

1 **Molecular Ecology Resources – Permanent Genetic Resources**

2 **Using a butterflyfish genome as a general tool for RAD-Seq studies in specialized reef**
3 **fish**

4
5 Running title: **RAD-Seq in butterflyfish**

6
7 Joseph D. DiBattista^{1,2*}, Pablo Saenz-Agudelo^{1,3}, Marek J. Piatek⁴, Xin Wang¹, Manuel
8 Aranda¹, and Michael L. Berumen¹

9
10 ¹*Red Sea Research Center, Division of Biological and Environmental Science and*
11 *Engineering, King Abdullah University of Science and Technology, Thuwal 23955, Saudi*
12 *Arabia,* ²*Department of Environment and Agriculture, Curtin University, PO Box U1987,*
13 *Perth, WA 6845, Australia,* ³*Instituto de Ciencias Ambientales y Evolutivas, Universidad*
14 *Austral de Chile, Valdivia 5090000, Chile,* ⁴*Computational Bioscience Research Center,*
15 *King Abdullah University of Science and Technology, Thuwal 23955, Saudi Arabia*

16
17 *Correspondence: Joseph DiBattista, Department of Environment and Agriculture, Curtin
18 University, PO Box U1987, Perth, WA 6845, Australia

19
20 E-mail: josephdibattista@gmail.com

21
22
23
24
25
26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44 **Keywords: adaptation, Chaetodontidae, ddRAD, ecological genomics, next-generation**
45 **sequencing, SNP**

46 **Abstract**

47 Data from a large-scale restriction site associated DNA (RAD-Seq) study of nine butterflyfish
48 species in the Red Sea and Arabian Sea provided a means to test the utility of a recently
49 published draft genome (*Chaetodon austriacus*) and assess apparent bias in this method of
50 isolating nuclear loci. We here processed double-digest restriction-site (ddRAD) associated
51 DNA sequencing data to identify single nucleotide polymorphism (SNP) markers and their
52 associated function with and without our reference genome to see if it improves the quality of
53 RAD-Seq markers. Our analyses indicate (1) a modest gap between the number of non-
54 annotated versus annotated SNPs across all species, (2) an advantage of using genomic
55 resources for closely related but not distantly related butterflyfish species based on the ability
56 to assign putative gene function to SNPs, and (3) an enrichment of genes among sister
57 butterflyfish taxa related to calcium transmembrane transport and binding. The latter result
58 highlights the potential for this approach to reveal insights into adaptive mechanisms in
59 populations inhabiting challenging coral reef environments such as the Red Sea, Arabian Sea,
60 and Arabian Gulf.

61

62 **Introduction**

63 Ecologists and evolutionary biologists often mandate genetic data to test hypotheses related
64 to recent population history (i.e. bottlenecks), connectivity among sites, and the spatial
65 distribution of genetic variation for effective species management. Moreover, the capacity to
66 rapidly genotype many individuals' at large numbers of genetic markers improves our ability
67 to estimate demographic parameters (e.g. gene flow, effective population size, admixture),
68 resolve phylogenetic placement, identify genes that may be under selection, and even identify
69 genes that facilitate adaptation to our rapidly changing environment. Now, with the

70 decreasing cost of high-throughput next-generation sequencing (NGS), data can be sampled
71 from across the genome to meet this demand (for review see Kosuri & Church 2014).

72 Restriction site associated DNA sequencing (RAD-Seq) methods, which use NGS to
73 target sequence data adjacent to restriction enzyme recognition sites (Davey *et al.* 2011),
74 have generated informative single nucleotide polymorphism (SNP) datasets in several fish
75 species. Examples include studies characterizing different trout species and their hybrids
76 (Hand *et al.* 2015), the genetic basis for different forms of stickleback (e.g. oceanic versus
77 freshwater; Hohenlohe *et al.* 2010; Catchen *et al.* 2013a), identifying cryptic lineages of
78 herring in the Baltic Sea (*Clupea harengus*; Corander *et al.* 2013), and resolving relationships
79 among iconic African cichlid species (Wagner *et al.* 2013; Henning *et al.* 2014). Most of
80 these taxa, however, are considered “model organisms”, and therefore unique in their ability
81 to map, align, and otherwise trace their SNPs back to annotated genomes.

82 Reef fish, on the other hand, have received much less attention, despite 5000 species
83 inhabiting tropical seas. We note that RAD-Seq methods have only been applied a few times
84 in reef fish, including surgeonfish (Gaither *et al.* 2015), clownfish (Saenz-Agudelo *et al.*
85 2015), hamlets (Puebla *et al.* 2014; Picq *et al.* 2016), angelfish (Tariel *et al.* 2016), grunts
86 (Bernal *et al.* 2016), groupers (Jackson *et al.* 2014), and parrotfish (Stockwell *et al.* 2016).
87 Most of these studies, however, are limited in scope (based on small sample sizes, few study
88 species, and minimal sampling sites) or fail to confidently assign gene function to SNPs of
89 interest; all rely exclusively on *de novo* assembly of raw sequence reads. This deficiency in
90 reef fish relative to other aquatic organisms can therefore be, at least in part, attributed to a
91 lack of publicly available genomic resources. All published genomes to date in this group are
92 based on cold water (*Takifugu rubripes*; van de Peer 2004), brackish water (*Tetraodon*
93 *nigroviridis*; van de Peer 2004), pelagic (*Thunnus orientalis*; Nakamura *et al.* 2013), and
94 phylogenetically distinct (*Rhincodon typus*; Read *et al.* 2015) species.

95 A recently published genome for the Red Sea butterflyfish (*Chaetodon austriacus*;
96 DiBattista *et al.* 2016a), an obligate corallivore often targeted by the ornamental fish trade
97 (Wabnitz 2003; Lawton *et al.* 2013), allows us to explore genes underpinning ecological
98 niche selection (Cole *et al.* 2008), evolutionary distinctness (Bellwood *et al.* 2010), but more
99 importantly biogeographic patterns related to the specialized and speciose Chaetodontidae
100 family (> 130 species). Indeed, the Red Sea has experienced intermittent historical isolation
101 (Bailey 2009), fluctuating environmental conditions (e.g. Raitso *et al.* 2013), and harbors
102 large numbers of endemic species (DiBattista *et al.* 2016b) whose origins are still subject to
103 debate (DiBattista *et al.* 2016c). The *C. austriacus* genome may therefore enhance RAD-seq
104 analysis in related species.

105 We therefore have two methodological aims in this study. First, we test whether using the
106 *C. austriacus* genome as a reference improves the quality of RAD-Seq markers for this
107 species, which includes assessing the apparent bias in our method of isolating nuclear loci.
108 Second, to increase the utility of this genome for broader topics of study, we further test the
109 capacity of this genome to serve as a scaffold for RAD-Seq markers in other butterflyfish
110 species. This is the first time comparisons between RAD-Seq analyses, with and without a
111 representative genome, have been conducted across such a large number of related, aquatic
112 organisms.

113

114 **Material and methods**

115 *Sample collection*

116 We collected fin or gill tissue from between 48 and 108 individuals for each of nine species
117 of butterflyfish sampled at sites in the Red Sea and Arabian Sea, including the genomically
118 enabled *C. austriacus* (see Fig.1 and Table 1). We also collected samples from a surgeonfish

119 species (*Ctenochaetus striatus*; $N = 120$ from 10 populations), which served as an outgroup
120 for our analyses. All samples were preserved in 95% ethanol prior to processing.

121

122 *Restriction site associated DNA sequencing method*

123 The RAD-Seq method involves digesting genomic DNA with restriction enzymes and
124 sequencing fragments of DNA adjacent to restriction sites. In our case, we used the double-
125 digest RAD tag method (ddRAD; Peterson *et al.* 2012), which uses one common and one rare
126 restriction enzyme cutter but no shearing step, although several other methods exist (Willette
127 *et al.* 2014; Andrews *et al.* 2016).

128 In brief, genomic DNA was extracted using a Qiagen DNeasy Blood and Tissue Kit
129 (Qiagen, Valencia, CA) following the manufacturer's protocol. Total DNA from each
130 extracted aliquot was quantified using a Qubit dsDNA HS Assay Kit (Invitrogen, Carlsbad,
131 CA). Genomic libraries were prepared from 500 ng of DNA per sample by: 1) digesting at
132 37°C using the restriction enzymes *SphI* and *MluCI* (NEB), 2) ligating to unique
133 combinations of custom adaptors, 3) pooling 16 individuals at a time, 4) size-selecting 300 to
134 500 bp fragments on agarose gels from each pool, 5) amplifying over 10 PCR cycles to
135 reduce clonality in 50 μ l reactions containing 25 μ l Illumina True Seq Master Mix, 20 μ l of
136 library DNA, and unique indexing primers that correspond to the standard Illumina
137 multiplexed sequencing protocol, and 6) combining pools, as appropriate, in equimolar
138 concentration to form nine genomic libraries, which were then run on nine lanes of an
139 Illumina HiSeq2000 (101 bp; v3 reagents) at the King Abdullah University of Science and
140 Technology (KAUST) Bioscience Core Laboratory using the SE option.

141

142 *RAD-Seq analysis*

143 RAD sequences were first analyzed *de novo* (without a reference genome) using Stacks *vers.*
144 1.44 (Catchen *et al.* 2011, 2013b). Raw reads for each individual in the RAD tag library were
145 trimmed at the end to a common length of 81 bp in FASTQ format using the
146 “process_radtags” pipeline since read quality was considerably lower from this position
147 onwards; reads with an average phred score < 20 (in a 5 bp sliding window) were discarded.
148 All remaining reads from each species were separately analyzed and SNPs were genotyped
149 using the “denovo_map.pl” pipeline. The parameters used in Stacks were: minimum number
150 of reads required to form a stack (m) was set to 3; maximum number of mismatches between
151 loci for individuals (M) was set to 4; maximum number of mismatches when aligning
152 secondary reads to primary stacks (N) was set to 2; maximum number of mismatches
153 between loci when creating a catalog (n) was set to 2. This parameter combination yielded the
154 closest number of SNPs for our focal species (*C. austriacus*) compared to the number of
155 SNPs recovered using the “ref_map.pl” pipeline in Stacks (see below and Appendix S1), and
156 was therefore used as the optimal settings for all the other species. Population filtering
157 parameters used in order to include a locus in the final data set were: 1) minor allele
158 frequency > 0.05 and 2) the locus had to be genotyped in at least 80% of individuals
159 (populations: -r) and present in all (strict quality filter) or all but one (relaxed quality filter)
160 populations (populations: -p)(for more detail see Saenz-Agudelo *et al.* 2015).

