

1 **Transcriptomics provides a robust framework for the relationships of the major**
2 **clades of cladobranch sea slugs (Mollusca, Gastropoda, Heterobranchia), but fails to**
3 **resolve the position of the enigmatic genus *Embletonia***

4

5 **Dario Karmeinski¹, Karen Meusemann^{1,2,3}, Jessica A. Goodheart⁴, Michael Schroedl^{5,6},**
6 **Alexander Martynov⁷, Tatiana Korshunova⁸, Heike Wägele¹, Alexander Donath^{1*}**

7

8 ¹ Zoological Research Museum Alexander Koenig, Centre for Molecular Biodiversity
9 Research, Adenauerallee 160, D-53113 Bonn, Germany

10 ² University of Freiburg, Institute of Biology I, Evolutionary Biology & Ecology, Hauptstr. 1, D-
11 79104 Freiburg (Brsg.), Germany

12 ³ Australian National Insect Collection, National Research Collections Australia, CSIRO
13 National Facilities and Collections, Clunies Ross Street, Acton, ACT 2601, Canberra,
14 Australia

15 ⁴ Scripps Institution of Oceanography, University of California, San Diego, La Jolla, 92037,
16 California, USA

17 ⁵ SNSB-Bavarian State Collection of Zoology, Münchhausenstr. 21, 81247 München,
18 Germany

19 ⁶ Ludwig Maximilians Universität München, GeoBioCenter LMU und Biozentrum,
20 Großhaderner Str. 2, 82152 Planegg-Martinsried, Germany

21 ⁷ Zoological Museum of the Moscow State University, Bolshaya Nikitskaya Str. 6, 125009
22 Moscow, Russia

23 ⁸ Koltzov Institute of Developmental Biology, Vavilova Str. 26, 119334 Moscow, Russia

24

25 *** Corresponding author: Alexander Donath (a.donath@leibniz-zfmk.de)**

26

27 **Abstract**

28 **Background:** Cladobranche sea slugs represent roughly half of the biodiversity of soft-
29 bodied, marine gastropod molluscs (Nudibranchia) on the planet. Despite their global
30 distribution from shallow waters to the deep sea, from tropical into polar seas, and their
31 important role in marine ecosystems and for humans (as bioindicators and providers of
32 medical drug leads), the evolutionary history of cladobranche sea slugs is not yet fully
33 understood. Here, we amplify the current knowledge on the phylogenetic relationships by
34 extending the cladobranche and outgroup taxon sampling using transcriptome data.

35 **Results:** We generated new transcriptome data for 19 species of cladobranche sea slugs and
36 two additional outgroup taxa. We complemented our taxon sampling with previously
37 published transcriptome data, resulting in a final supermatrix covering 56 species from all but
38 one accepted cladobranche superfamilies. Transcriptome assembly using six different
39 assemblers, selection of those assemblies providing the largest amount of potentially
40 phylogenetically informative sites, and quality-driven compilation of data sets resulted in
41 three different supermatrices: one with a full coverage of genes per species (446 single-copy
42 protein-coding genes) and two with a less stringent coverage (667 genes with 98.9%
43 partition coverage and 1,767 genes with 86% partition coverage, respectively). We used
44 these supermatrices to infer statistically robust maximum-likelihood trees. All analyses,
45 irrespective of the data set, indicate maximum statistical support for all major splits and
46 phylogenetic relationships on family level. The only discordance between the inferred trees
47 is the position of *Embletonia pulchra*. Extensive testing using Four-cluster Likelihood
48 Mapping, Approximately Unbiased tests, and Quartet Scores revealed that its position is not
49 due to any informative phylogenetic signal, but caused by confounding signal.

50 **Conclusions:** Our data matrices and the inferred trees inferred can serve as a solid
51 foundation for future work on the taxonomy and evolutionary history of Cladobranchia. The
52 correct placement of *E. pulchra*, however, proves challenging, even with large data sets.
53 Moreover, quartet mapping shows that confounding signal present in the data is sufficient to

54 explain the inferred position of *E. pulchra*, again leaving its phylogenetic position as an
55 enigma.

56 **Keywords**

57 Phylogenomics, Cladobranchia, RNA-Seq, Transcriptomes, Phylogeny, Embletoniidae

58 **Background**

59 Marine Heterobranchia (Gastropoda) have become a major focus as bioindicators to monitor
60 the health of coral reefs [1–7]. They mainly prey on a high variety of marine sessile
61 organisms, from algae to sponges, cnidarians, bryozoans and tunicates, and very often take
62 up the chemical compounds of the food for their own defence. These “stolen” compounds
63 have become of high interest for pharmacists in finding new drug leads for medical
64 applications [8–10]. However, they are also of high interest in understanding the evolution of
65 photosymbiosis and the role of “stolen” chloroplasts or even whole algal cells incorporated in
66 the slugs’ body, which help the slugs survive starving periods or otherwise increase fitness
67 [11–14]. Within marine Heterobranchia, the shell-less Nudibranchia have developed a
68 variety of biological strategies that make them unique within Metazoa. Of particular interest
69 is the sequestration of cnidocysts from the cnidarian prey, storing them in special
70 morphological structures (cnidosacs) in exposed body areas, and the ability to mature the
71 stolen cnidocysts (cleptocnides) in the cnidosac [15–18]. This unique defence system seems
72 to have evolved only in one of the major nudibranch clades, the Cladobranchia, within which
73 there are likely two independent origins [18].

74

75 Nudibranchia, with the two clades Cladobranchia and Anthobranchia, form a monophyletic
76 group that is well explained by morphological features [19]. Recently, the sister group
77 relationship to Pleurobranchomorpha (Pleurobranchida) as well as monophyly of
78 Nudibranchia was confirmed by transcriptomic data [20]. The monophyly of Nudibranchia
79 has also been confirmed in various molecular analyses using larger taxon sets, albeit small

80 gene sets (see review in [21]). However, few studies have used both morphological and
81 molecular methods to obtain and explain phylogenetic relationships within Cladobranchia. A
82 comprehensive study of Anthobranchia (Doridida) applying both molecular phylogenetic and
83 ontogenetic data was published recently [22]. Similar studies are still lacking for
84 Cladobranchia.

