (eRNA): A case study on zebrafish tank water with perspectives for the  
744 future development of eRNA analysis on aquatic vertebrates. *Environmental DNA*. 2021;3:14–  
745 21.
- 746 67. Laroche O, Wood SA, Tremblay LA, Lear G, Ellis JI, Pochon X. Metabarcoding monitoring  
747 analysis: the pros and cons of using co-extracted environmental DNA and RNA data to assess  
748 offshore oil production impacts on benthic communities. *PeerJ*. 2017;5:e3347.

- 749 68. Pochon X, Zaiko A, Fletcher LM, Laroche O, Wood SA. Wanted dead or alive? Using  
750 metabarcoding of environmental DNA and RNA to distinguish living assemblages for biosecurity  
751 applications. *PLoS ONE*. 2017;12:e0187636.
- 752 69. Harris M. Assessing the Persistence of Environmental DNA and Environmental RNA for  
753 Zooplankton Biodiversity Monitoring by Metabarcoding. McGill University; 2019.  
754 [https://search.proquest.com/openview/547572df2ecd232f9071d0fa45507688/1?cbl=44156&](https://search.proquest.com/openview/547572df2ecd232f9071d0fa45507688/1?cbl=44156&loginDisplay=true&pq-origsite=gscholar)  
755 [oginDisplay=true&pq-origsite=gscholar](https://search.proquest.com/openview/547572df2ecd232f9071d0fa45507688/1?cbl=44156&loginDisplay=true&pq-origsite=gscholar).
- 756 70. Cristescu ME. Can Environmental RNA Revolutionize Biodiversity Science? *Trends in Ecology*  
757 *& Evolution*. 2019;34:694–7.
- 758 71. Thalmann O, Hebler J, Poinar HN, Pääbo S, Vigilant L. Unreliable mtDNA data due to nuclear  
759 insertions: a cautionary tale from analysis of humans and other great apes: NUMTS IN APES.  
760 *Molecular Ecology*. 2004;13:321–35.
- 761 72. Harrison PM. Identification of pseudogenes in the *Drosophila melanogaster* genome.  
762 *Nucleic Acids Research*. 2003;31:1033–7.
- 763 73. Antunes A, Ramos MJ. Discovery of a large number of previously unrecognized  
764 mitochondrial pseudogenes in fish genomes. *Genomics*. 2005;86:708–17.
- 765 74. Bengtsson-Palme J, Ryberg M, Hartmann M, Branco S, Wang Z, Godhe A, et al. Improved  
766 software detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi and  
767 other eukaryotes for analysis of environmental sequencing data. *Methods in Ecology and*  
768 *Evolution*. 2013;4:914–9.
- 769 75. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al. Pfam: the protein  
770 families database. *Nucl Acids Res*. 2014;42:D222–30.
- 771 76. Douglas GM, Maffei VJ, Zaneveld J, Yurgel SN, Brown JR, Taylor CM, et al. PICRUSt2 for  
772 prediction of metagenome functions. *Nature Biotechnology*. 2020;38:685–8.
- 773 77. Callahan BJ, Grinevich D, Thakur S, Balamotis MA, Yehezkel TB. Ultra-accurate Microbial  
774 Amplicon Sequencing Directly from Complex Samples with Synthetic Long Reads. preprint.  
775 *Microbiology*; 2020. doi:10.1101/2020.07.07.192286.
- 776 78. Tedersoo L, Tooming-Klunderud A, Anslan S. PacBio metabarcoding of Fungi and other  
777 eukaryotes: errors, biases and perspectives. *New Phytol*. 2018;217:1370–85.
- 778 79. Wurzbacher C, Larsson E, Bengtsson-Palme J, Van den Wyngaert S, Svantesson S,  
779 Kristiansson E, et al. Introducing ribosomal tandem repeat barcoding for fungi. 2018.  
780 doi:10.1101/310540.
- 781

782 **Supplementary Material**

783

784 **Table S1. Description of the datasets analyzed in Part A and Part B.**

Experiment	Dataset	Proportion of dataset comprised of pseudogenes (%)	Average gene length (bp)	Average pseudogene length (bp)	Gene GC content (%)	Pseudogene GC content (%)
Part A	Simulated DNA barcode dataset	19	659.6	508.1	32.0	30.8
Part B	Control mock community with 100,000 randomly sampled sequences from BOLD	0	615	NA	31	NA
Part B	Mock community with decreased GC content	19	615	615	31	29
Part B	Mock community with increased indels	19	615	607	31	31
Part B	Control mock community with half-length sequences	0	307** - 308*	NA	30*-32**	NA
Part B	Mock community with half-length sequences and decreased GC content	19	307** - 308*	308	30*-32**	28-29
Part B	Mock community with half-length sequences and increased indels	19	307** - 308*	304	30*-32**	31-32
Part B	Control mock community with twice 100,000 randomly sampled BOLD sequences	0	622	NA	31	NA
Part B	Mock community with twice the	38	622	622	31	28

	number of pseudogenes with decreased GC content					
Part B	Mock community with twice the number of pseudogenes with increased indels	38	622	614	31	32
Part B	Control mock community of 100,000 randomly sampled sequences from BOLD	0	622	NA	31	NA
Part B	Mock community with halved number of pseudogenes with decreased GC content	9.5	622	623	31	28
Part B	Mock community with halved number of pseudogenes with increased indels	9.5	622	615	31	32

785 \* 5' fragment

786 \*\* 3' fragment

787

788

789



790 **Table S2. Primers used in the freshwater benthos COI metabarcode dataset used**  
 791 **in Part C (Hajibabaei et al., 2019 PLoS ONE).**

792

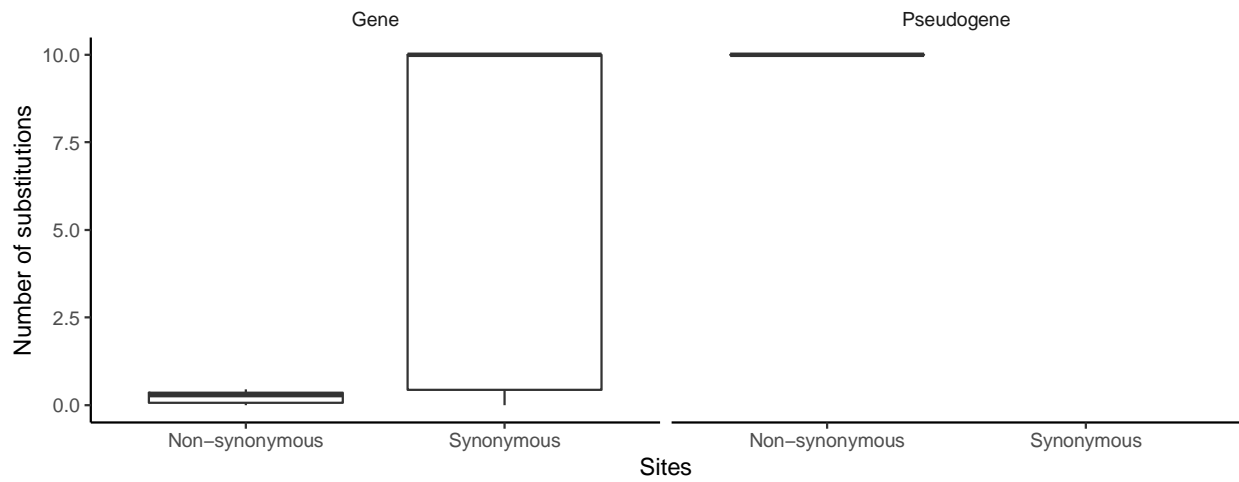
Amplicon	Primer	Target	Primer sequence (5'-3')	Reference
BR5	B	Freshwater benthic macroinvertebrates	CCIGAYATRGCITTYCCICG	Hajibabaei et al., 2012
	ArR5	Tropical arthropods	GTRATIGCICCIACIARIACIG G	Gibson et al. 2014*
F230R	LCO1490	Metazoan macroinvertebrates	GGTCAACAAATCATAAAGAT ATTGG	Folmer et al., 1994
	230_R	Arthropods	CTTATRTTRTTTATICGIGGR AAIGC	Gibson et al., 2015
ml-jg	mlCOLintF	Metazoa	GGWACWGGWTGAACWGT WTAYCCYCC	Leray et al., 2013
	lgHCO2198	Marine invertebrates	TAIACYTCIGGRTGICCRAAR AAYCA	Geller et al., 2013
BF1	BF1	Freshwater macroinvertebrates	ACWGGWTGRACWGTNTAY CC	Elbrecht and Leese, 2017
	BR2	Freshwater macroinvertebrates	TCDGGRTGNCCRAARAAYC A	Elbrecht and Leese, 2017
BF2	BF2	Freshwater macroinvertebrates	GCHCCHGAYATRGCHTTYC C	Elbrecht and Leese, 2017
	BR2	Freshwater macroinvertebrates	TCDGGRTGNCCRAARAAYC A	Elbrecht and Leese, 2017
fwh1	fwhF1	Freshwater macroinvertebrates	YTCHACWAAYCAYAARGAY ATYGG	Vamos et al., 2017
	fwhR1	Freshwater macroinvertebrates	ARTCARTTWCCRAAHCHC C	Vamos et al., 2017