161 RAD sequences were also analyzed using the *C. austriacus* genome as a reference
162 (DiBattista *et al.* 2016a). For this approach, trimmed RAD reads from each species (and thus
163 SNPs) were first aligned to the *C. austriacus* scaffolds using Bowtie *vers.* 1.0.1 (Langmead *et*
164 *al.* 2009). We configured Bowtie to report only one (i.e. the best) alignment per read and
165 allowed up to 3 mismatches in the default seed length. Aligned reads were then analyzed with
166 ref_map.pl using the same minimum number of reads to build a stack as outlined above for
167 denovo_map.pl; population filtering parameters also remained the same (see Appendix S1).

168 This step additionally allowed us to estimate homology levels for all studied species against
169 the *C. austriacus* genome.

170 Complimentary Maximum Likelihood (ML) phylogenetic trees were built with RAD
171 sequences using an unpartitioned GTRCAT approximation model in RAxML *vers.* 2.0
172 (Stamatakis 2006) as implemented in Geneious *vers.* 10.0.2 (Drummond *et al.* 2009)(also see
173 Ree & Hipp 2015). The input for these analyses consisted of a full sequence (.phylip)
174 alignment generated in STACKS (denovo_map: -m 3 -M 4 -N 2 -n 6) by compiling four
175 individuals per species, resulting in 38,070 total bp representing 2846 shared SNP loci (2170
176 fixed SNPs and 676 variable SNPs within species). Only loci present in all 10 species and at
177 least 3 individuals per species were included. Branch support was based on 100,000 rapid-
178 bootstrap replicates with *Ct. striatus* used as an outgroup. Trees were visualized using
179 FigTree *vers.* 1.4.2 (<http://tree.bio.ed.ac.uk/software/figtree/>).

180 We next explored the form and function of SNPs that mapped to the *C. austriacus*
181 scaffolds in several ways. First, we identified SNPs from *C. austriacus* and its closest
182 relative, *C. melapterus*, which mapped to regions containing gene models. These represent
183 the only two species where using the *C. austriacus* reference genome increased the number of
184 SNPs identified (see below). The associated genes were then used to look for potential
185 enrichment of specific gene functions using GO enrichment analyses. We identified 3775
186 SNPs from the *C. austriacus* dataset mapping to 3056 different genes in the reference
187 genome of *C. austriacus*, and 2341 SNPs from the *C. melapterus* dataset mapping to 2011
188 different genes of the reference genome. GO enrichment analyses were performed using
189 TopGo *vers.* 2.22.0 from the R Bioconductor package (Alexa & Rahnenfuhrer 2010), the
190 weight01 model, a cut-off of $P < 0.01$, and specific background gene sets covered by stacks
191 from the ref_map.pl analysis. These analyses were used to find which GO terms were over-
192 represented (or under-represented) using annotations for that gene set. Second, we estimated

193 the number of SNPs from *C. austriacus*, with corresponding tolerance intervals (95%
194 confidence in 99% of the sample), that mapped to the assembled genome as a function of
195 scaffold size. This particular analysis highlights scaffolds that host a statistically higher
196 number of SNPs for their length. Frequency distributions for the number of SNPs and
197 scaffold lengths were also calculated.

198

199 **Results**

200 We used Stacks to assemble and genotype SNPs in two different ways. First, we used
201 `denovo_map.pl` to assemble and genotype SNPs *de novo* for ten different species, including
202 *C. austriacus*, the only butterflyfish with a fully sequenced genome. We also used the
203 `ref_map.pl` pipeline to assemble and map reads from each of these species to the *C.*
204 *austriacus* scaffolds, which are now publicly available (NCBI Bioproject PRJNA292048;
205 <http://caus.reefgenomics.org>; see Liew *et al.* 2016). Using an optimal parameter combination
206 (see Appendix S1) and `denovo_map.pl`, we found 289,504 loci, with 8,842 (strict quality
207 filter) or 10,711 (relaxed quality filter) variable loci containing at least one SNP site within or
208 between individuals for *C. austriacus* (Table 1). On average, there were 262,247 loci across
209 all butterflyfish species (range: 189,554 to 363,221), and between 476 (*Chaetodon*
210 *trifascialis*) and 10,257 (*Chaetodon paucifasciatus*) loci per species passing our filters based
211 on strict criteria, and between 1,271 (*Chaetodon trifascialis*) and 13,539 (*Chaetodon*
212 *paucifasciatus*) loci per species passing our filters based on relaxed criteria (Table 1).
213 Average depth of coverage ranged from 10.2X (*Chaetodon fasciatus*) to 16.2X (*Chaetodon*
214 *paucifasciatus*) for `denovo_map.pl` (average depth of coverage = 12.3X), and from 12.1X
215 (*Chaetodon fasciatus*) to 16.1X (*Chaetodon paucifasciatus*) for `ref_map.pl` (average depth of
216 coverage = 13.5X).

217 We note a 38% (strict quality filter) or 56% (relaxed quality filter) increase in the number
218 SNPs identified using ref_map.pl versus denovo_map.pl for *C. melapterus* (Table 1), the
219 closest relative to *C. austriacus*. In contrast, there was a 2- to 6-fold decrease (strict quality
220 filter) or 4- to 12-fold decrease (relaxed quality filter), in the number of SNPs identified for
221 the remaining butterflyfish species using ref_map.pl (excluding *C. austriacus*; Table 1). This
222 downward trend was more pronounced when we considered fish outside of the
223 Chaetodontidae family (i.e. *Ct. striatus*, family Acanthuridae), with a 52-fold (strict quality
224 filter) or a 56-fold (relaxed quality filter) decrease in the number of SNPs identified. Indeed,
225 our surgeonfish outgroup, *Ct. striatus*, had only 8 (strict quality filter) or 27 (relaxed quality
226 filter) useable SNP loci after filtering when using ref_map.pl.

227

228 *Gene function under putative selection*

229 In order to determine if variable SNPs in *C. austriacus* and its closest relative, *C. melapterus*,
230 were associated with specific biological processes, we performed GO enrichment analyses on
231 SNPs lying within gene models. To this end, we first analysed if SNPs were randomly
232 distributed with respect to gene function. Analysis of the GO terms covered by the species-
233 specific SNPs showed a strong bias towards particular gene functions, with the strongest
234 enrichment for genes involved in the positive regulation of GTPase activity, axon guidance,
235 and chloride transmembrane transport, among others.

236 To determine if this observed bias in gene function was a consequence of our RAD-seq
237 method and thus chosen restriction enzyme, we analysed the gene models covered by our loci
238 in comparison to all proteins. Briefly, the *SphI* recognition site (GCATGC) contains an ATG
239 followed by a C, which in coding sequences translates to a methionine (M) followed by a
240 leucine (L). Consequently, using this restriction enzyme could enrich for genomic regions
241 encoding proteins starting with or harboring internal Methionine-Leucine (ML) amino acid

242 stretches. We therefore analyzed the frequency of ML amino acid stretches in all gene models
243 for *C. austriacus* and compared them to the frequency within the set of gene models
244 associated with our loci. This analysis revealed a highly significant enrichment of such genes
245 within our loci, thus identifying this as a source of the previously observed bias (Chi-square
246 test: $P < 0.0001$). Based on this result we generated species-specific GO term backgrounds
247 for our subsequent enrichment analyses to account for the observed bias, which may be an
248 important consideration for other RAD-seq approaches.

249 SNP specific GO enrichment analyses identified 29 and 37 biological processes in *C.*
250 *austriacus* and *C. melapterus*, respectively, that were significantly enriched ($P < 0.01$), of
251 which six were shared between the species. These shared terms included calcium ion
252 transmembrane transport, axon guidance, neuron cell-cell adhesion, positive regulation of
253 excitatory postsynaptic potential, and positive regulation of cardiac muscle cell proliferation.
254 Variable SNPs exclusively found in *C. austriacus* were predominantly associated with genes
255 involved in response to regulation of cytosolic calcium concentrations, morphogenesis, and
256 locomotory behavior (Figure 2a), whereas SNPs exclusively found in *C. melapterus* were
257 enriched for functions involved in ion transmembrane transport, synaptic transmission, and
258 social behavior (Figure 2b). Moreover, we found that only 385 SNPs (i.e. 4%) identified with
259 the *denovo_map.pl* for *C. austriacus* approach failed to map to the reference genome,
260 indicating that both approaches produced representative SNPs.

261 We next mapped SNP loci identified from *C. austriacus* with *ref_map.pl* to the gap-
262 closed scaffolds in order to track the distribution of SNP density across its own genome. We
263 found that (as expected) the number of SNPs mapping to the reference genome increased
264 with scaffold size ($R^2 = 0.71961$; Fig. 3c). That said, 3055 scaffolds had ≤ 10 SNPs that
265 mapped back to them versus 115 scaffolds with ≥ 10 SNPs that mapped back to them (Fig.
266 3a) despite a large mean scaffold size (~ 150 kb; Fig. 3b). These two categories (i.e. > 10

267 SNPs per scaffold versus < 10 SNPs per scaffold) represent sets of SNPs above and below
268 the tolerance threshold, respectively, with a clear difference in GO enrichment between them
269 (see Appendix S2).

270

271 *Phylogenetic overview*

272 In order to test the phylogenetic generality of these analyses, we assigned SNPs identified
273 within each butterflyfish species, and the outgroup surgeonfish, to annotated versus non-
274 annotated regions of the *C. austriacus* scaffolds. We found that approximately 16% more
275 SNPs could be assigned to non-annotated versus annotated regions of the genome, and that
276 the absolute number of SNPs assigned varied among species (Fig. 4a). An estimate of the
277 percent homology to the *C. austriacus* reference genome for all RAD data shows that *C.*
278 *austriacus* and its closest relative, *C. melapterus*, mapped at a level of 86% and 87%
279 respectively, with the rest of the butterflyfish mapping between 19% and 28% similarity; the
280 outgroup surgeonfish had the lowest level of homology (~2%) (Fig. 4b). This phylogenetic
281 decay is consistent with the inferred relationships among butterflyfish species revealed by a
282 Maximum Likelihood consensus tree built with shared SNP loci (Fig. 4c).