85

86 Pola and Gosliner [23] tried to resolve the phylogeny of Cladobranchia using one nuclear
87 and two mitochondrial genes: the study resulted in a topology that primarily consisted of an
88 unresolved comb. Bleidissel [24] analysed the Aeolidida within the Cladobranchia, based on
89 three genes (18S, 16S, and CO1), in order to investigate the evolution of the incorporation of
90 algae from the genus *Symbiodinium* in certain sea slugs. In this study, for the first time, the
91 paraphyly of the aeolidid family Facelinidae was shown. Similar to morphological data, the
92 success of retrieving more reliable relationships based on few molecular markers increases,
93 when working on family level. Recently, by the inclusion of the type species of the genus
94 *Facelina*, the “true” family Facelinidae was revealed and the name Myrrhinidae resurrected
95 for the second “facelinid” clade [25]. Korshunova and colleagues [26] studied the
96 relationships within the former Flabellinidae, including representatives of many Aeolidida.
97 The authors provided much evidence for the paraphyly of the former Flabellinidae, which
98 they then split into five different families.

99

100 Recent analyses, using a large transcriptomic data set, provided the first robust cladobranch
101 tree that enabled the study of evolution of food preferences [27, 28]. In a subsequent study,
102 a broader data set with nearly 90 taxa was used to examine the evolution of the cnidosac
103 [18], which is the main defence system of Aeolidida [29]. Similar defence structures have
104 evolved independently in *Hancockia* [15], a genus assigned to Dendronotida [18]. However,
105 the authors based their interpretations on a phylogenetic tree with partly low statistical
106 support. Moreover, a few taxa showed relatively long branches compared to other members
107 of the family (*Cerberilla*) or even the same genus (*Janolus*). Therefore, bias due to possible

108 long branch artefacts cannot be excluded. A reduced data set was used by Goodheart and
109 Wägele [30] to study the taxonomic relationship of an enigmatic pelagic cladobranch, the
110 genus *Phylliroe*, to analyse morphological traits enabling a shift from a benthic life style into
111 a pelagic form. With this study presented here, using an extended data set including 40
112 publicly available transcriptomes and combining them with 21 newly sequenced
113 transcriptomes, we provide robust support for yet unresolved relationships and reconsider
114 the phylogenetic position of the genus *Embletonia*, which has been assigned to various
115 groups in the past without any current consensus [16, 18, 24, 31, 32]. Robustly resolved and
116 reliably inferred phylogenetic trees that are not affected by confounding signal, but driven by
117 “true” phylogenetic signal, are one prerequisite for answering questions about the
118 evolutionary history of taxa and biological phenomena, such as the aforementioned evolution
119 of the cnidosac and photosymbiosis. Therefore, only trees that reflect most likely the “true”
120 history of species allow the inference of biological traits to understand biodiversity and its
121 origin. Inferred trees resulting from methodological or computational inadequacy can lead to
122 erroneous hypotheses (see, e.g., [33]). Taxa that diversified quickly and/or underwent rapid
123 radiation events within a short period of time are especially difficult to analyse (see, e.g., [34]
124 and several examples in [35]). Rapid radiation might also be the reason why for some
125 marine Heterobranchia it seems so difficult to reliably infer a species tree [21, 23, 36, 37].

126

127 In order to obtain a statistically highly supported tree and to check whether ambiguous splits
128 in this tree might be based on confounding and thus erroneous signal, we performed a
129 thorough study on 57 cladobranch and four outgroup transcriptomes. A comparison of the
130 results of various *de novo* transcriptome assemblers allowed us to specifically select those
131 assemblies that showed the highest sequence coverage with respect to a reference
132 orthologue set. After accounting for possible influences on phylogenetic inference, e.g.,
133 among-lineage heterogeneity and rejecting stationary, homogenous and time-reversible
134 conditions we compiled three final data sets including 56 out of originally 61 species
135 (discarding three species with a low coverage of the orthologue set as well as two species

136 due to model violation (see Methods): 1) the full data set allowing gene partitions to be
137 missing for single species, 2) a smaller intermediate data set in terms of the number of
138 genes, but with less missing data, and 3) a strict data set only including gene partitions for
139 which all species were present. In addition to careful preparation and processing of the data
140 throughout all steps of the analyses, i.e. evaluating the most appropriate assemblies,
141 identifying single-copy protein-coding orthologs, a thorough check of multiple sequence
142 alignments, and optimization and evaluation of the final data sets, we comprehensively
143 examined the ambiguously inferred position of *Embletonia* for alternative topologies with
144 approximately unbiased (AU) tests [38], Four-cluster Likelihood Mapping [39, 40], and
145 quartet puzzling [41] approaches.

146 **Results and Discussion**

147 ***Data preparation prior to phylogenetic analyses***

148 A list with details on the 21 species with newly sequenced transcriptome data is provided in
149 Supplementary Table S1, Additional File 2. Accession numbers for all species are given in
150 Supplementary Table S2, Additional File 2.

151

152 ***Transcriptome sequencing and data processing***

153 Paired-end sequencing resulted in approximately 7.5 Gbases of raw data per sample. For
154 the newly generated transcriptomes, the number of complete read pairs ranged from
155 20,266,817 in *Calmella cavolini* to 43,524,035 in *Facelina rubrovittata* with a median of
156 24,882,673 (*Hancockia cf. uncinata*). After trimming of possible adapter sequences and
157 sequence regions of low quality, the average read length of complete read pairs ranged from
158 118.1 bp in *Hermisenda emurai* to 139.6 bp in *Doto* sp. with a median of 133.8 bp in
159 *Polycera quadrilineata* (Supplementary Table S3, Additional File 2). Details on sequence
160 processing is provided in the Supplementary Text, Additional File 1. Transcriptome assembly
161 using six different *de novo* assemblers per data set resulted in a total number of 366

162 assemblies, i.e. six assemblies for each of the 61 transcriptomic data sets (see
163 Supplementary Text, Additional File 1 and Supplementary Table S4, Additional File 2).

164

165 *Evaluation of transcriptome assemblies, orthology prediction, and alignment procedures*

166 Evaluation of assembled transcriptomes and subsequently applying BUSCO version 3.0.0

167 [42] with the Metazoa set including 978 orthologs revealed a median of 731 (75%) complete

168 BUSCO genes per sample (maximum: 943 complete BUSCO genes, fragmented: 27,

169 missing: 8 in *Caloria elegans* assembled with BinPacker; minimum: 158 complete BUSCO

170 genes, fragmented: 123 missing: 697 in *Doris kerguelenensis* assembled with BinPacker).