793 \* This primer sequence was published based on its alignment to the plus strand but is  
 794 shown here in the 5'-3' orientation

795

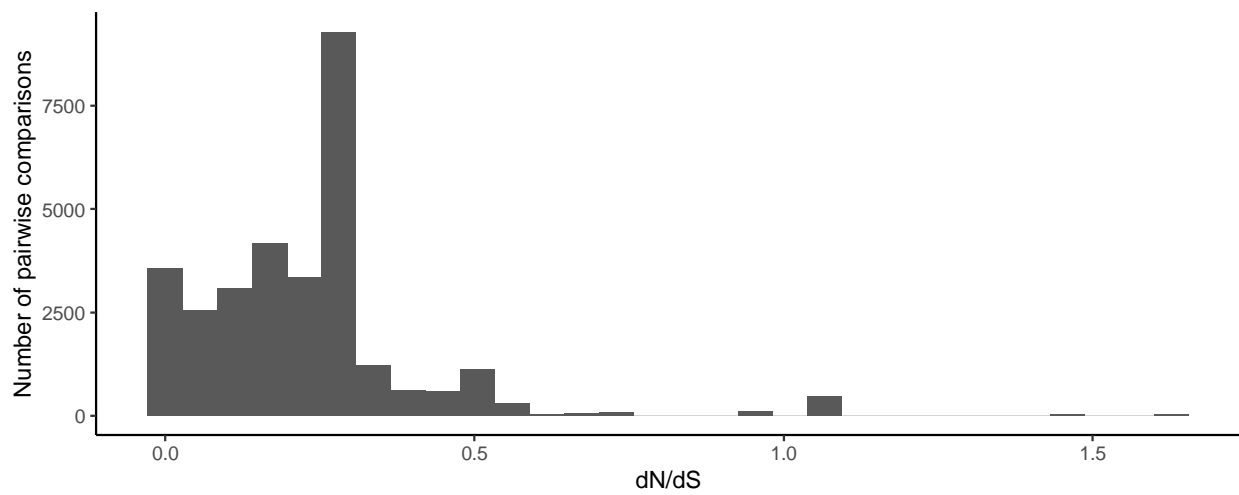
796 **Fig S1. COI gene sequences accumulate substitutions in synonymous sites.** For  
797 10 species with annotated COI genes and pseudogenes, we did a pairwise comparison  
798 of nucleotide substitutions in non-synonymous and synonymous sites: a) COI barcode  
799 sequences tend to accumulate substitutions in synonymous sites. In contrast, COI  
800 pseudogenes tend to accumulate substitutions in non-synonymous sites. After filtering  
801 out pairwise comparisons between species with  $< 0.01$  substitutions in synonymous  
802 sites (sequences too similar to yield a reliable dN/dS estimate) or  $> 2$  substitutions in  
803 synonymous sites (sequences that have accumulated too many substitutions to yield a  
804 reliable dN/dS estimate), it was only possible to analyze dN/dS ratios for COI barcode  
805 sequences. b) Most pairwise comparisons of COI gene sequences resulted in dN/dS  
806 ratios  $< 1$  consistent with purifying selection pressure and the conservation of a protein  
807 sequence.

808 a)



809

810 b)