283

284 **Discussion**

285 We compared analysis of RAD-Seq data using a *de novo* approach (i.e. without a genome) to
286 analysis using a reference genome for a series of Red Sea and Arabian Sea butterflyfish. Our
287 analyses indicate (1) a modest gap between the number of non-annotated versus annotated
288 SNPs across all species, (2) an advantage of using genomic resources for closely related but
289 not distantly related butterflyfish species based on the ability to assign putative gene function
290 to SNPs, and (3) an enrichment of genes among sister butterflyfish taxa related to calcium
291 transmembrane transport and binding. The RAD-Seq dataset we present here thus provides

292 the most comprehensive estimate of genetic polymorphism to date for any reef fish, and
293 further creates the potential to elevate the Chaetodontidae family to a ‘model group’ for
294 future genomic studies. Overall, our results suggest that RAD-Seq approaches inherently
295 provide valuable insight to population genomic questions for non-model organisms (i.e. those
296 lacking reference genomes), but that detailed analysis of RAD-Seq data is improved when a
297 reference genome is available from a closely related species (as per Pecoraro *et al.* 2015) or
298 when apparent bias in the method of SNP isolation can be identified and accounted for.

299 Most RAD tags are thought to be randomly distributed across the genome based on the
300 stochastic placement of restriction cut sites (Davey *et al.* 2013), and indeed we found no
301 evidence for what we here refer to as SNP “hotspots” (e.g. Myers *et al.* 2005; also see Fig. 3).
302 The genome assembly of *C. austriacus*, however, consists of gap-closed scaffolds and does
303 not include a linkage map (DiBattista *et al.* 2016a); the exact location of each SNP relative to
304 autosomal (or sex-linked) chromosomes is therefore not yet known. Moreover, we enabled
305 the “select one SNP per read” filtering option in Stacks in order to comply with population
306 genetic assumptions of independent loci, our intended downstream application for these data,
307 which further reduced our ability to detect multiple SNPs sitting within 81 bp of each other.
308 Despite this uncertainty, a large portion of the identified SNPs could be assigned to annotated
309 regions of the genome for most of the butterflyfish species we considered (average number of
310 annotated SNPs = 1199; Fig. 4a). We did, however, detect a strong bias with respect to
311 covered gene functions (i.e. proteins containing ML stretches), which apparently stems from
312 our choice of restriction enzyme. Restriction enzymes that target gene rich regions (e.g. Roda
313 *et al.* 2013), or even specific primers flanking variable SNPs (Campbell *et al.* 2014), have
314 previously been used when more specificity is required. This refinement provides an
315 advantage for researchers who may wish to select a smaller number of SNP sites for
316 genotyping with known spacing, location or function, but also a direct application to projects

317 estimating connectivity and adaptability of fishes. However, such approaches are likely to
318 introduce their own biases that need to be considered in subsequent analyses as we have
319 shown here. Moreover, if a mutation has occurred in the cut-site of some (but not all)
320 individuals, allelic dropout can occur, which has been shown to further bias the inference of
321 genetic variability (Gautier *et al.* 2013). Without an available reference genome, these types
322 of methodological biases are missed.

323 We observed a clear advantage to using genomic resources for closely related but not
324 distantly related butterflyfish species in terms of the ability to assign SNPs a putative
325 function. This advantage would likely extend to other sister taxa of reef fish, and as such it
326 would be ideal for studies focused on between-species comparisons *within* monophyletic
327 lineages. Similar advantages were recently shown for the Pacific bluefin tuna (*Thunnus*
328 *orientalis*; Pecoraro *et al.* 2015). For example, in our case, the availability of the *C.*
329 *austriacus* genome sequence allowed for functional interrogation of the data using GO
330 annotations within sister taxa, which revealed that variable SNPs were significantly enriched
331 within certain gene functions (Fig. 2). These functions included genes associated with
332 neurological processes, which might indicate potential differences in behavior or sensory
333 perception for these specialized reef fish *within* or *between* the different species, and
334 therefore represent plausible candidate genes for future functional interrogation. This is
335 consistent with the hypothesis that detecting and interpreting visual and olfactory cues are
336 highly important for butterflyfish given that most of these fish are monogamous and need to
337 maintain pair bonds, recognize mates, and defend territories from conspecifics (e.g. Boyle &
338 Tricas 2014). More importantly we found that genes associated with calcium transmembrane
339 transport and binding were significantly overrepresented in both *C. austriacus* and *C.*
340 *melapterus*, suggesting potential convergent adaptations to the prevailing Red Sea and
341 Arabian Sea environments, respectively. A plausible alternative is that these SNPs are not

342 themselves adaptive but linked with other regions of the genome that were not sequenced,
343 resulting in a consistent association and statistically significant enrichment of sequence
344 variation and specific gene functions. If we assume that these functions are under selection,
345 calcium/chlorine ratios are remarkably similar in the Red Sea and adjacent Indian Ocean
346 (Krumgalz 1982), although the Gulf of Aqaba is much lower, suggesting that adaptation to
347 these conditions may dictate which gene functions are conserved. The assimilation of calcium
348 from seawater is an important process for the formation of structural skeleton in corals,
349 protective shells for many of the marine invertebrates, as well as promotes optimal larval
350 development in fish via ossification and gill formation (e.g. Malvezzi *et al.* 2015). Further
351 study of additional fish from the Red Sea and Arabian Sea region could test the hypothesis
352 that genetic adaptations allow species to handle the unique chemical challenges present in
353 regional waters, and particularly could provide some insight to the mechanisms underlying
354 high rates of endemism in the region.

355 We note that the percentage of reads mapping to the reference genome was actually
356 higher for *C. melapterus* (87.5%) versus the genomically enabled *C. austriacus* (86.4%) (Fig.
357 4b). Although this result may appear unusual, these species are closely related to each other
358 (Fig. 4c), have only recently diverged (~50 Kya), share mitochondrial DNA haplotypes (see
359 Waldrop *et al.* 2016), and likely hybridize in regions of overlap based on similar behavior
360 observed in nearby seas (DiBattista *et al.* 2015). Recent and/or on-going genetic exchange
361 suggests that introgression across the nuclear genome must be high for these two species (see
362 Waldrop *et al.* 2016). Moreover, despite phylogenetic similarity, *C. melapterus* had lower
363 levels of polymorphism (0.78% versus 2.3%, respectively) and fewer restriction sites
364 recovered (5913 versus 8812 variable loci passing strict criteria, respectively) compared to *C.*
365 *austriacus*, and yet the average depth of coverage for the former increased from 12.2X to
366 14.9X when mapping to the reference genome (Table 1). We may therefore have achieved

367 marginally better mapping results for *C. melapterus* simply because we had a lower probably
368 of discarding reads with ref_map.pl based on increased homogeneity and higher per locus
369 coverage. We additionally attribute the large range of pre- and post-filtered loci among
370 species to differences in sample size, in addition to stochastic factors such as variable DNA
371 quality for each sample and library preparation protocol.

372

373 *Conclusion*

374 The rapid identification of a large number of SNPs within populations likewise holds promise
375 for a number of evolutionary-oriented studies. While the number publicly available genomes
376 is steadily increasing, our results suggest that increasing phylogenetic distance decreases the
377 utility of genomes, but that this remains an important means to identify apparent bias in SNP
378 isolation approaches as well as assess conserved biological processes. Substantial advances in
379 our understanding of evolutionary mechanisms may therefore be possible without the
380 substantial investment of resources required for full comparative genomic studies.

381

382 **Acknowledgements.** This research was supported by the KAUST Office of Competitive
383 Research Funds (OCRF) under Award No. CRG-1-2012-BER-002 and baseline research
384 funds to M.L.B., as well as a National Geographic Society Grant 9024-11 to J.D.D. For
385 support in Socotra, we kindly thank the Ministry of Water and Environment of Yemen, staff
386 at the Environment Protection Authority (EPA) Socotra, and especially Salah Saeed Ahmed,
387 Fouad Naseeb and Thabet Abdullah Khamis, as well as Ahmed Issa Ali Affrar from Socotra
388 Specialist Tour for handling general logistics. For logistic support elsewhere, we thank Eric
389 Mason at Dream Divers in Saudi Arabia; the Red Sea State Government and The Red Sea
390 University in Sudan, as well as Equipe Cousteau including N. Hussey, S. Kessel, C.
391 Scarpellini, and M. Younis; Nicolas Prévot at Dolphin Divers and the crew of the M/V Deli

392 in Djibouti; the KAUST Coastal and Marine Resources Core Lab and Amr Gusti; and the
393 Ministry of Agriculture and Fisheries in Oman including Abdul Karim. For specimen
394 collections, we thank Tilman Alpermann, Giacomo Bernardi, Brian Bowen, Camrin Braun,
395 Richard Coleman, Michelle Gaither, Jean-Paul Hobbs, Jennifer McIlwain, Gerrit Nanninga,
396 Mark Priest, Luiz Rocha, Tane Sinclair-Taylor, and members of the Reef Ecology Lab at
397 KAUST. For assistance with bench work at KAUST we thank Craig Michell. We also
398 acknowledge important contributions from Luiz Rocha and David Catania for assistance with
399 specimen archiving at the California Academy of Sciences; the KAUST Bioscience Core
400 Laboratory with Sivakumar Neelamegam and Hicham Mansour for their assistance with
401 Illumina sequencing.

402

403 **References**

404 Alexa A, Rahnenfuhrer J (2010) topGO: Enrichment analysis for Gene Ontology. R package
405 version 2.22.0.

406

407 Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA (2016) Harnessing the power
408 of RADseq for ecological and evolutionary genomics. *Nature Reviews Genetics*, **17**, 81-92.

409

410 Bailey R (2009) *Ecosystem Geography: From Ecoregions to Sites*. 2nd edition, New York,
411 Springer Publishing.

412

413 Bellwood DR, Klanten S, Cowman PF, Pratchett MS, Konow N, van Herwerden L (2010)
414 Evolutionary history of the butterflyfishes (f: Chaetodontidae) and the rise of coral feeding
415 fishes. *Journal of Evolutionary Biology*, **23**, 335-349.