171 All quality assessment results of the transcriptomes using BUSCO are summarised in

172 Supplementary Table S5, Additional File 2.

173 We additionally evaluated the quality of all transcriptomes separately for each assembly

174 method based on the results of orthology prediction and identified single-copy protein-coding

175 genes with our custom-made orthologue set comprising 1,992 orthologues (see Methods

176 and Supplementary Text, Additional File 1). Results were ranked based on the cumulative

177 length of transcripts that were successfully assigned to the reference genes used to identify

178 single-copy orthologues (OGs) in the transcriptomes (see Supplementary Table S6,

179 Additional File 2). The cumulative lengths ranged from 82,409 bp in *Pseudobornella*

180 *orientalis* (the genus was recently resurrected by Korshunova and colleagues [43]) (IDBA-

181 Tran, 458 genes successfully assigned) to 784,043 bp in *Caloria elegans* (Shannon, 1,904

182 genes successfully assigned). The median was 472,305 bp for the cumulative length and

183 1,577 for the number of successfully assigned genes. The best assembly (according to the

184 largest cumulative length) out of the six available per sample was selected as the

185 representative transcriptome for the respective species. This transcriptome was used for all

186 further downstream analyses and submitted to NCBI (see Supplementary Tables S2 and S7,

187 Additional File 2). In order to reduce the amount of missing data in subsequent analyses we

188 excluded three samples for which less than 60% of OGs included in the search had been

189 identified: *Pseudobornella orientalis*, *Dermatobranchus* sp., and *Tritoniopsis frydis*.

190 Furthermore, we only kept OGs for which at least 50% of the investigated 58 species had a
191 positive hit. This resulted in 1,767 OGs that we subsequently used to generate multiple
192 sequence alignments (MSAs) on amino acid level. Checking the MSAs for outlier sequences
193 (i.e. putatively misaligned or misassigned amino acid sequences), we identified 897
194 sequences in 112 MSAs that were subsequently removed. Outliers were found in sequences
195 from all remaining 58 species with the highest number of 30 outlier sequences in
196 *Limenandra confusa* and the lowest number of eight outlier sequences in *Doris*
197 *keruelensis* (median: 15 outliers, all details are provided in the Supplementary Text,
198 Additional File 1 and Supplementary Table S8, Additional File 2).

199 Alignment masking resulted in masking of alignment sites in 1,519 out of 1,767 genes
200 (Supplementary Text, Additional File 1) leaving ~ 71% of aligned unmasked sites for
201 subsequent analyses.

202

203 *Compilation, evaluation and optimization of data sets*

204 Analysing the concatenated supermatrix using MARE v. 1.2-rc [44], AliStat v. 1.6 [45] for
205 information content and data coverage, and SymTest v. 2.0.47 [46] for putative violating of
206 stationary, (time-)reversible and homogenous (SRH) model conditions [47, 48] using the
207 implemented Bowker's matched pairs test of symmetry [49] led to the results shown in
208 Supplementary Figures S1 and S2, Additional File 3.

209 With respect to the amount and distribution of missing data we initially compiled two data
210 sets as described in the methods section. The data set allowing for the highest amount of
211 missing data, termed "original unreduced data set", was not further reduced after
212 concatenation and comprised 58 species, 771,739 aligned amino acid positions and 1,767
213 gene partitions. The second data set with a full gene coverage for all 58 species (termed
214 "original reduced data set") comprised 143,859 aligned amino acid positions and 364 gene
215 partitions. Analysing both data sets for violation of SRH model conditions with SymTest
216 revealed that two species strongly violated the SRH conditions: *Calmella cavolini* and *Doris*
217 *keruelensis* (Supplementary Figure S2, Additional File 3). The latter transcriptome, which

218 likely belongs to an *Architectonia* species (personal communication Vanessa Knutson), was
219 probably mislabeled in the repository from which it was downloaded. Therefore, the
220 sequences belonging to these two species were removed entirely from all MSAs from the
221 original unreduced data set. This newly created data set (termed “unreduced data set”)
222 spanned a superalignment length of 771,706 amino acid positions including 1,767 gene
223 partitions.

224 To reduce the amount of missing data, we compiled an “intermediate” data set featuring only
225 those gene partitions for which at least one representative of the selected taxa was present
226 (see Materials and methods, Supplementary Text, Additional File 1, and Supplementary
227 Table S9, Additional File 2). This data set (termed “intermediate data set”) spanned a
228 superalignment length of 271,732 amino acid positions and included 667 gene partitions.
229 The third and most strict data set with full gene coverage for each of the 56 species (termed
230 “strict data set”) had a superalignment length of 170,140 amino acid positions and included
231 446 gene partitions. Details on data matrix diagnostics are provided in the Supplementary
232 Text, Additional File 1, Supplementary Table S10, Additional File 2, and Supplementary
233 Figures S3-S7, Additional File 3.

234

235 ***Phylogenetic relationships of sea slug taxa***

236 All analyses irrespective of the data set indicate maximum statistical support for all major
237 splits and phylogenetic relationships on family level (Fig. 1, Supplementary Figures S8-S12,
238 Additional File 3). Notably, low statistical support was inferred with regard to the
239 phylogenetic position of the genus *Embletonia*. In the following, we discuss taxa
240 relationships using the names according to the latest changes [50] that are implemented in
241 World Register of Marine Species [51, 52], although we disagree with several assignments
242 as discussed below.

243

244 *Phylogenetic relationships of major taxa and sea slug families*

245 Out of the seven accepted superfamilies of Cladobranchia, we were able to include
246 members of six superfamilies, whereas a representative of the rare Doridoxoidea was not
247 available to us. We inferred Aeolidida, Aeolidioidea (sensu WoRMS), Proctonotoidea, and
248 Dendronotoidea, with representatives of various families and genera, as being monophyletic.
249 This was fully supported by the quartet scores [41] for Aeolidida, Aeolidioidea, and
250 Proctonotoidea, and strongly supported for Dendronotoidea (see QuartetSampling scores,
251 splits 1-3 and 8 in Fig. 1 and Supplementary Table S11, Additional File 2). Arminoidea and
252 Tritonioidea are only represented by one genus each. Therefore, their assumed monophyly
253 still has to be tested by including relevant genera like *Doridomorpha* in Arminoidea, or
254 *Tochuina* in Tritonioidea.

255 Our analyses revealed the following ambiguities: *Flabellina affinis* (Flabellinidae), which is
256 currently regarded as a representative of Fionoidea [18], is inferred as sister taxon to
257 Aeolidioidea with maximal statistical support. Quartet sampling, on the other hand, showed
258 only medium support (split 4 in Fig. 1, Supplementary Table S11, Additional File 2) with the
259 large majority of quartets (67%) supporting the focal branch (Aeolidioidea + *Flabellina*
260 *affinis*), but the strong skew in discordance (quartet differential (QD) = 0) indicating the
261 possibility of a single different evolutionary history supported by all remaining quartets.

262 The family Flabellinopsidae is currently listed as a member of the Aeolidioidea in WoRMS
263 [52] with *Flabellinopsis iodinea* (Flabellinopsidae) being sister to all remaining taxa in this
264 large clade, confirming previous results [18, 26–28]. Again, this position is statistically
265 maximally supported by classic support values in our study and quartet puzzling scores
266 confirmed this position (split 5, Fig. 1) with strong support (94% of the non-uncertain
267 quartets). Although a strong skew in discordance (QD = 0) indicates the possible presence
268 of an alternative quartet relationship, this result is rather less meaningful due to the low
269 number of discordant trees (5% of the non-uncertain quartets). Thus, our results on
270 Flabellinidae and Flabellinopsidae partly contradict recent analyses and subsequent
271 systematic assignments.