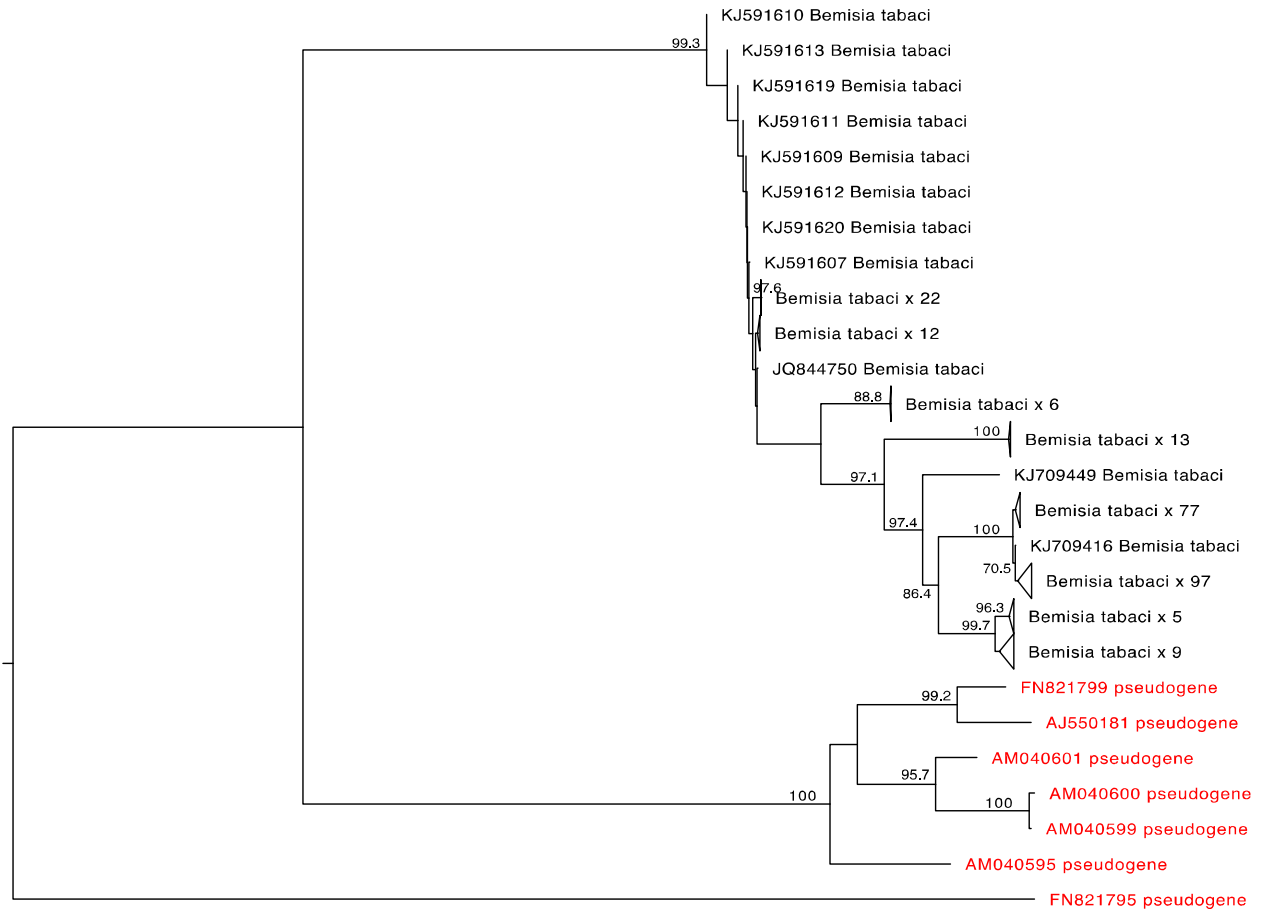
811

812

813

814

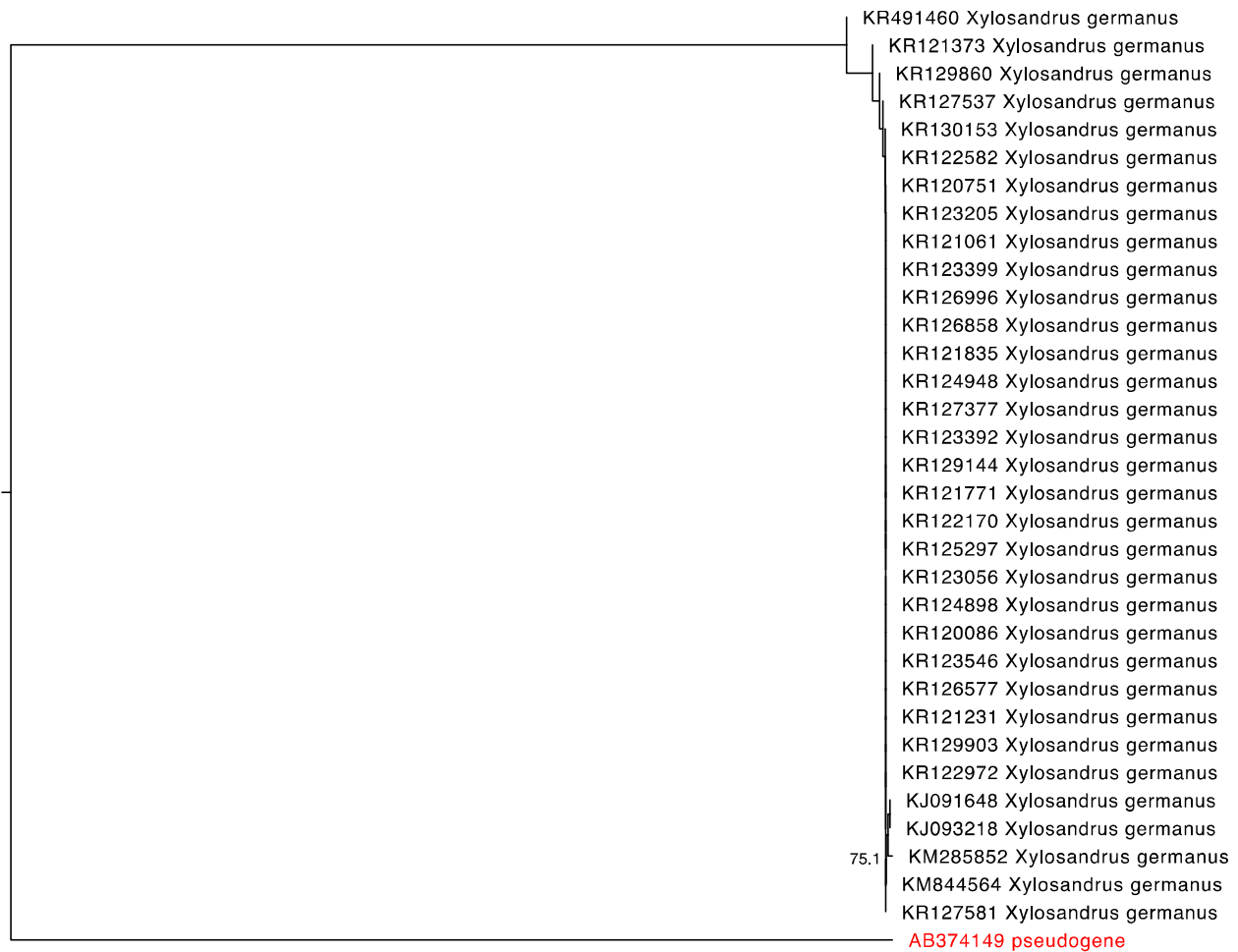
815 **Fig S2. *Bemisia tabaci* COI pseudogenes cluster together on long branches. A**  
816 mid-point rooted neighbor joining phylogram using the Kimura 2-parameter model of  
817 nucleotide substitution included gene and known pseudogene sequences. Sequences  
818 annotated in GenBank as a nuclear copy of a mitochondrial gene are shown in red.  
819 Nodes with greater than 70% bootstrap support are labelled.



820

821

822 **Fig S3. A single *Xylosandrus germanus* COI pseudogene sequence is found on a**  
823 **long branch.** A mid-point rooted neighbor joining phylogram using the Kimura 2-  
824 parameter model of nucleotide substitution included COI gene sequences as well as a  
825 sequence annotated in GenBank as a nuclear copy of a mitochondrial gene (red).  
826 Nodes with greater than 70% bootstrap support are labelled.  
827



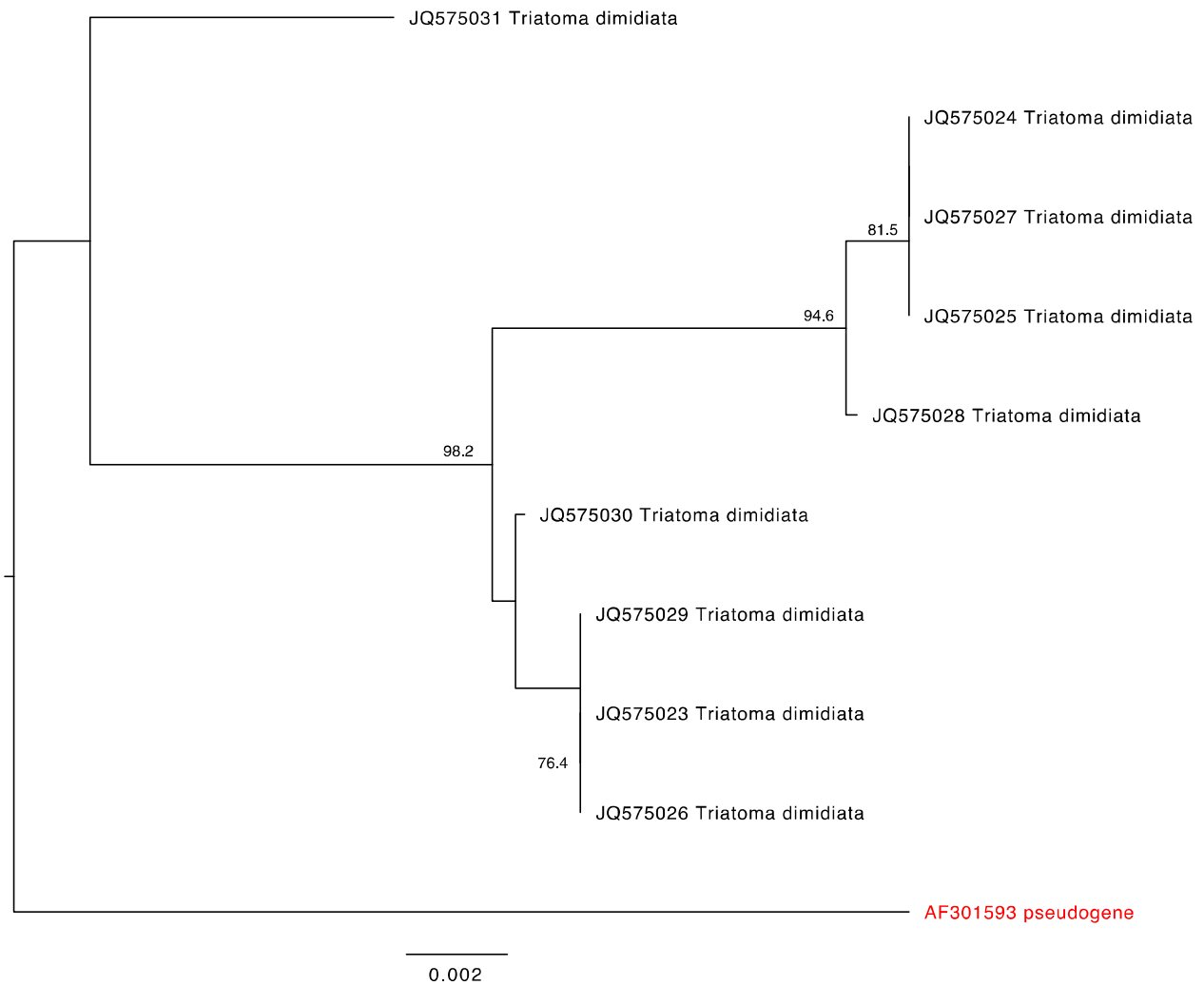
828

829

830

831

832 **Fig S4. A single *Triatoma dimidiata* COI pseudogene sequence is found on a long**  
833 **branch.** A mid-point rooted neighbor joining phylogram using the Kimura 2-parameter  
834 model of nucleotide substitution included COI gene sequences as well as a sequence  
835 annotated in GenBank as a nuclear copy of a mitochondrial gene (red). Nodes with  
836 greater than 70% bootstrap support are labelled.  
837