416

417 Bernal MA, Gaither MR, Simison WB, Rocha LA (2016) Introgression and selection shaped
418 the evolutionary history of sympatric sister-species of coral reef fishes (genus: *Haemulon*).
419 *Molecular Ecology*. DOI: 10.1111/mec.13937.

420

421 Boyle KS, Tricas TC (2014) Discrimination of mates and intruders: visual and olfactory cues
422 for a monogamous territorial coral reef butterflyfish. *Animal Behaviour*, **92**, 33-43.

423

424 Campbell NR, Harmon SA, Narum SR (2014) Genotyping-in-Thousands by sequencing (GT-
425 seq): A cost effective SNP genotyping method based on custom amplicon sequencing.
426 *Molecular Ecology Resources*, **15**, 855-867.

427

428 Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011) Stacks: building
429 and genotyping loci de novo from short-read sequences. *G3: Genes, Genomes, Genetics*, **1**,
430 171-182.
431
432 Catchen J, Bassham S, Wilson T, Currey M, O'Brien C, Yeates Q, Cresko WA (2013a) The
433 population structure and recent colonization history of Oregon threespine stickleback
434 determined using restriction-site associated DNA-sequencing. *Molecular Ecology*, **22**, 2864-
435 2883.
436
437 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA (2013b) Stacks: an analysis
438 tool set for population genomics. *Molecular Ecology*, **22**, 3124-3140.
439
440 Cole AJ, Pratchett MS, Jones GP (2008) Diversity and functional importance of coral-feeding
441 fishes on tropical coral reefs. *Fish and Fisheries*, **9**, 286-307.
442
443 Corander J, Majander KK, Cheng L, Merilä J (2013) High degree of cryptic population
444 differentiation in the Baltic Sea herring *Clupea harengus*. *Molecular Ecology*, **22**, 2931-2940.
445
446 Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML (2011) Genome-
447 wide genetic marker discovery and genotyping using next-generation sequencing. *Nature*
448 *Reviews Genetics*, **12**, 499-510.
449
450 Davey JW, Cezard T, Fuentes-Utrilla P, Eland C, Gharbi K, Blaxter ML (2013) Special
451 features of RAD Sequencing data: implications for genotyping. *Molecular Ecology*, **22**, 3151-
452 3164.
453
454 DiBattista JD, Wang X, Saenz-Agudelo P, Piatek M, Aranda M, Berumen ML (2016a) Draft
455 genome of an iconic Red Sea reef fish, the blacktail butterflyfish (*Chaetodon austriacus*):
456 current status and its characteristics. *Molecular Ecology Resources*. Online Early,
457 doi:10.1111/1755-0998.12588.
458
459 DiBattista JD, Roberts M, Bouwmeester J, Bowen BW, Coker DF, Lozano-Cortés DF, Choat
460 JH, Gaither MR, Hobbs JP, Kahil M, Kochzius M, Myers R, Paulay G, Robitzsch V, Saenz-
461 Agudelo P, Salas E, Sinclair-Taylor TH, Toonen RJ, Westneat M, Williams S, Berumen ML
462 (2016b) A review of contemporary patterns of endemism for shallow water reef fauna in the
463 Red Sea. *Journal of Biogeography*, **43**, 423-439.
464
465 DiBattista JD, Choat JH, Gaither MR, Hobbs JP, Lozano-Cortés DF, Myers R, Paulay G,
466 Rocha LA, Toonen RJ, Westneat M, Berumen ML (2016c) On the origin of endemic species
467 in the Red Sea. *Journal of Biogeography*, **43**, 13-30.
468
469 DiBattista JD, Rocha LA, Hobbs JP, He S, Priest MA, Sinclair-Taylor TH, Bowen BW,
470 Berumen ML (2015) When biogeographical provinces collide: hybridization of reef fishes at
471 the crossroads of marine biogeographical provinces in the Arabian Sea. *Journal of*
472 *Biogeography*, **42**, 1601-1614.
473
474 Drummond AJ, Ashton B, Cheung M, Heled J, Kearse M, Moir R, Stones-Havas S, Thierer
475 T, Wilson A (2009) Geneious v4.8. Available at: <http://www.geneious.com/>
476

477 Gaither MR, Bernal MA, Coleman RR, Bowen BW, Jones SA, Simison WB, Rocha LA
478 (2015) Genomic signatures of geographic isolation and natural selection in coral reef fishes.
479 *Molecular Ecology*, **24**, 1543-1557.
480

481 Gautier M, Gharbi K, Cezard T, Foucaud J, Kerdelhué C, Pudlo P, Cornuet JM, Estoup A
482 (2013) The effect of RAD allele dropout on the estimation of genetic variation within and
483 between populations. *Molecular Ecology*, **22**, 3165-3178.
484

485 Hand BK, Hether TD, Kovach RP, Muhlfeld CC, Amish SJ, Boyer MC, O'Rourke SM,
486 Miller MR, Lowe WH, Hohenlohe PA, Luikart G (2015) Genomics and introgression:
487 Discovery and mapping of thousands of species-diagnostic SNPs using RAD sequencing.
488 *Current Zoology*, **61**, 146-154.
489

490 Henning F, Lee HJ, Franchini P, Meyer A (2014) Genetic mapping of horizontal stripes in
491 Lake Victoria cichlid fishes: benefits and pitfalls of using RAD markers for dense linkage
492 mapping. *Molecular Ecology*, **23**, 5224-5240.
493

494 Hohenlohe PA, Bassham S, Etter PD, Stiffler N, Johnson EA, Cresko WA (2010) Population
495 genomics of parallel adaptation in threespine stickleback using sequenced RAD tags. *PLoS*
496 *Genetics*, **6**, e1000862.
497

498 Jackson AM, Semmens BX, de Mitcheson YS, Nemeth RS, Heppell SA, Bush PG, Aguilar-
499 Perera A, Claydon JAB, Calosso MC, Sealey KS, Schärer MT, Bernardi G (2014) Population
500 structure and phylogeography in Nassau grouper (*Epinephelus striatus*), a mass-aggregating
501 marine fish. *PLoS ONE*, **9**, e97508.
502

503 Krumgalz BS (1982) Calcium distribution in the world ocean waters. *Oceanologica Acta*, **5**,
504 121-128.
505

506 Kosuri S, Church, GM (2014) Large-scale *de novo* DNA synthesis: technologies and
507 applications. *Nature Methods*, **11**, 499-507.
508

509 Langmead B, Trapnell C, Pop M, Salzberg SL (2009) Ultrafast and memory-efficient
510 alignment of short DNA sequences to the human genome. *Genome Biology*, **10**, R25.
511

512 Lawton RJ, Pratchett MS, Delbeek JC (2013) Harvesting of butterflyfishes for aquarium and
513 artisanal fisheries. In: *Biology of Butterflyfishes* (eds: Pratchett MS, Berumen ML, Kapoor
514 BG), 269-291.
515

516 Liew YJ, Aranda M, Voolstra CR (2016) Reefgenomics.Org – a repository for marine
517 genomics data. *Database*, 1-4.
518

518 Malvezzi AJ, Murray CS, Feldheim KA, DiBattista JD, Garant D, Gobler CJ, Chapman DD,
519 Baumann H (2015) A quantitative genetic approach to assess the evolutionary potential of a
520 coastal marine fish to ocean acidification. *Evolutionary Applications*, **8**, 352-362.
521

521 Nakamura Y, Mori K, Saitoh K, Oshima K, Mekuchi M, Sugaya T, et al. (2013) Evolutionary
522 changes of multiple visual pigment genes in the complete genome of Pacific bluefin tuna.
523 *Proceedings of the National Academy of Sciences USA*, **110**, 11061-11066.
524

525 Myers S, Bottolo L, Freeman C, McVean G, Donnelly P (2005) A fine-scale map of
526 recombination rates and hotspots across the human genome. *Science*, **310**, 321-324.

527
528 Pecoraro C, Babbucci M, Villamor A, Franch R, Papetti C, Leroy B, Ortega-Garcia S, Muir J,
529 Rooker J, Arcoha F, Murua H, Zudaire I, Chassot E, Bodin N, Tinti F, Bargelloni L, Cariani
530 A (2015) Methodological assessment of 2b-RAD genotyping technique for population
531 structure inferences in yellowfin tuna (*Thunnus albacares*). *Marine Genomics*, Online Early.
532
533 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012) Double digest RADseq: an
534 inexpensive method for *de novo* SNP discovery and genotyping in model and non-model
535 species. *PLoS ONE*, **7**, e37135.
536
537 Picq S, McMillan WO, Puebla O (2016) Population genomics of local adaptation versus
538 speciation in coral reef fishes (*Hypoplectrus* spp, Serranidae). *Ecology and Evolution*, **6**,
539 2109-2124.
540
541 Puebla O, Bermingham E, McMillan WO (2014) Genomic atolls of differentiation in coral
542 reef fishes (*Hypoplectrus* spp., Serranidae). *Molecular Ecology*, **23**, 5291-5303.
543
544 Raitos DE, Pradhan Y, Brewin RJW, Stenchikov G, Hoteit I (2013) Remote sensing the
545 phytoplankton seasonal succession of the Red Sea. *PloS ONE*, e64909.
546
547 Read TD, Petit III RA, Joseph SJ, Alam MT, Weil R, Ahmad M, Bhimani R, Vuong JS,
548 Haase CP, Webb DH, Dove AD (2015) Draft sequencing and assembly of the genome of the
549 world's largest fish, the whale shark: *Rhincodon typus* Smith 1828. *PeerJ PrePrints*, e1036.
550
551 Ree RH, Hipp AL (2015) Inferring phylogenetic history from restriction site associated DNA
552 (RADseq). *Next-generation Sequencing in Plant Systematics*, 181-204.
553
554 Roda F, Ambrose L, Walter GM, Liu HL, Schaul A, Lowe A, Pelsler PB, Prentis P, Rieseberg
555 LH, Ortiz-Barrientos D (2013) Genomic evidence for the parallel evolution of coastal forms
556 in the *Senecio lautus* complex. *Molecular Ecology*, **22**, 2941-2952.
557
558 Saenz-Agudelo P, DiBattista JD, Piatek MJ, Gaither MR, Harrison HB, Nanninga GB,
559 Berumen ML (2015) Seascape genetics along environmental gradients in the Arabian
560 Peninsula: insights from ddRAD sequencing of anemonefishes. *Molecular Ecology*, **24**,
561 6241-6255.
562
563 Stamatakis A (2006) RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses
564 with thousands of taxa and mixed models. *Bioinformatics*, **22**, 2688-2690.
565
566 Stockwell BL, Larson WA, Waples RK, Abesamis RA, Seeb LW, Carpenter KE (2016) The
567 application of genomics to inform conservation of a functionally important reef fish (*Scarus*
568 *niger*) in the Philippines. *Conservation Genetics*, **17**, 239-249.
569
570 Tariel J, Longo GC, Bernardi G (2016) Tempo and mode of speciation in *Holacanthus*
571 angelfishes based on RADseq markers. *Molecular Phylogenetics and Evolution*, **98**, 84-88.
572
573 van de Peer Y (2004) Tetraodon genome confirms *Takifugu* findings: most fish are ancient
574 polyploids. *Genome Biology*, **5**, 250.
575