272

273 Within Aeolidioidea, the families Myrrhinidae and Aeolidiidae form a monophyletic sister
274 group relationship in our study, thus confirming the results of [28] and [18]. This is also
275 consistent with recent morphological and molecular analyses [53].

276 The majority of the family Facelinidae is inferred as being monophyletic, but the facelinid
277 species *Noumeaella rubrofasciata* groups with Myrrhinidae in published analyses [18, 28] as
278 well as in our study with nearly maximal 'classical' statistical support. However, quartet
279 puzzling only shows weak support for this relationship (38% of the non-uncertain quartets;
280 see split 6 in Fig. 1 and Supplementary Table S11, Additional File 2). In fact, the quartet
281 frequencies show no clear signal since all three quartet topologies are roughly equally
282 supported (27% of the non-uncertain quartets support the second possible quartet topology,
283 36% support the third; QD = 0.85). Thus, the assignment of this species to Facelinidae [50]
284 or Myrrhinidae (our results) should be reconsidered in future studies. Interestingly, this
285 species did not cluster with other *Noumeaella* species in a three-gene analysis of Aeolidida
286 by Schillo and colleagues [37].

287 Fionoidea in the sense of Bouchet and colleagues [50] is paraphyletic, mainly due to the
288 position of *Flabellina affinis* and *Embletonia*, the latter is discussed below.

289

290 Within Fionoidea, the family Trinchesiidae represented here with three genera, is
291 monophyletic. Unidentiidae is sister to all remaining taxa within Fionoidea. Previously,
292 Korshunova and colleagues [26] inferred this family as sister taxon to Facelinidae and
293 Aeolidiidae. Quartet puzzling analyses, however, do not unambiguously support the
294 relationship of the Unidentiidae as sister to all other Fionoidea. There is rather weak support
295 (52% of the non-uncertain quartets) for said topology and the support for the other two
296 possible quartet topologies is almost similar (QD = 0.99), which indicates that no alternative
297 history is favoured (see split 7 in Fig. 1 and Supplementary Table S11, Additional File 2). In
298 this context, the results of Goodheart and colleagues [18] are quite noteworthy, because in
299 their study, Unidentiidae is the sister taxon of *Embletonia* and the clade *Embletonia* +
300 Unidentiidae is sister to all remaining Fionoidea. Results by Martynov and colleagues [53]

301 suggest a sister group relationship to other aeolidacean families, which is incompatible with
302 our results (see below).

303

304 The family Samlidae, represented by *Luisella babai*, is considered as being part of Fionoidea
305 [18]. In our study, however, it is inferred as sister to all remaining Aeolidida in all analyses
306 with maximum 'classical' statistical support as well as very strong quartet support (see split 8
307 in Fig. 1 and Supplementary Table S11, Additional File 2): About 98% of the quartets
308 supported this relationship, without evidence for alternative quartet topologies (QD = 1),
309 confirming previous results by Korshunova and colleagues [26].

310

311 With regard to Proctonotoidea, Tritonioidea, and Dendronotoidea, our results confirm the
312 findings published by Goodheart and colleagues [18] with the family Embletoniidae being the
313 only exception, as we will discuss below.

314

315 ***The phylogenetic position of Embletoniidae remains ambiguous***

316 The monogeneric family Embletoniidae, which currently only comprises two recognized
317 species, *Embletonia pulchra* and *E. gracilis*, has experienced a vivid history since the first
318 description of the genus *Embletonia* by Alder and Hancock [54], with *Pterochilus pulcher*
319 Alder and Hancock, 1844 as type species. The authors considered this species as a link
320 between cladobranch aeolids and panpulmonate sacoglossans, two taxa that are not closely
321 related to each other, but show many convergent characters. Pruvot-Fol [31], who named
322 the family for the first time, included members of Trinchesiidae, but assigned the whole clade
323 as a "section" to the dendronotoid family Dotidae. The two recognized members of
324 *Embletonia* share some characters with members of Fionoidea or Aeolidioidea, e.g., the
325 reduction of the lateral teeth, the absence of rhinophoral sheaths [56], and the presence of a
326 cnidosac at the end of the cerata, a synapomorphy of Aeolidida [19], which additionally
327 favours a position within this clade. However, Martin and colleagues [16] and Goodheart and

328 colleagues [18] have shown that this cnidosac differs to a great extent from the typical
329 aeolidid cnidosac by lacking a proper sac-like structure with musculature around it, as well
330 as a connection to the digestive gland, which is necessary for taking up sequestered
331 cnidocysts. Nevertheless, cnidocysts were found in the structures investigated by Goodheart
332 and colleagues [18]. The authors explain this atypical situation with a loss of characters or as
333 constituting a transitional form in the evolution of the cnidosac. Most recently, Martynov and
334 colleagues [53] provided evidence for paedomorphic processes, which would explain a
335 regressive evolution within Embletoniidae. This phenomenon is quite common in various
336 unrelated taxa inhabiting soft-bottom interstitial environments. *Embletonia* feeds on
337 hydrozoans, which is a typical food source of many aeolidids, but also of some
338 dendronotoidea. Unique to this genus are the cerata, which show bi- to quadrifid apices.
339 Highly structured cerata are not known from any aeolidids. However, the digestive gland
340 reaches far into these cerata, a character less pronounced in Proctonotoidea, and only
341 present in one further non-aeolidid group, the genus *Hancockia*.

342 *Embletonia* also shares traits that are characteristic for non-aeolidid groups, a reason why
343 Pruvot-Fol [31] included the genus into the family Dotidae (Dendronotoidea). This
344 assignment to Dotidae, as well as grouping with Trinchesiidae was, however, rejected later
345 by Schmekel [32], and the closer relationship to Dendronotoidea was emphasized by Miller
346 and Willan [57]. The primary connecting character is the lack of oral tentacles, which are
347 considered to be a synapomorphy of the Aeolidida [19]. Furthermore, their oral gland ducts
348 do not open into the oral tube by two separate ducts, but fuse into one common duct, which
349 is described for Proctonotoidea. Proctonotoidea mainly feed on bryozoans, however, a few
350 members also rely on hydrozoan prey, similar to *Embletonia*.