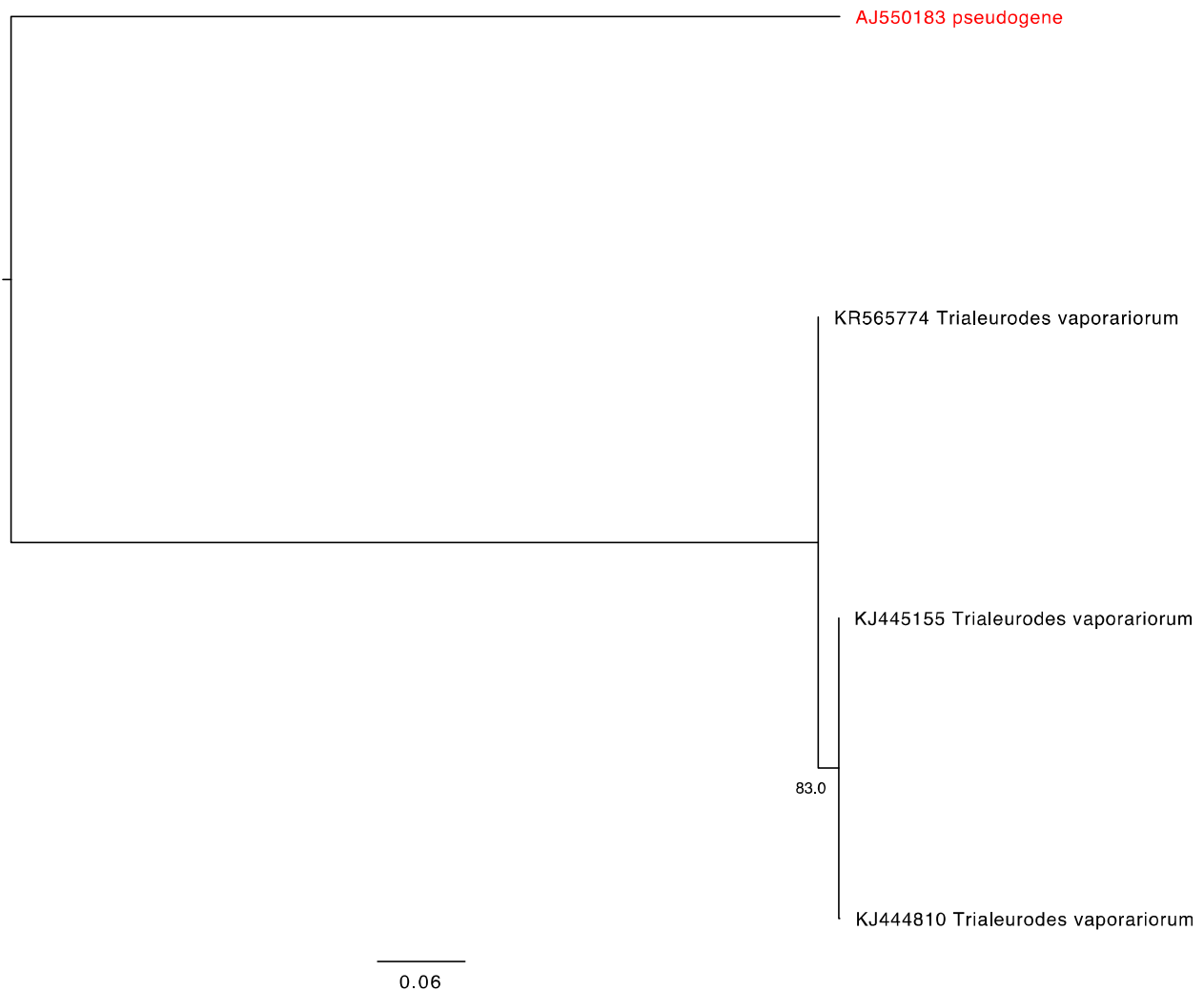
838

839

840

841

842 **Fig S5. A single *Trialeurodes vaporariorum* COI pseudogene sequence is found**  
843 **on a long branch.** A mid-point rooted neighbor joining phylogram using the Kimura 2-  
844 parameter model of nucleotide substitution included COI gene sequences as well as a  
845 sequence annotated in GenBank as a nuclear copy of a mitochondrial gene (red).  
846 Nodes with greater than 70% bootstrap support are labelled.  
847



848

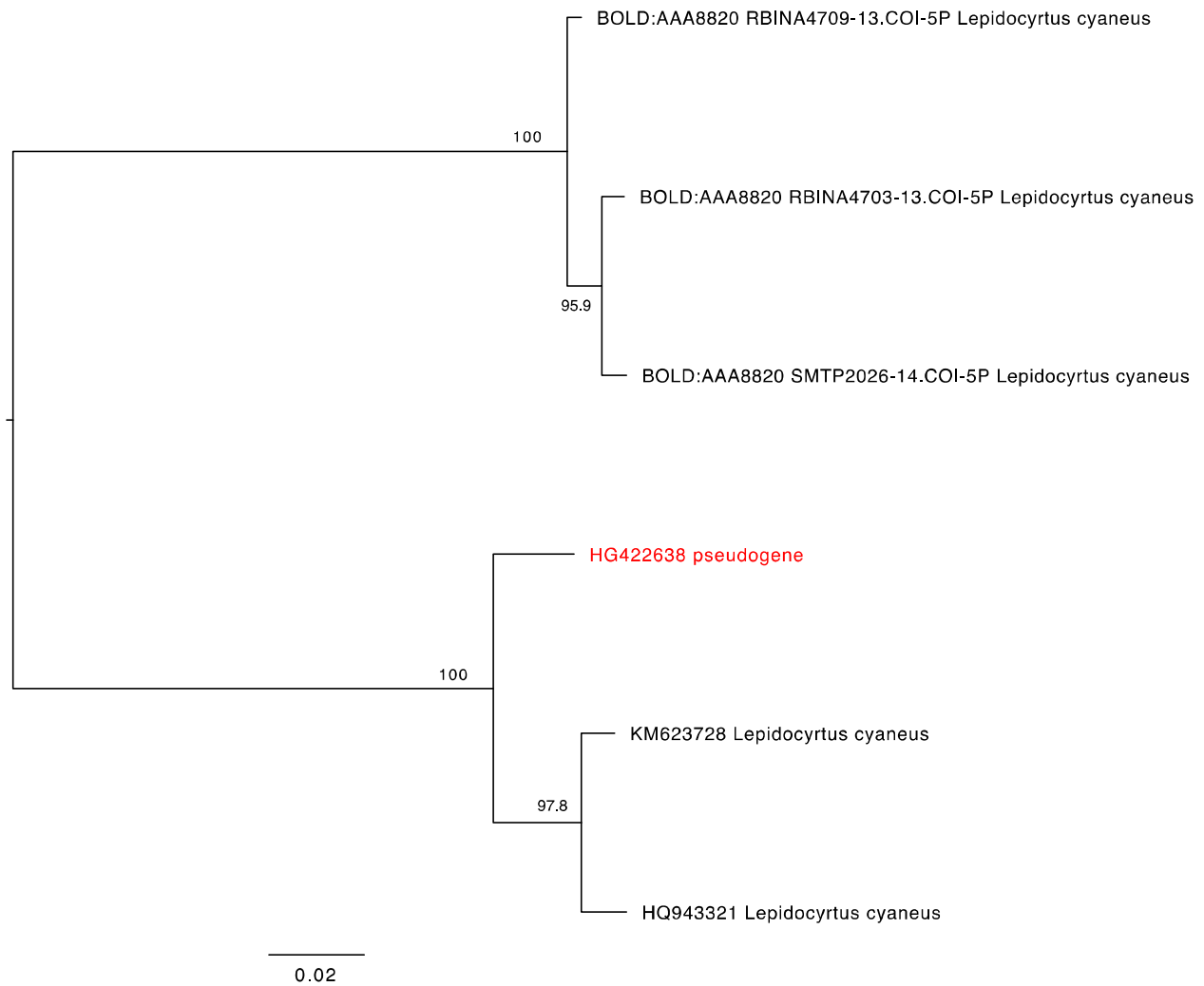
849

850 **Fig S6. *Melissotarsus insularis* COI gene and annotated pseudogene sequences**  
851 **are often found in intermixed clusters.** A mid-point rooted neighbor joining  
852 phylogram using the Kimura 2-parameter model of nucleotide substitution included COI  
853 gene sequences as well as sequences annotated in GenBank as a nuclear copy of a  
854 mitochondrial gene (red). Nodes with greater than 70% bootstrap support are labelled.  
855 Clusters of nearly identical sequences were collapsed.  
856