576 Wabnitz C (2003) *From Ocean to Aquarium: The Global Trade in Marine Ornamental*
577 *Species* (No. 17) Cambridge, UNEP/Earthprint.
578
579 Wagner CE, Keller I, Wittwer S, Selz OM, Mwaiko S, Greuter L, Sivasundar A, Seehausen O
580 (2013) Genome-wide RAD sequence data provide unprecedented resolution of species
581 boundaries and relationships in the Lake Victoria cichlid adaptive radiation. *Molecular*
582 *Ecology*, **22**, 787-798.
583
584 Waldrop E, Hobbs JP, Randall JE, DiBattista JD, Rocha LA, Kosaki RK, Berumen ML,
585 Bowen BW (2016) Phylogeography, population structure and evolution of coral-feeding
586 butterflyfishes (Subgenus *Corallochaetodon*). *Journal of Biogeography*, **43**, 1116-1129.
587
588 Willette DA, Allendorf FW, Barber PH, Barshis DJ, Carpenter KE, Crandall ED, Cresko,
589 WA, Fernandez-Silva, I, Matz, MV, Meyer, E, Santos, MD, Seeb LW, Seeb JE (2014) So,
590 you want to use next-generation sequencing in marine systems? Insight from the Pan-Pacific
591 Advanced Studies Institute. *Bulletin of Marine Science*, **90**, 79-122.
592

593 **Data Accessibility**

594 Raw RAD reads (FASTQ format) are available in the NCBI repository, BioProject:
595 PRJNA292048.
596 Genome browser URL: <http://caus.reefgenomics.org>
597 Filtered RAD reads in .vcf format are available from Dryad: doi:10.5061/dryad.f09rh.

598

599 **Author contributions**

600 J.D.D., P.S.-S., and M.L.B. conceived of and designed the RAD-Seq study. J.D.D., P.S.-S.,
601 and M.J.P. produced and analyzed the RAD libraries. X.W. and M.A. performed GO
602 enrichment analyses. All authors developed the manuscript and approve of the final paper.

603

604

Table 1. Stacks results of ddRAD data for nine species of Red Sea and Arabian Sea resident butterflyfish (and an outgroup surgeonfish) before and after quality filtering using *denovo_map.pl* and *ref_map.pl* options in Stacks. In all cases, 12 individuals were sampled per population. Population filtering parameters included: 1) minor allele frequency > 0.05, 2) the locus had to be genotyped in at least 80% of individuals and 3) the locus was present in all (or all but one) populations, which are represented by numbers outside and inside the parentheses, respectively, for “# of variable loci passing filter”.

Species	Sample size	Number of populations (geographic range of sampling)	denovo_map.pl				ref_map.pl			
			# of reads used	av. depth of coverage	# of loci	# of variable loci passing filter	# of reads used	av. depth of coverage	# of loci	# of variable loci passing filter
<i>Chaetodon austriacus</i> (exquisite butterflyfish)	84	7 (Gulf of Aqaba to South Farasan Banks)	106,207,151	13.7	289,504	8,842 (10,711)	91,785,299	13.4	194,235	8,812 (10,780)
<i>Chaetodon fasciatus</i> (Red Sea racoon butterflyfish)	96	8 (Gulf of Aqaba to Djibouti)	93,749,728	10.2	262,590	1,343 (2,650)	20,137,784	12.1	42,029	568 (697)
<i>Chaetodon larvatus</i> (hooded butterflyfish)	60	5 (Thuwal to Djibouti)	85,686,503	12.6	298,347	10,028 (12,393)	23,585,269	13.1	64,641	1,871 (2,179)
<i>Chaetodon melapterus</i> (Arabian butterflyfish)	72	6 (Djibouti to Muscat)	97,774,030	12.2	273,706	2,275 (4,384)	85,591,171	14.9	184,187	5,913 (7,761)
<i>Chaetodon mesoleucos</i> (white-face butterflyfish)	48	4 (Thuwal to Djibouti)	54,334,800	11.9	218,077	7,608 (11,151)	14,493,353	13.3	50,966	1,487 (1,806)
<i>Chaetodon paucifasciatus</i> (Eritrean butterflyfish)	72	6 (Gulf of Aqaba to South Farasan Banks)	100,273,743	16.2	363,221	10,257 (13,539)	18,726,021	16.1	44,073	955 (1,145)
<i>Chaetodon pictus</i> (horseshoe butterflyfish)	60	5 (Djibouti to Masirah Island)	46,939,818	10.9	225,976	2,073 (4,131)	9,907,802	12.8	40,368	552 (759)
<i>Chaetodon semilarvatus</i> (bluecheek butterflyfish)	96	8 (Jazirat Burqan to Djibouti)	89,373,543	12.3	189,554	1,270 (2,053)	19,177,775	13.9	37,766	239 (338)
<i>Chaetodon trifascialis</i> (chevron butterflyfish)	108	9 (Gulf of Aqaba to Masirah Island)	110,578,849	10.7	239,244	476 (1,271)	25,527,940	12.2	45,670	74 (331)
<u>Outgroup</u> <i>Ctenochaetus striatus</i> (striated surgeonfish)	120	10 (Gulf of Aqaba to Al Hallaniyats)	126,536,786	9.8	516,163	415 (1,508)	2,573,428	12.4	8,516	8 (27)

Abbreviation: av = average; SNP, single nucleotide polymorphism.

Fig. 1 Map of the Red Sea and Arabian Sea, including study species and collection sites.

Fig. 2 Top ten highest enriched ($P < 0.01$) gene ontologies for (a) *Chaetodon austriacus* and (b) *C. melapterus* based on annotated SNPs produced using a ddRAD protocol that have protein and GO information available. Gene ontology categories include molecular function, biological process, and cellular component.

Fig. 3 (a) Frequency of SNPs per scaffold for *Chaetodon austriacus* mapping to the assembled genome, (b) frequency distribution of scaffold size (in million base pairs), and (c) number of SNPs as a function of scaffold size using a ddRAD protocol and “ref_map.pl” option in Stacks. Mean = black dashed line; ± 2 standard deviations = green dashed lines; tolerance interval = red dashed lines.

Fig. 4 (a) Number of annotated versus non-annotated SNPs (total number of SNPs, i.e. annotated plus non-annotated, is in brackets) and (b) percentage of reads mapping to *Chaetodon austriacus* reference genome using Bowtie for Red Sea and Arabian Sea resident butterflyfish (and an outgroup surgeonfish). Mapping parameters used for Bowtie were $-n = 3$ and $-k = 1$. A SNP was considered annotated if located in any region of the scaffold identified as protein coding; a non-annotated SNP was any SNP that did not meet this criteria.