351

352 Few studies addressed the phylogenetic relationship of Embletoniidae using molecular data
353 [18, 24, 53]. Bleidissel [24] focussed on Aeolidida and included *Embletonia*, because of its
354 putative assignment to this group. Bleidissel's analyses, based on three genes, inferred a
355 sister group relationship of Embletoniidae with Notaeolidiidae, with the latter again being

356 sister to all remaining Aeolidida. In the only study based on a large data set, *Embletonia* was
357 inferred, along with *Unidentia*, within Aeolidida as sister to the remaining Fionioidea, thus
358 excluding a closer relationship with *Notaeolidia* [18]. Martin and colleagues [16] included
359 characters of the cnidosac into the morphological data matrix published by Wägele and
360 Willan [19], and their analysis resulted in an assignment of *Embletonia* to Aeolidida (tree not
361 shown in the publication). Likewise, our unpublished morphological analyses render
362 *Embletonia* as a sister taxon to Aeolidida. However, it is more likely the lack of data that
363 constrains the position than apomorphic characters of high phylogenetic information.

364

365 In our analyses comprising the unreduced and strict data set, *Embletonia pulchra* is inferred
366 as sister to Proctonotoidea, but with negligible support in the strict data set (65 BS, 50.1
367 aLRT, 1 aBayes). When assuming that *Embletonia* is a sister taxon of the Proctonotoidea
368 (see split 9 in Fig. 1 and position i in Fig. 2) and taking into consideration the studies on the
369 evolution of prey preferences [28] and cnidocyst incorporation [18], we have to conclude that
370 (1) feeding on Hydrozoa is an old trait within Cladobranchia and has not changed in
371 *Embletonia* (in contrast to Proctonotoidea) and (2) the evolution of the cnidosac might have
372 started in the stemline of the clade Aeolidida/Proctonotoidea/Embletoniidae, with *Janolus*
373 and *Dirona* probably representing a condition where the ability to store cnidocysts was lost
374 due to a food switch to bryozoan prey. Both, an independent evolution of cnidosacs and
375 cnidocyst storage (in the genus *Hancockia*) as well as a loss or strong reduction of these
376 complex structures has occurred within Dendronotoidea [18].

377

378 In our results from the intermediate data set, *Embletonia* is a sister group to all remaining
379 Aeolidida, but with even less support (51 BS, 33.1 aLRT, 1 aBayes). Considering this
380 relationship as a possible evolutionary scenario (Fig. 2, position ii, results on the
381 intermediate data set) means that the evolution of the cnidosac would have had to start in
382 the stemline of Embletoniidae/Aeolidida, while the typical character of the Dendronotoidea,
383 the rhinophoral sheaths, had already been lost and oral tentacles had not yet evolved.

384 However, both discussed possibilities (see Fig. 2, positions i and ii) are neither supported
385 statistically by classical bootstrap values, nor by our quartet analyses: Frequencies of the
386 three possible quartet topologies are almost equal (33% vs. 35% vs. 31% of all non-
387 uncertain quartets, split 9 in Fig. 1 and Supplementary Table S11, Additional File 2), which
388 indicates a highly complex evolution or rapid radiation.

389 Morphological analyses of important characters, like the positions of the anus, jaws, and
390 radula also contradict both relationships discussed above with apomorphic features lacking
391 for both hypotheses [53]. Instead, Embletoniidae shows an uniserial radula with central teeth
392 more similar to various aeolidids.

393

394 ***Evaluation of alternative positions of Embletoniidae and possible confounding signal***

395 To gain more insights into one of the obtained positions of *Embletonia* and to investigate
396 alternative positions (see Fig. 2), further analyses were conducted. Note, that we consider
397 the strict data set as most reliable, since it has full gene coverage for all species, following
398 the rationale of Dell’Ampio and colleagues [58] and Misof and colleagues [40], who showed
399 that inferred positions with high statistical support can be simply due to non-phylogenetic
400 signal, e.g., the distribution of missing data. However, we also performed some of the
401 analyses on the intermediate data set.

402 We applied approximately unbiased (AU) tests [38] for alternative positions of *Embletonia* on
403 the intermediate and strict data set. An AU test always takes the complete tree topology into
404 account and not only single splits. Further, it does not test whether or not confounding signal
405 is inherent in the data set, e.g., due to non-randomly distributed data and/or among-lineage
406 heterogeneity violating SRH conditions. We therefore also applied Four-cluster Likelihood
407 Mapping (FcLM) [39] along with a permutation approach on the strict data set. By testing all
408 three possible quartet topologies around *Embletonia* we evaluated whether or not there was
409 an alternative signal. Further, we checked for any sign of confounding signal (see [40]). To
410 this end, we defined four groups (Supplementary Table S12, Additional File 2) considering

411 group 4 as outgroup. We performed separate analyses for two outgroup variations: first, with
412 19 species including Anthobranchia and Pleurobranchomorpha and second, only with the 15
413 remaining cladobranch species. We drew quartets on the original data set and on three
414 artificial data sets, from which any existing phylogenetic signal was removed in three
415 different ways (see Materials and methods, Supplementary Text, Additional File 1, and [40]):
416 (a) by destroying the phylogenetic signal but leaving the distribution of missing data and the
417 compositional heterogeneity, which can lead to violating SRH conditions, untouched; (b) by
418 leaving the distribution of missing data untouched but making the data set completely
419 homogenous (no SRH model violation possible), and (c) by randomizing the missing data
420 distribution and making the data set completely homogenous. For all details see
421 Supplementary Text, Additional File 1.

422

423 Interestingly, the results of the phylogenetic trees and the results of the FcLM
424 (Supplementary Table S13, Additional File 2) and AU tests (Supplementary Table S14,
425 Additional File 2) were quite contradicting:

426 (i) Although the ML trees of the unreduced and strict data sets suggest that *Embletonia* is
427 sister to Proctonotoidea and although the AU test was unable to reject this topology ($p >$
428 0.05), it received the lowest proportion of quartets ($< 20\%$) in the FcLM approach. Thus, this
429 relationship can only be explained by confounding signal (see original and permutation
430 results in Supplementary Table S13, Additional File 2).

431 (ii) Although the best ML tree of the intermediate data set suggests *Embletonia* to be sister
432 to all remaining Aeolidida, a position that is not rejected by the AU test ($p > 0.05$), the FcLM
433 results indicate only minimal support for such a relationship: the proportion of supporting
434 quartets, excluding those that can be explained by confounding signal, was only around 3%.
435 This also implies that AU tests, irrespective of whether or not a topology for the data set is
436 significantly rejected, cannot be used to check if the signal is confounding.

437 (iii) A sister group relationship of *Embletonia* to a clade Aeolidida + Proctonotoidea, which
438 received strongest support in the FcLM analyses (8-16% of all quartets after excluding the

439 proportion of supporting quartets that can be explained by confounding signal, see
440 Supplementary Table S13, Additional File 2), was equally rejected by the AU test.