858 **Fig S7. A single *Lepidocyrtus cyaneus* COI pseudogene sequence clusters with**  
859 **other gene sequences.** A mid-point rooted neighbor joining phylogram using the  
860 Kimura 2-parameter model of nucleotide substitution included COI gene sequences as  
861 well as a sequence annotated in GenBank as a nuclear copy of a mitochondrial gene  
862 (red). Nodes with greater than 70% bootstrap support are labelled.  
863

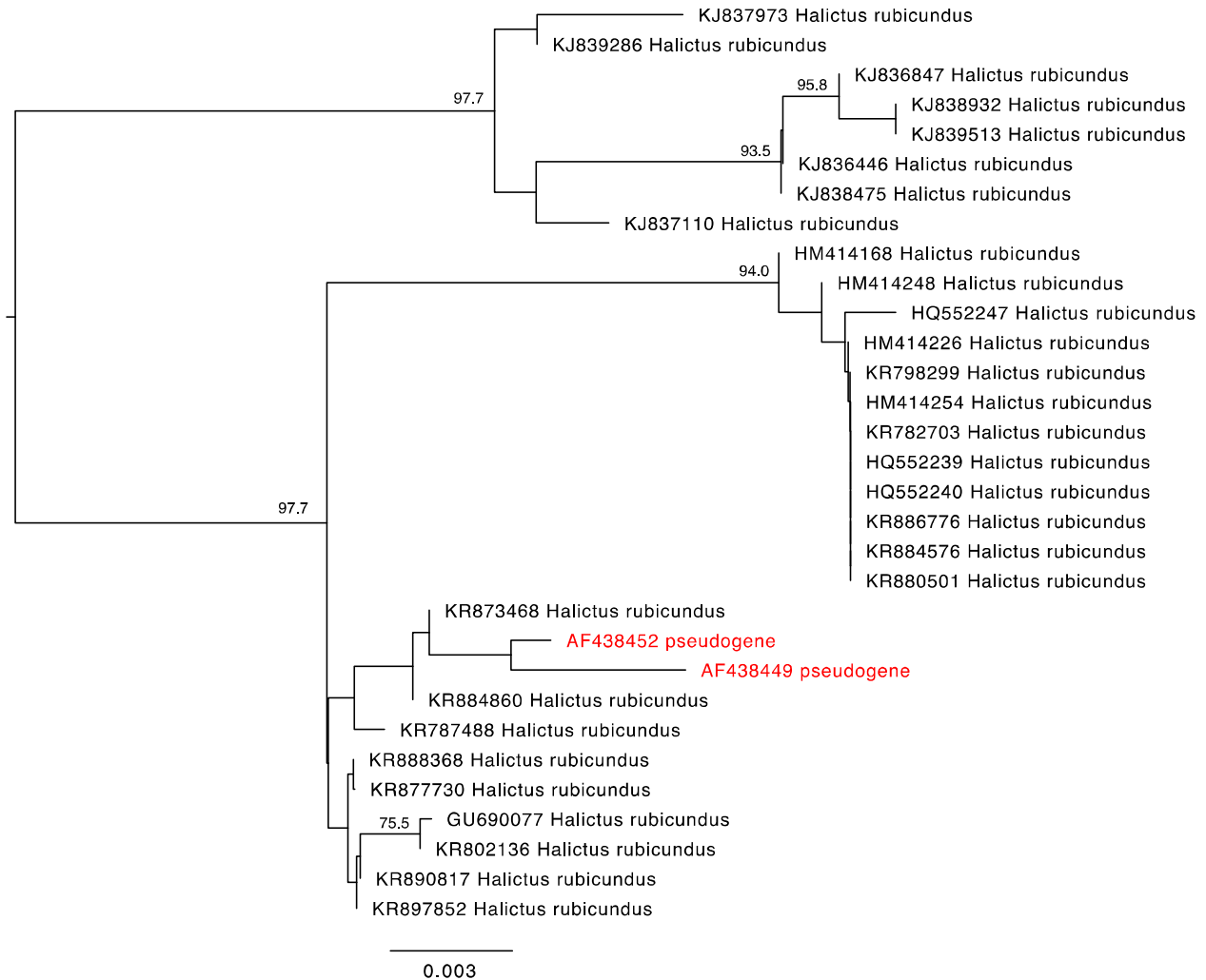


864

865

866

867 **Fig S8. Two *Halictus rubicundus* COI pseudogene sequences cluster together**  
868 **near other gene sequences.** A mid-point rooted neighbor joining phylogram using the  
869 Kimura 2-parameter model of nucleotide substitution included COI gene sequences as  
870 well as two sequences annotated in GenBank as a nuclear copy of a mitochondrial  
871 gene (red). Nodes with greater than 70% bootstrap support are labelled.  
872

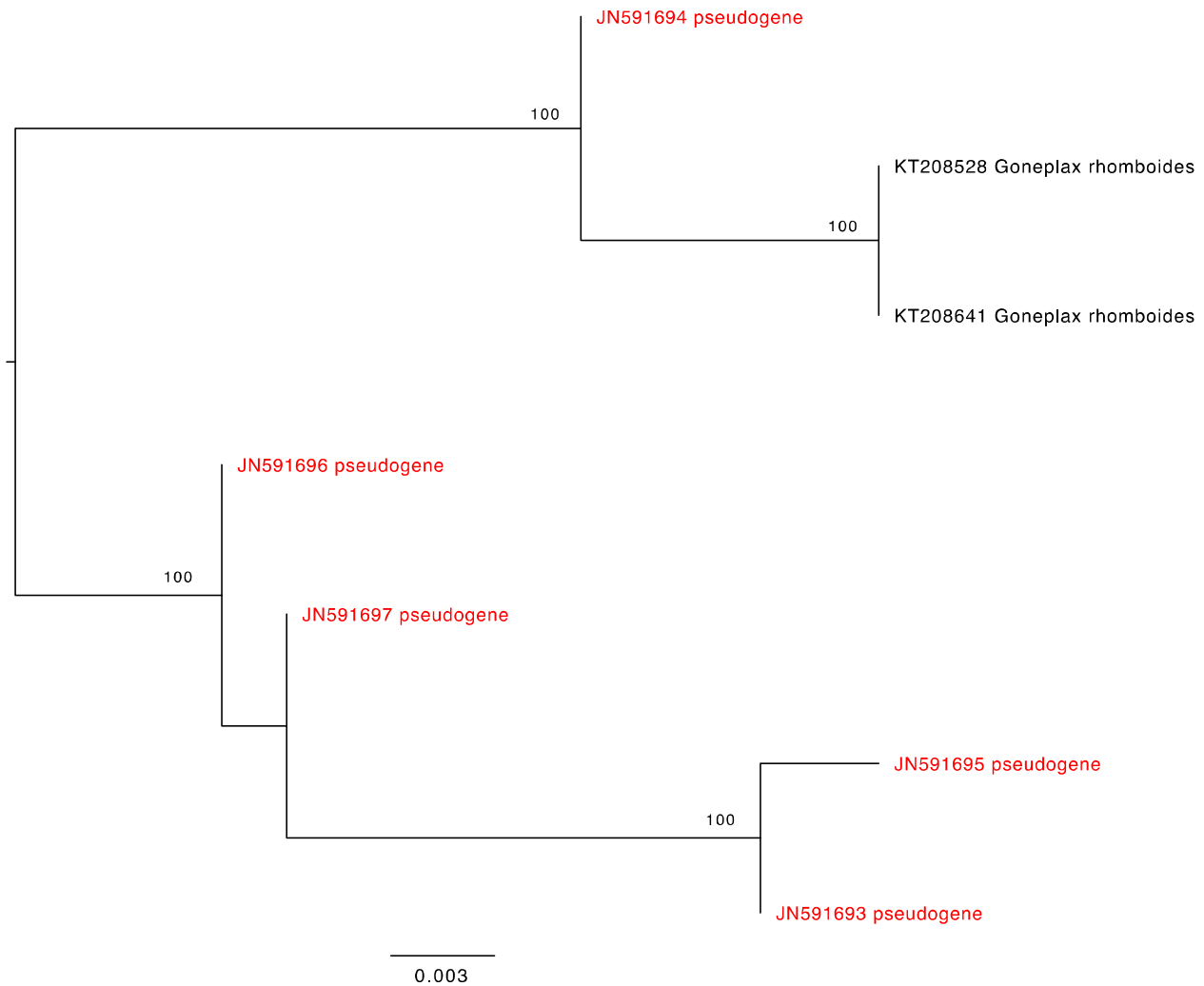


873

874

875

876 **Fig S9. Several *Goneplax rhomboides* COI pseudogene sequences cluster**  
877 **together.** A mid-point rooted neighbor joining phylogram using the Kimura 2-parameter  
878 model of nucleotide substitution included COI gene sequences as well as sequences  
879 annotated in GenBank as a nuclear copy of a mitochondrial gene (red). Nodes with  
880 greater than 70% bootstrap support are labelled.  
881