Fig. 1

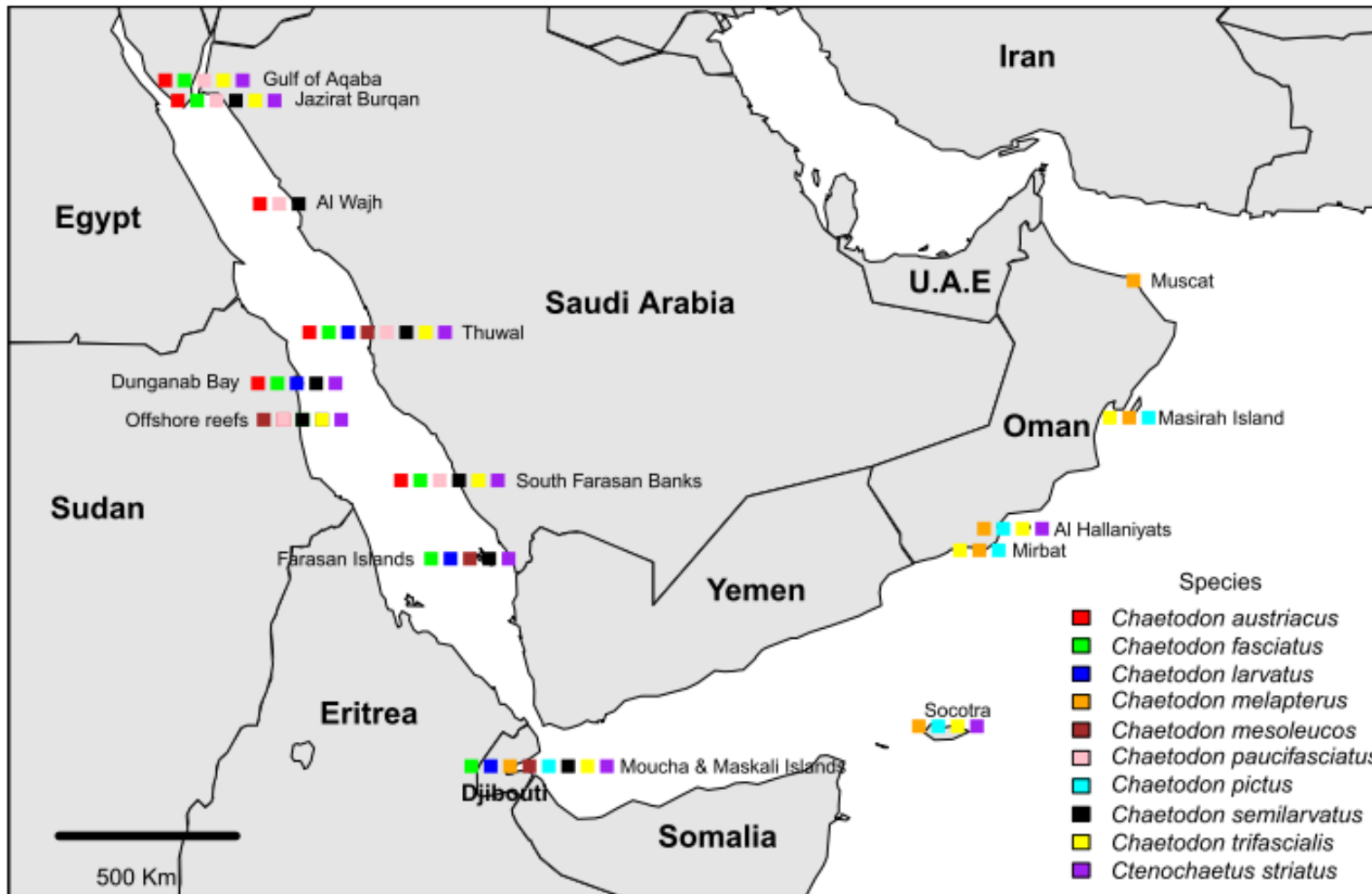


Fig. 2a

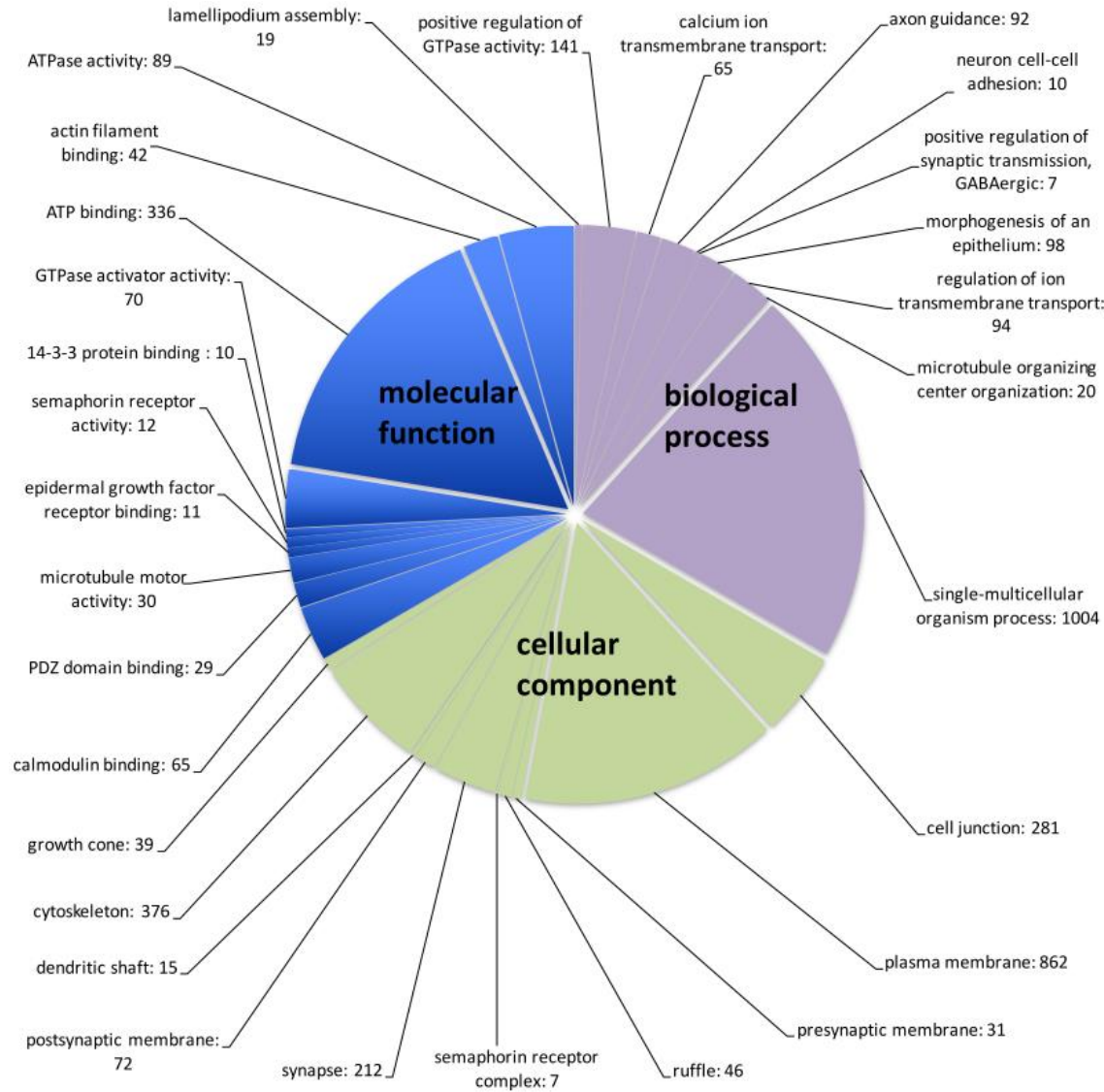


Fig. 2b

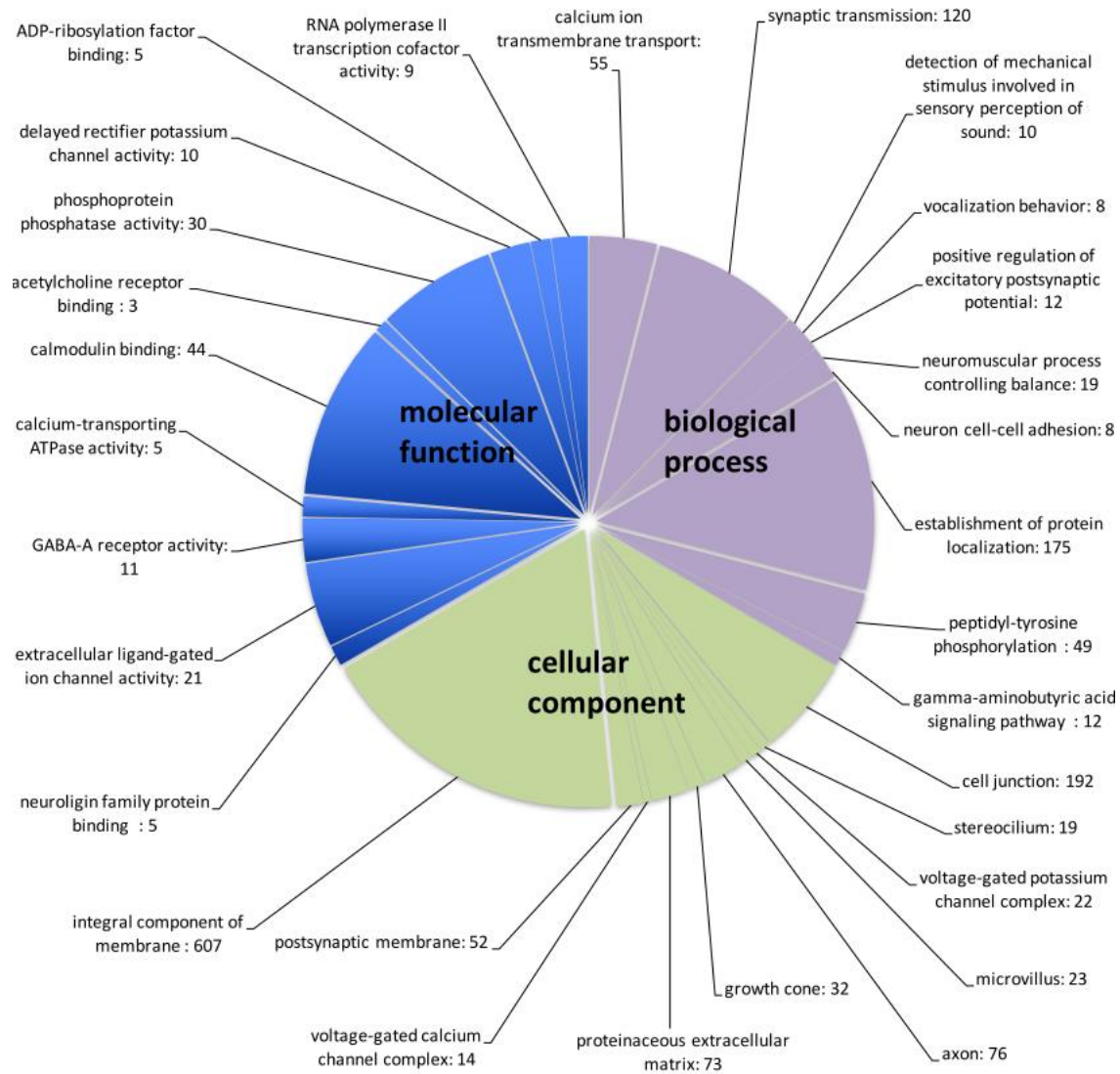