441

442 There is only very little signal that is not confounding (around 3-8%, compare quartets of
443 original with permuted approaches, Supplementary Table S13, Additional File 2), which
444 would support either *Embletonia* + Aeolidida (position ii in Fig. 2) or *Embletonia* as sister to a
445 clade Aeolidida + Proctonotoidea (position iii in Fig. 2). Thus, these results clearly indicate
446 that the position of *Embletonia* as a sister taxon of Proctonotoidea is not due to any
447 informative phylogenetic signal, but only due to confounding signal in our data set, and again
448 leaves the phylogenetic position of *Embletonia* as an enigma.

449

450 In order to analyse further possibilities of putative relationships of *Embletonia*, we tested four
451 alternative positions (iv - vii, see Fig. 2) of *Embletonia*, which have been discussed in the
452 literature before, by applying the AU test on the strict data set (see Fig. 2 and see below).
453 Note that none of these positions were inferred in any of our ML analyses.

454 (iv) Since *Embletonia* exhibits characters, which are shared with the Dendronotoidea, we
455 analysed a putative sister group relationship with this superfamily.

456 (v) Although an assignment to Tritonioidea is very unlikely, because *Embletonia* does not
457 share all the characters special for this superfamily, the position of the Arminoidea is variable
458 within the various published phylogenies [18, 59, 60] when including this superfamily.
459 Nevertheless, we tested this possibility.

460 The last two tests imply a closer relationship of *Embletonia* with Fionoidea, a relationship
461 that was assumed in former times and reflects the current systematics [50]. Therefore, we
462 tested (vi) a position of *Embletonia* as sister to Fionoidea and (vii) *Embletonia* as sister to
463 Unidentiidae and this clade being again sister to the remaining Fionoidea in restricted sense
464 [18, 53].

465

466 AU tests significantly rejected ($p < 0.05$) all four alternative positions (iv - vii, see Fig. 2) of
467 *Embletonia* (see Supplementary Table S14, Additional File 2).

468

469 Despite our extensive molecular data sets and tests, we still cannot unambiguously assign
470 *Embletonia* to one of the superfamilies in our tree. Beyond only small putative phylogenetic
471 signal as indicated by our FcLM analyses, which is also in line with the negligible support
472 considering classical statistical support, a reason could be the lack of relevant taxa in our
473 data set that could positively influence the position of Embletoniidae in the cladobranch tree
474 (e.g., Doridomorpha, Charcotiidae, Notaeolidiidae). Interestingly, morphological traits are
475 also confounding and do not yet allow for an unambiguous assignment. Because of its
476 unresolved position, several evolutionary traits within the Cladobranchia cannot be
477 satisfactorily explained.

478

479 **Conclusions**

480 Due to the high number of orthologous single-copy genes that could be successfully
481 extracted from the transcriptomes, the high information content and up to full gene coverage
482 of the supermatrices, and the high resolution of all three phylogenies, we conclude that the
483 use of transcriptomic data is a valuable tool for analysing phylogenetic relationships within
484 Cladobranchia. Nevertheless, analyses of large data sets can be error-prone to systematic
485 bias and classical support values might be inflated as has been shown and discussed [61–
486 64]. Beyond careful data processing prior to phylogenetic tree inference, additional thorough
487 tests, e.g., AU tests, quartet approaches like FcLM and quartet puzzling as well as checks
488 for confounding signals on a variety of different data matrices become more and more
489 indispensable. Our study has revealed that, despite previous efforts, the position of some
490 families within this group, especially the Embletoniidae, requires further investigation and
491 possibly taxonomic revision. In future studies, the present data set should be extended by
492 increasing the number of group-specific orthologous single-copy genes and by including

493 Charcotiidae, Notoalidiidae and other relevant species to shed light on the relationships
494 between families and superfamilies in Cladobranchia in order to draw a more complete
495 image of the evolution of this enigmatic group.

496

497 **Materials and methods**

498 An overview of the complete workflow is displayed in Fig. 3. Major steps are described here
499 while all details and settings can be found in the Supplementary Text, Additional File 1.

500

501 *Taxon sampling and sampling of transcriptome data*

502 For this study, we used recently published transcriptome data and generated new
503 transcriptome data for 21 species. We collected 19 species of Cladobranchia and two more
504 distantly related species of heterobranch sea slugs from different locations in the
505 Mediterranean Sea and the Sea of Japan (Supplementary Table S1, Additional File 2). The
506 specimens were preserved in RNAlater (Qiagen) or IntactRNA (Evrogen) and stored at -80
507 °C. The specimens collected on Elba island (Supplementary Table S1, Additional File 2)
508 were stored at -20 °C for approximately two weeks and then transferred to -80 °C until RNA
509 extraction. RNA extraction was performed using the Macherey & Nagel NucleoSpin RNA II
510 kit. Preparation and amplification of the cDNA libraries were performed by StarSeq GmbH,
511 Mainz using the Illumina TruSeq Stranded RNA HT kit. Paired-end sequencing was also
512 conducted at StarSeq with a read length of 150 base pairs on an Illumina NextSeq 500
513 sequencing platform. Raw reads were submitted to the NCBI SRA database. All accession
514 numbers are provided in Supplementary Table S2, Additional File 2.

515 Our newly generated transcriptome samples were combined with the published
516 transcriptome data of another 40 samples that we downloaded from the NCBI SRA database
517 (Supplementary Table S2, Additional File 2) [27, 28, 65, 66]. The published data comprised
518 37 species of Cladobranchia as well as two dorids, *Prodoris clavigera* and *Doris*

519 *kerguelenensis*, and one pleurobranchid, *Pleurobranchaea californica* (Supplementary Table
520 S2, Additional File 2).

521

522 *De novo transcriptome assembly*

523 All raw sequence reads of published and newly generated samples were quality-checked
524 prior to and after adapter trimming using FastQC Version 0.11.5 [67]. Adapter trimming and
525 quality filtering were performed with Trimmomatic v0.36 [68] using a custom adapter file (see
526 Additional File 4).

527

528 Data from altogether 61 samples were assembled using six assembly tools: BinPacker v. 1.1
529 [69], IDBA-Tran v. 1.1.1 [70], Shannon v 0.0.2 [71], SOAPdenovo-Trans v. 1.04 [72], Trans-
530 ABySS v. 1.5.5 [73], and Trinity v. 2.4.0 [74]. All assemblers were run with default settings
531 and all paired-end reads that survived the trimming process were used as input. We
532 additionally provided surviving single-end reads to those assemblers that were capable of
533 processing them (IDBA-Tran, SOAPdenovo-Trans, and Trans-ABYSS).