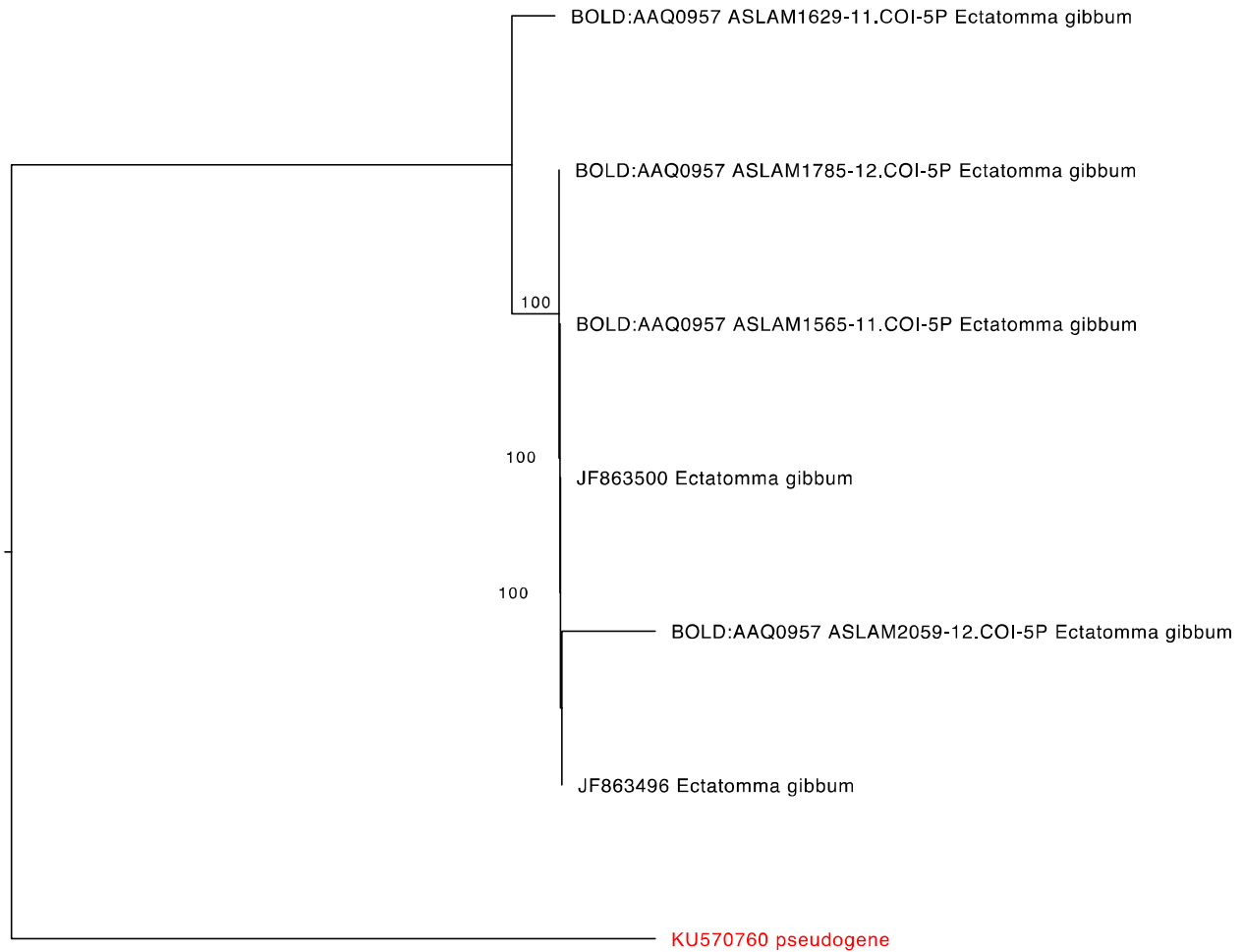


882

883

884

885 **Fig S10. A single *Ectatomma gibbum* COI pseudogene sequence is found on its**  
886 **own branch.** A mid-point rooted neighbor joining phylogram using the Kimura 2-  
887 parameter model of nucleotide substitution included COI gene sequences as well as a  
888 sequence annotated in GenBank as a nuclear copy of a mitochondrial gene (red).  
889 Nodes with greater than 70% bootstrap support are labelled.  
890