Fig. 3a

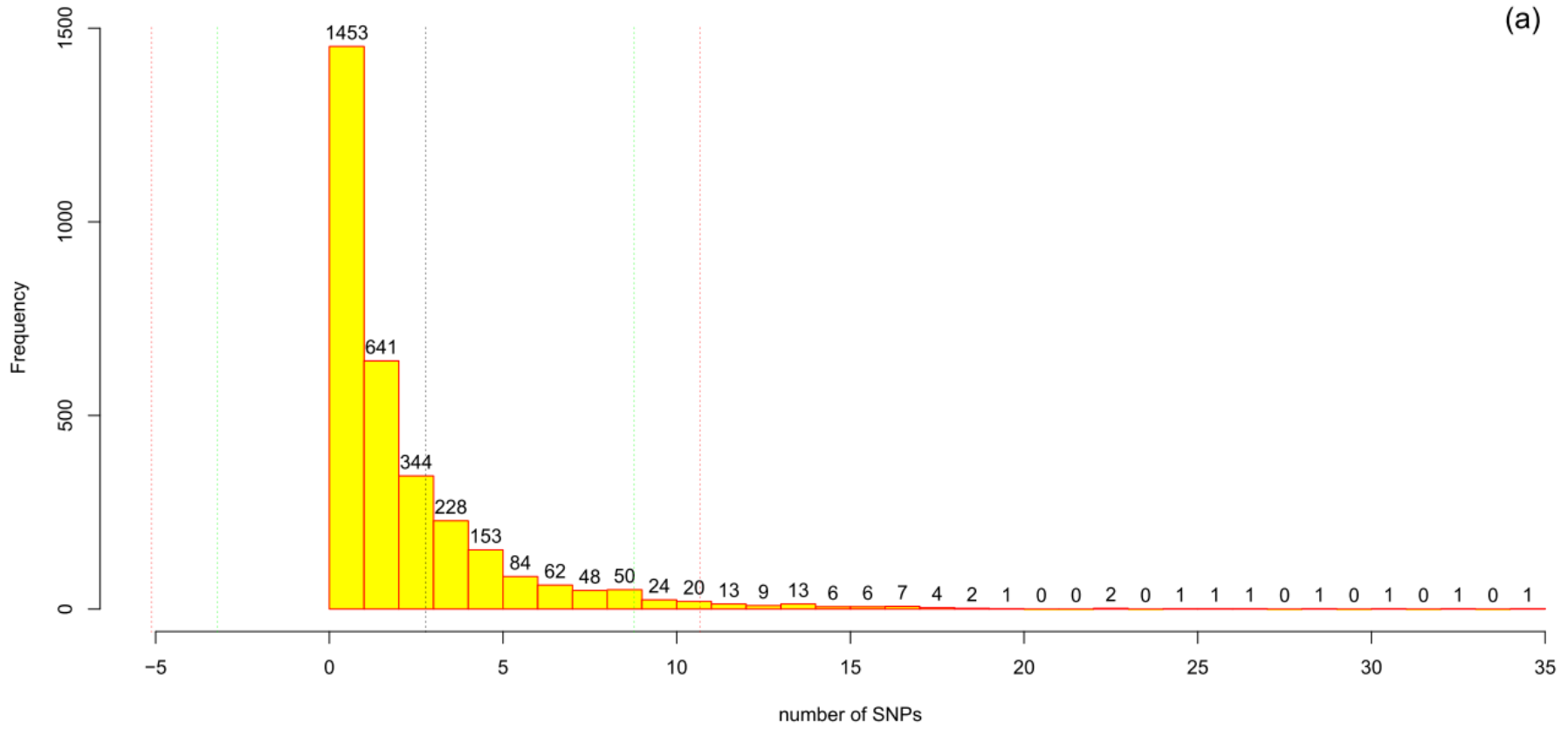


Fig. 3b

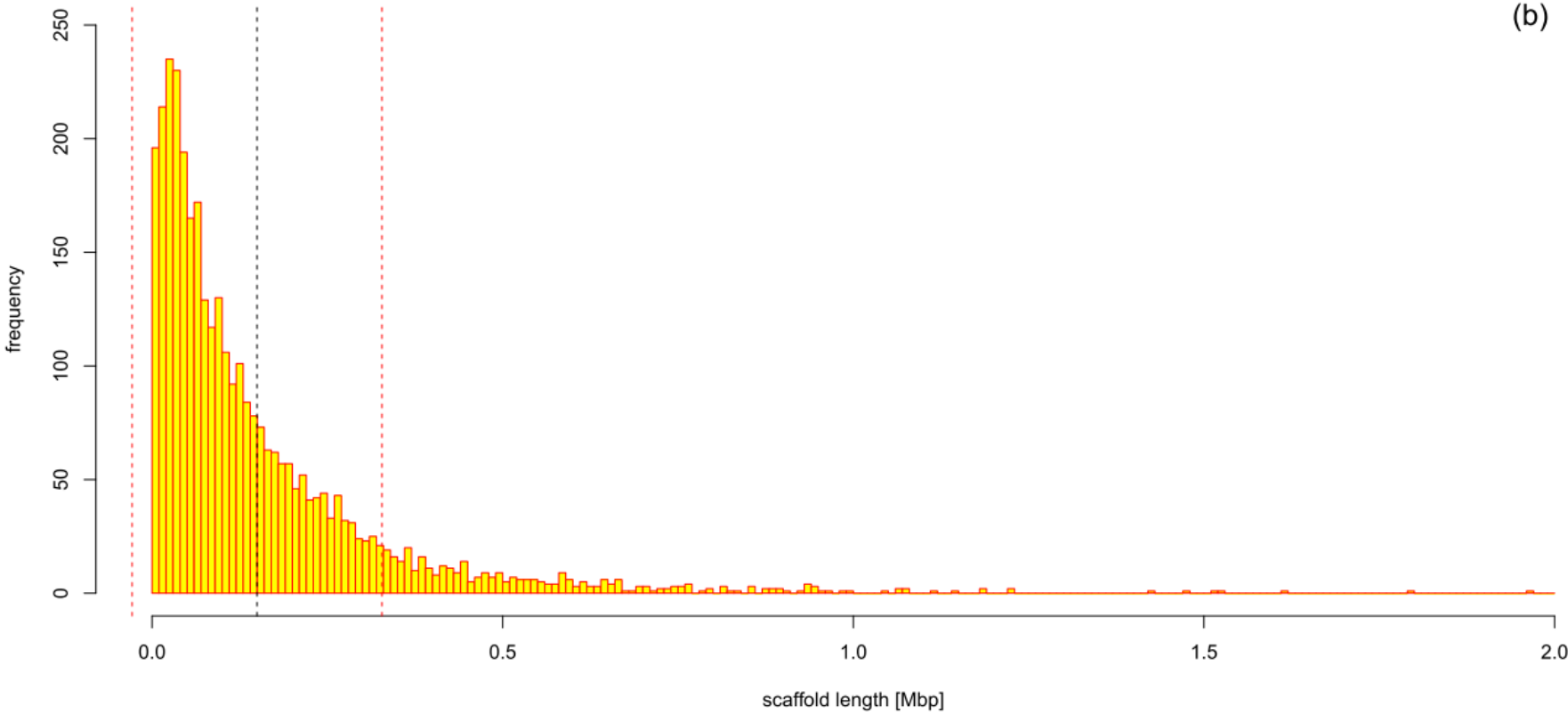


Fig. 3c

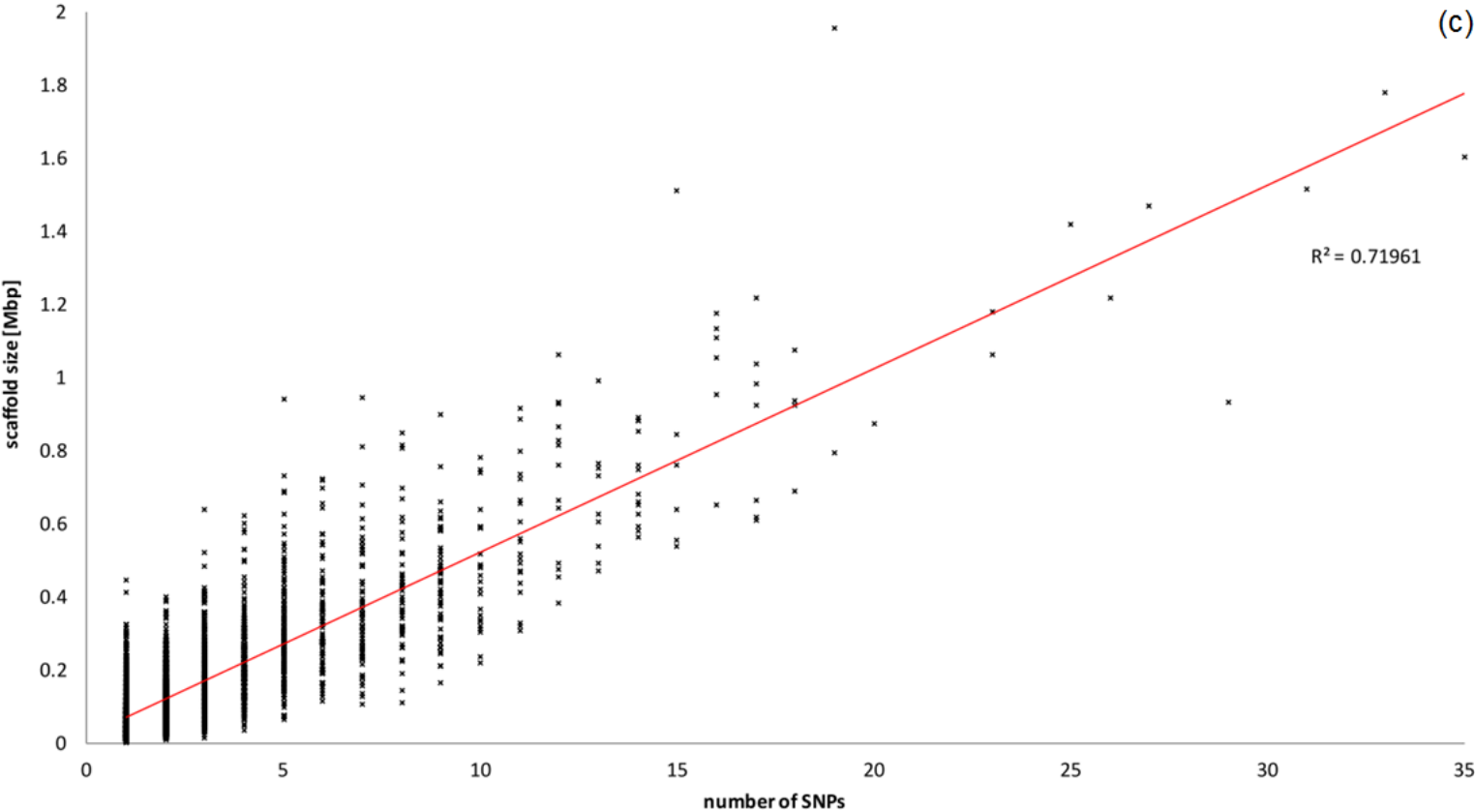