534 Following identification of the best transcriptome assembly per species (see below), possible
535 foreign contaminants were identified upon submission of the newly sequenced
536 transcriptomes to NCBI Transcriptome Shotgun Assembly (TSA) database and subsequently
537 removed from the sequences. Details are provided in the Supplementary Text, Additional
538 File 1 and in Supplementary Table S7, Additional File 2. The five alternative assemblies for
539 each sample that has been sequenced in frame of this study are provided in Additional File
540 5.

541

542 *Orthology prediction and generation of data matrices*

543 We designed a custom-made orthologue set by selecting all genes that were listed by
544 OrthoDB version 9 [75] to be single-copy at the hierarchical level “Lophotrochozoa” and
545 downloaded the respective table with the IDs of the orthologue groups (called OGs
546 hereinafter). We additionally downloaded the official gene sets of three species with well-

547 sequenced and annotated genomes, which we selected as reference species (i.e.
548 *Biomphalaria glabrata*, Official Gene Set (OGS) version 1.2 vectorbase [76], *Crassostrea*
549 *gigas*, OGS version Sep-2012 (ENA genebuild) [77], and *Lottia gigantea*, OGS version Jan-
550 2013 (JGI genebuild) [78]. We excluded five genes from this set due to defective sequence
551 headers, leading to a custom-made orthologue set comprising 1,992 orthologues. Orthology
552 prediction was performed using Orthograph v.0.6.2 [79], for which we used the
553 aforementioned orthologue set (Additional File 6). Details are provided in the Supplementary
554 Text, Additional File 1. To reduce the amount of missing data per species, three
555 transcriptome assemblies that covered less than 60% of the orthologue set were excluded
556 from further analyses: *Pseudobornella orientalis* (53% of the orthologue set missing),
557 *Dermatobranchus* sp. (46% missing), and *Tritoniopsis frydis* (51% missing). We then
558 removed all OGs for which less than 50% of the investigated species had a positive hit. This
559 resulted in 1,767 OGs for further analyses.

560

561 The quality of all transcriptome assemblies was further assessed with BUSCO v3.0.0 using
562 the metazoa_odb9 reference set genes comprising 978 BUSCO groups [42] and default
563 settings (Supplementary Table S5, Additional File 2). Because BUSCO's general metazoa
564 data set is not very specific for nudibranchs and since there is no way to easily compile a
565 nudibranch-specific reference data set (R. Waterhouse, personal communication), we
566 devised a method that makes use of the output generated by Orthograph. For each
567 Orthograph run, we calculated the number of sequences that were assigned to OGs by
568 Orthograph as well as the cumulative length of these sequences. With the aim to maximize
569 the amount of data, the latter was used as a criterion to determine the best assembly for
570 each species (for details see Supplement Text, Additional File 1, Supplementary Table S6,
571 Additional File 2, and Additional File 7).

572

573 Multiple sequence alignments on translational level were generated using DIALIGN-TX
574 Version 1.0.2 [80] and checked for outlier sequences using a newly implemented version of

575 the outlier script described in [40] (see Supplementary Text, Additional File 1 for details;
576 unfiltered alignments are provided in Additional File 8). Sequences identified as outliers as
577 well as all sequences belonging to the three reference taxa were removed from the
578 alignments (Additional File 9).

579

580 The amino acid multiple sequence alignments were examined with the program Aliscore
581 version 2.0 [81, 82] in order to identify ambiguous or randomly similar aligned sites. All
582 positions flagged by Aliscore (~ 29% of the originally aligned sites, see Supplementary Text,
583 Additional File 1) were discarded using AliCut version 2.31 [83] (Additional File 10). The
584 resulting masked amino acid alignments were concatenated into a supermatrix along with
585 the creation of a partition file using FASconCAT-G version 1.04 [84].

586

587 *Compilation, evaluation and optimization of data sets*

588 This amino acid supermatrix, with 58 species and including 1,767 genes, was analysed
589 using the software tool MARE version 1.2-rc [44] in order to assess the potential information
590 content (IC) of each gene partition, the overall information content of the matrix, and the
591 coverage in terms of gene partitions. The tool AliStat version 1.6 [45] was used to calculate
592 alignment diagnostics and the software SymTest version 2.0.47 [46–48] was used to analyse
593 the compositional heterogeneity of the supermatrix in order to detect possible violations of
594 stationary, (time-)reversible, and homogeneous (SRH) conditions [49].

595

596 To reduce especially among-lineage heterogeneity (see Results and Discussion), we
597 excluded the species *Doris kerguelenensis* and *Calmella cavolini* from our data (see
598 Supplementary Figures S2 and S7, Additional File 3).

599

600 We repeated analyses with MARE, AliStat, and SymTest and compiled three final data sets,
601 allowing different levels of missing data (Supplementary Table S10, Additional File 2): an
602 unreduced data set with 56 species and all 1,767 gene partitions with 771,707 aligned

603 amino-acid sites and allowing ~ 39% missing data; an intermediate data set in which data for
604 at least one representative of the defined groups (Supplementary Table S9, Additional File
605 2) had to be present, which led to a data matrix of 56 species and 667 gene partitions
606 (271,732 aligned sites) with 98% gene coverage and 18% of missing data, and our most
607 strict data set only including genes present in all 56 species. This led to a data matrix with
608 170,140 aligned sites, 446 gene partitions and less than 13% of missing data. Missing data
609 can lead to confounding signals in phylogenetic inference [40, 44, 58]. We therefore consider
610 our strict data set as most reliable. Details are provided in the Supplementary Text,
611 Additional File 1. The three supermatrices are provided in Additional File 11.

612

613 *Phylogenetic tree inference*

614 For all three data sets, maximum likelihood (ML) trees were calculated using IQ-TREE
615 version 1.6.12 [85]. The best fitting amino acid models for each partition were identified
616 using ModelFinder [86], which was run using an edge-link partitioned approach [87]. Out of
617 20 tree searches per data set, we selected the best ML tree according to the best log-
618 likelihood. Statistical support was derived from non-parametric bootstrap replicates ensuring
619 bootstrap convergence. Additionally, we calculated SH-like approximate likelihood ratio test
620 support [88] and approximate Bayes test support [89]. The best ML tree of each of the three
621 data sets was tested for the presence of rogue taxa using RogueNaRok v.1.0 [90]. Details
622 for each step including used settings are provided in the Supplementary Text, Additional File
623 1.

624

625 *Testing for alternative topologies*

626 *Quartet puzzling*

627 To analyse phylogenetic discordance, we applied the Quartet Sampling (QS) method [41],
628 which aims to identify the lack of branch support due to low phylogenetic information,
629 discordance due to lineage sorting or introgression, and misplaced or erroneous taxa (rogue

630 taxa). Details on the analysis and interpretation of scores are provided in the Supplementary
631 Text, Additional File 1 and Supplementary Table S11, Additional File 2.