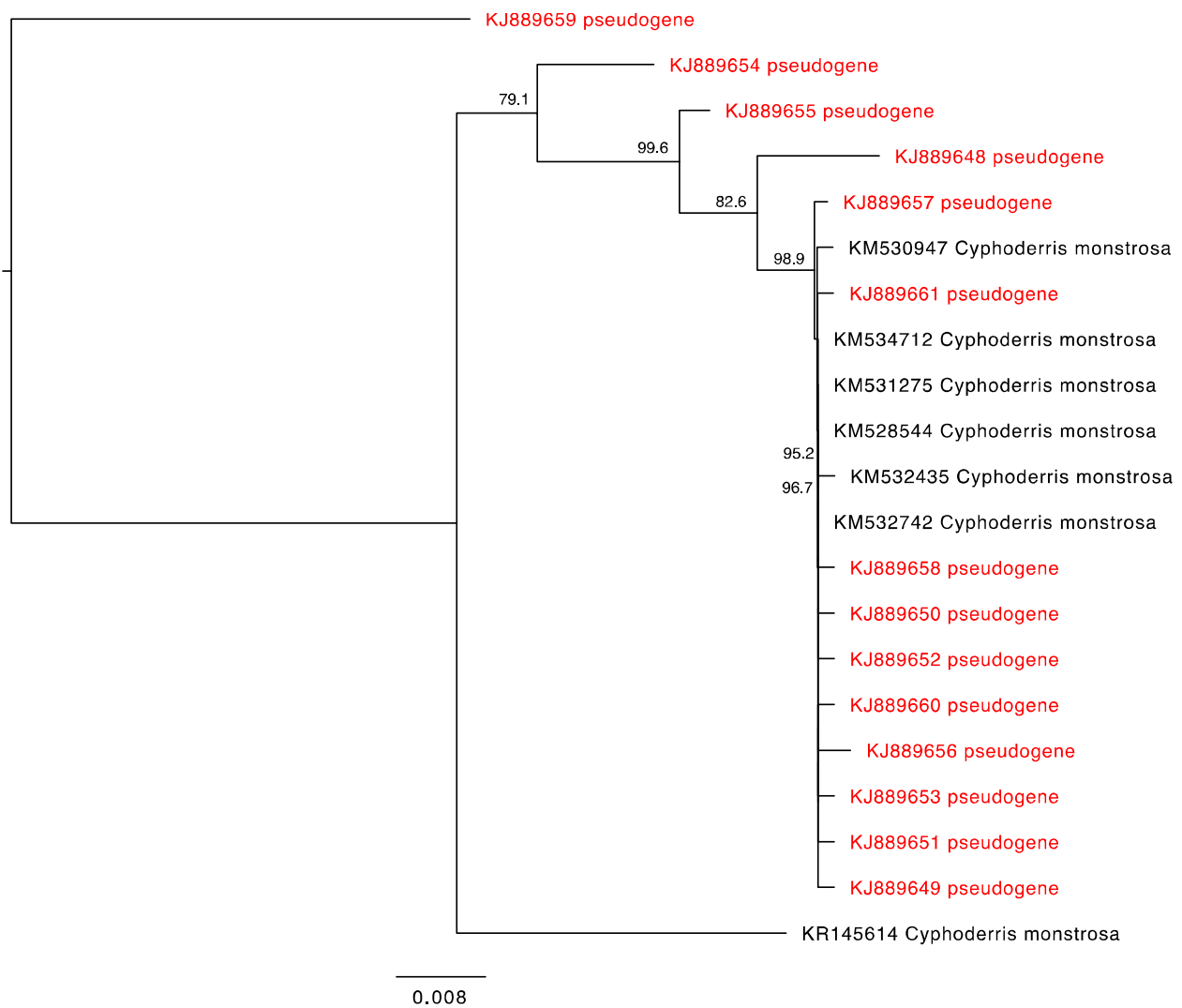


891

0.001

892

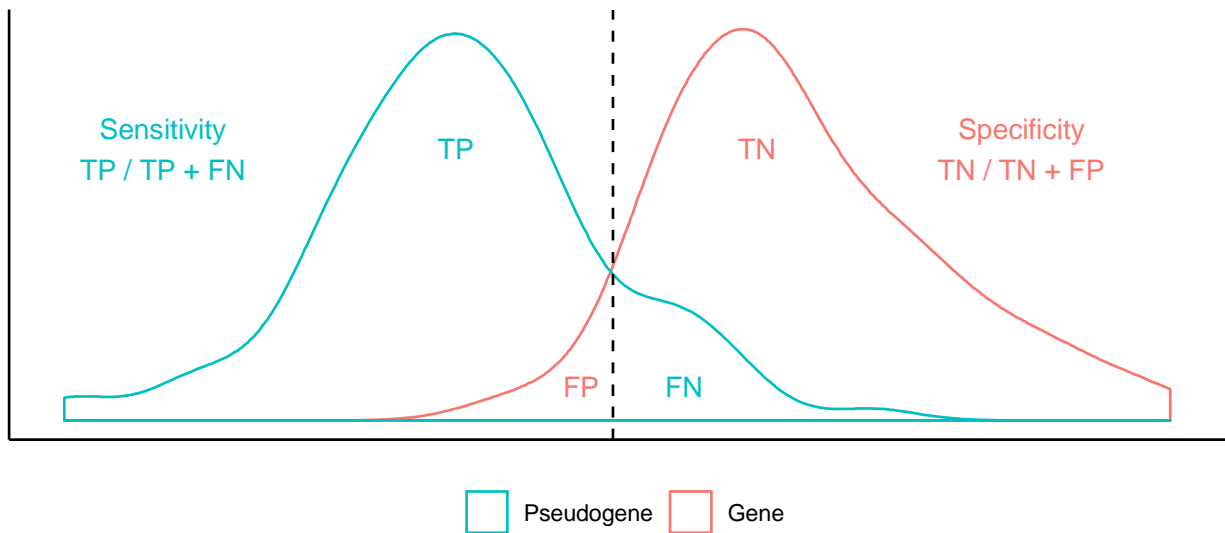
893 **Fig S11. *Cyphoderris monstrosa* COI gene and annotated pseudogene sequences**  
894 **sometimes cluster with regular gene sequences.** A mid-point rooted neighbor  
895 joining phylogram using the Kimura 2-parameter model of nucleotide substitution  
896 included COI gene sequences as well sequences annotated in GenBank as a nuclear  
897 copy of a mitochondrial gene (red). Nodes with greater than 70% bootstrap support are  
898 labelled.  
899



900

901

902 **Fig S12. Sensitivity and specificity were used to assess the effectiveness of our**  
903 **two pseudogene filtering approaches.** The vertical dashed line represents a  
904 threshold used to delimit COI pseudogene sequences. The ability to detect  
905 pseudogenes represents the positive condition. Correctly removed pseudogenes are  
906 true positives (TP). Incorrectly filtered COI gene sequences (genes) represents false  
907 positives (FP). Correctly retained genes represents true negatives (TN). Incorrectly  
908 retained pseudogenes represents false negatives (FN).  
909