Fig. 4a

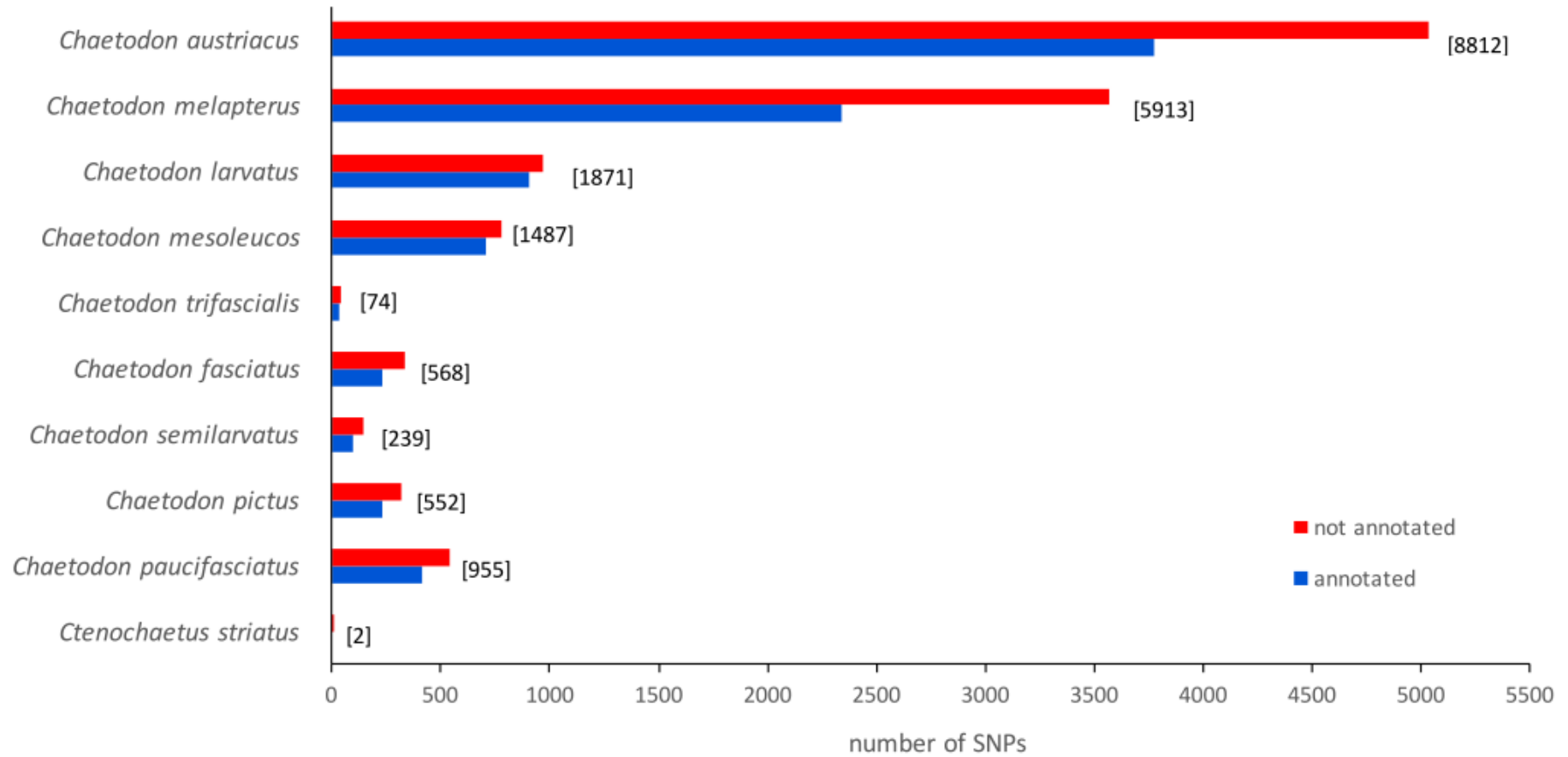
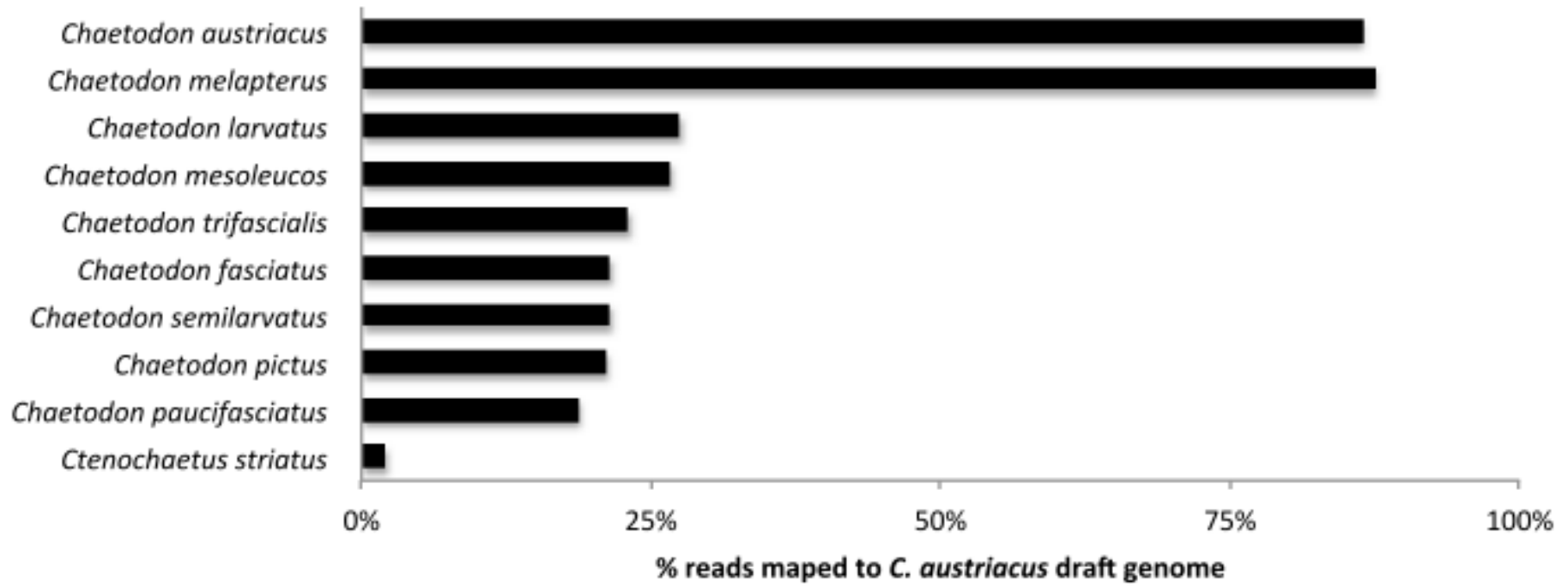


Fig. 4b



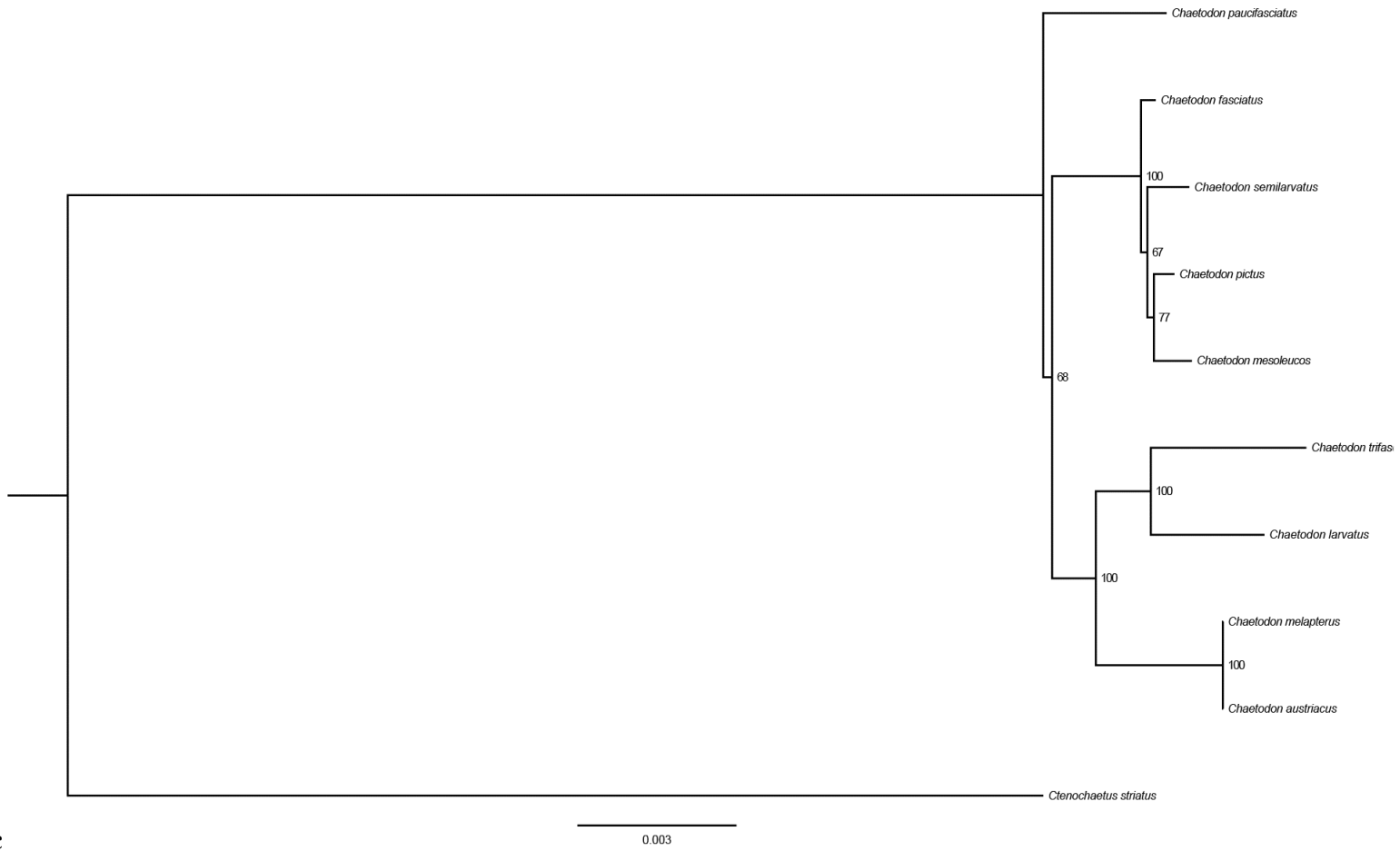


Fig. 4c