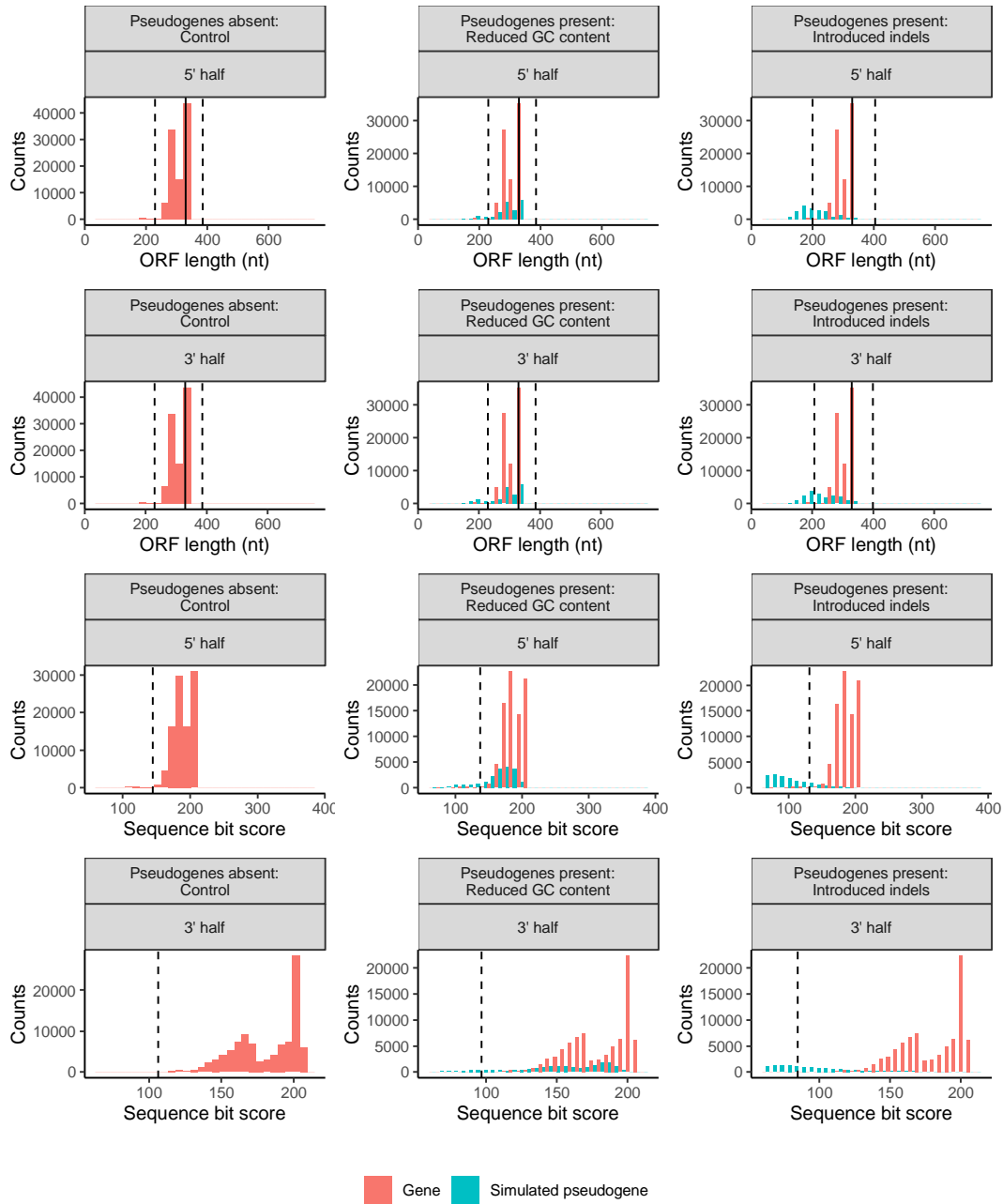


910

911

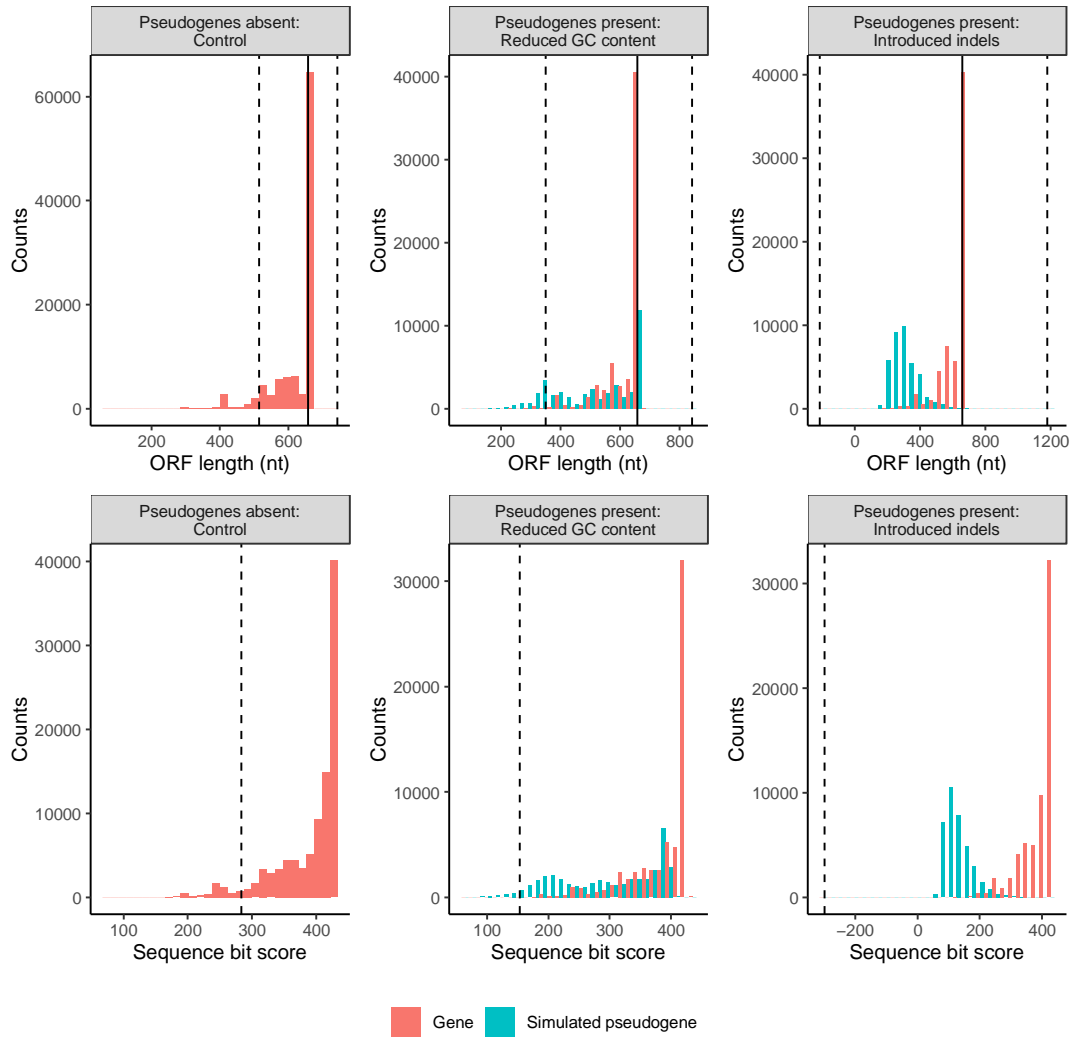
912 **Fig S13. Halving COI sequence lengths results in fewer simulated pseudogenes**  
913 **removed compared with full length COI barcode sequences.** Each column shows  
914 the results from a particular simulation: a controlled community with pseudogenes  
915 absent, a community with simulated pseudogenes with a reduced GC content, and a  
916 community with simulated pseudogenes with introduced indels. The top two panels  
917 show the length variation of sequences in the longest retained open reading frame for  
918 short sequences sampled from the 5' and 3' end of COI barcode sequences. The solid  
919 vertical line indicates half the length of a typical COI barcode at 329 bp. The two  
920 vertical dashed lines shows the boundaries for identifying ORFs with outlier lengths.  
921 The bottom two panels show the nucleotide bit score for short sequences sampled from  
922 the 5' and 3' ends of COI barcode sequences. The dashed vertical line shows the  
923 boundary for identifying sequences with unusually short scores.





924

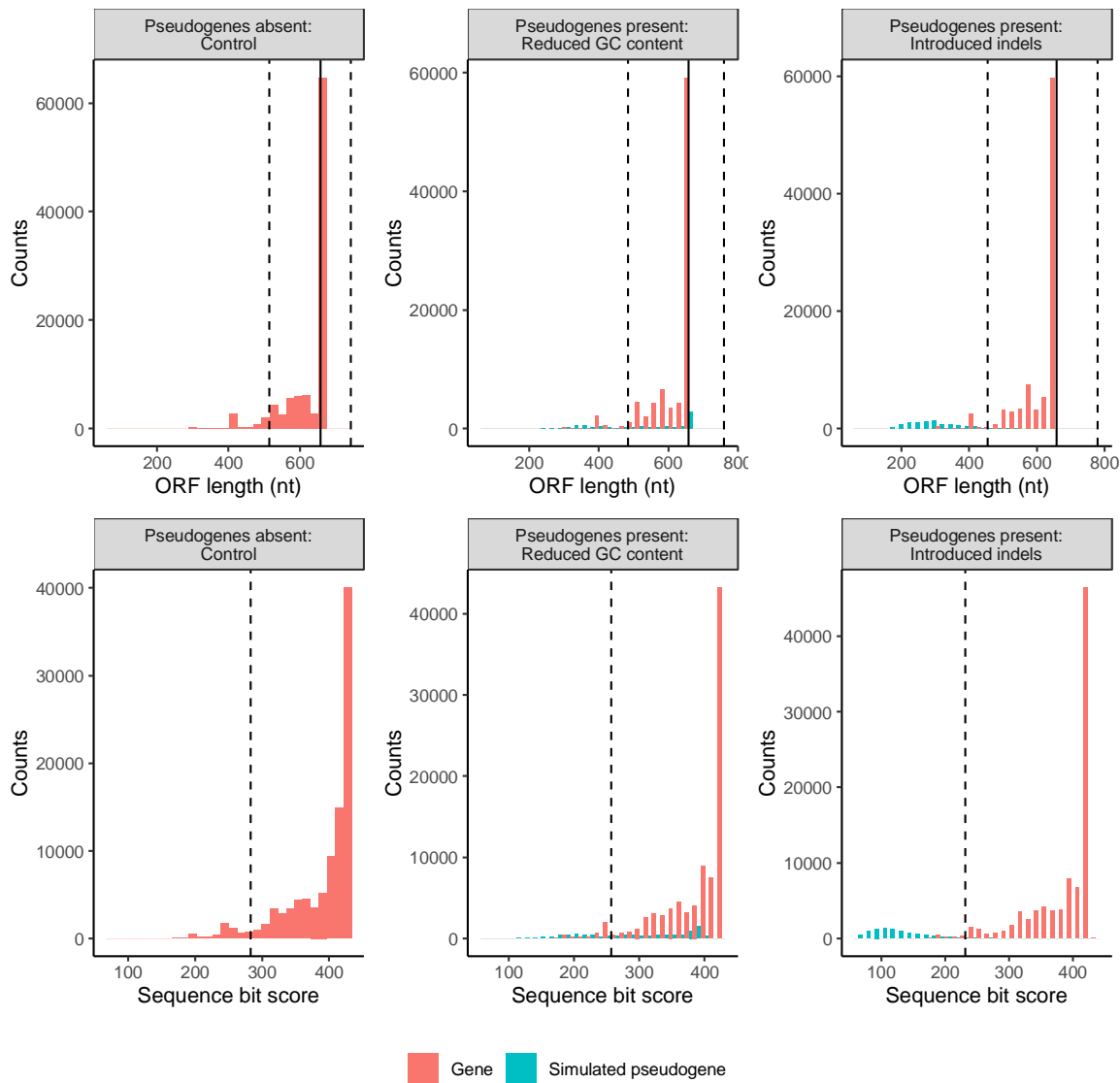
925 **Fig S14. Doubling the proportion of mutated sequences greatly reduces the**  
926 **number of simulated pseudogenes removed.** Each column shows the results from a  
927 particular simulation: a controlled community with pseudogenes absent, a community  
928 with pseudogenes that have a reduced GC content, and a community with  
929 pseudogenes where we have introduced indels. The top panel shows the length  
930 variation of sequences in the longest retained open reading frame. The solid vertical  
931 line indicates the length of a typical COI barcode at 658 bp. The two vertical dashed  
932 lines shows the boundaries for identifying ORFs with outlier lengths. The bottom panel  
933 shows the sequence bit score variation. The vertical dashed line shows the boundary  
934 for identifying sequences with small outlier scores.



935

936

937 **Fig S15. Halving the proportion of mutated sequences increases the number of**  
938 **simulated pseudogenes removed.** Each column shows the results from a particular  
939 simulation: a controlled community with pseudogenes absent, a community with  
940 pseudogenes that have a reduced GC content, and a community with pseudogenes  
941 where we have introduced indels. The top panel shows the length variation of  
942 sequences in the longest retained open reading frame. The solid vertical line indicates  
943 the length of a typical COI barcode at 658 bp. The two vertical dashed lines shows the  
944 boundaries for identifying ORFs with outlier lengths. The bottom panel shows the  
945 sequence bit score variation. The vertical dashed line shows the boundaries for  
946 identifying sequences with short outliers scores.



947

948