632

633 Testing the position of *Embletonia*

634 Since the inferred position of *Embletonia pulchra* was not stable comparing the best ML
635 trees of the intermediate and strict data set, we tested various possible topologies with AU
636 tests (see Fig. 2) [38] as implemented in IQ-TREE version 1.6.12 (see Results and
637 Discussion, Supplementary Text, Additional File 1, and Additional File 12). To further
638 analyse whether or not the placement of *Embletonia* in our best tree inferred from the strict
639 data set was influenced by confounding signal and violating SRH conditions, and whether or
640 not there was putative phylogenetic signal for alternative positions of *Embletonia*, we
641 additionally performed Four-cluster Likelihood Mapping (FcLM), which is outlined in the
642 results section and in detail in the Supplement Text, Supplementary File 1 (see also
643 Additional File 13). In summary, we tested the following seven alternative hypotheses
644 concerning the position of *Embletonia*:

645 i) *Embletonia* is sister to Proctonotoidea (AU test + FcLM)

646 ii) *Embletonia* is sister to all Aeolidida (AU test + FcLM)

647 iii) *Embletonia* is sister to (Aeolidida, Proctonotoidea) (AU test + FcLM)

648 iv) *Embletonia* is sister to Dendronotoidea (AU test)

649 v) *Embletonia* is sister to Arminoidea (AU test)

650 vi) *Embletonia* is sister to Fionoidea (AU test)

651 vii) *Embletonia* is sister to Unidentiidae and this clade is sister to remaining Fionoidea (AU
652 test).

653

654 **Abbreviations**

655 aLRT: approximate likelihood ratio test, AU test: approximately unbiased test, BS: bootstrap,

656 FcLM: Four-cluster Likelihood Mapping, IC: information content, ML: maximum likelihood,

657 MSA: multiple sequence alignment, OG: orthologue group, OGS: official gene sets, QD:
658 Quartet differential, QS: Quartet Sampling, SRH: stationary, (time-)reversible and
659 homogenous, TSA database: Transcriptome Shotgun Assembly database, WoRMS:

660 **Declarations**

661 *Ethics approval and consent to participate*

662 Not applicable

663 *Consent for publication*

664 Not applicable

665 *Availability of data and materials*

666 The data sets and scripts supporting the conclusions of this article are available via
667 Figshare, "[UNIQUE PERSISTENT IDENTIFIER AND WEB LINK TO DATA SET(S) WILL
668 BE PROVIDED UPON ACCEPTANCE OF THE ARTICLE]."

669 *Competing interests*

670 The authors declare no competing interests.

671 *Funding*

672 AD and HW received funding by the German Research Foundation (DFG DO 1781/1-1, DFG
673 WA 618/18-1). AM received support by the MSU Zoological Museum (AAAA-A16-
674 116021660077-3). The work of TK was conducted under the IDB RAS Government basic
675 research program in 2020 No. 0108-2019-0002.

676 *Authors' contributions*

677 AD, DK, and HW designed the study. DK, HW, MS, AM, and TK collected and provided
678 material. DK, KM, AD, and HW performed all data analyses. DK and AD developed scripts.
679 AD performed sequence data management. MS, AM, and TK provided pictures. JG provided

680 access to unpublished data. All authors contributed in writing the manuscript, with AD, DK,
681 KM, and HW taking the lead. All authors read and approved the final manuscript.

682 *Acknowledgements*

683 KM thanks Ondrej Hlinka (CSIRO) for HPC support and Thomas Wong (ANU) for kindly
684 providing Unique Tree. KM and DK thank Robert Waterhouse for details on official gene sets
685 used in OrthoDB 8. DK thanks Malte Petersen for suggestions on using Orthograph.

686

687 **References**

- 688 1. Nimbs MJ, Larkin M, Davis TR, Harasti D, Willan RC, Smith SDA. Southern range
689 extensions for twelve heterobranch sea slugs (Gastropoda: Heterobranchia) on the eastern
690 coast of Australia. *Mar Biodivers Rec.* 2016;9:27.
- 691 2. Eisenbarth J-H, Undap N, Papu A, Schillo D, Dialao J, Reumschüssel S, et al. Marine
692 Heterobranchia (Gastropoda, Mollusca) in Bunaken National Park, North Sulawesi,
693 Indonesia - A follow-up diversity study. *Diversity.* 2018;10:127.
- 694 3. Kaligis F, Eisenbarth J-H, Schillo D, Dialao J, Schäberle TF, Böhringer N, et al. Second
695 survey of heterobranch sea slugs (Mollusca, Gastropoda, Heterobranchia) from Bunaken
696 National Park, North Sulawesi, Indonesia - how much do we know after 12 years? *Mar*
697 *Biodivers Rec.* 2018;11:2.
- 698 4. Nimbs M, Smith S. Beyond Capricornia: Tropical sea slugs (Gastropoda, Heterobranchia)
699 extend their distributions into the Tasman Sea. *Diversity.* 2018;10:99.
- 700 5. Ompi M, Lumoindong F, Undap N, Papu A, Wägele H. Monitoring marine Heterobranchia
701 in Lembeh Strait, North Sulawesi (Indonesia), in a changing environment. *AACL Bioflux.*
702 2019;12:664–77.
- 703 6. Undap N, Papu A, Schillo D, Ijong FG, Kaligis F, Lepar M, et al. First survey of
704 heterobranch sea slugs (Mollusca, Gastropoda) from the Island Sangihe, North Sulawesi,
705 Indonesia. *Diversity.* 2019;11:170.

- 706 7. Papu A, Undap N, Martinez NA, Segre MR, Datang IG, Kuada RR, et al. First Study on
707 Marine Heterobranchia (Gastropoda, Mollusca) in Bangka Archipelago, North Sulawesi,
708 Indonesia. *Diversity*. 2020;12:52.
- 709 8. Fisch K, Hertzner C, Böhringer N, Wuisan Z, Schillo D, Bara R, et al. The potential of
710 Indonesian heterobranchs found around Bunaken Island for the production of bioactive
711 compounds. *Mar Drugs*. 2017;15:384.
- 712 9. Gavagnin M, Carbone M, Ciavatta ML, Mollo E. Natural products from marine
713 heterobranchs: an overview of recent results. *Chem J Mold*. 2019;14:9–31.
- 714 10. Avila C. Terpenoids in marine heterobranch molluscs. *Mar Drugs*. 2020;18:162.
- 715 11. Burghardt I, Wägele H. The symbiosis between the ‘solar-powered’ nudibranch *Melibe*
716 *engeli* Risbec, 1937 (Dendronotoidea) and *Symbiodinium* sp. (Dinophyceae). *J Molluscan*
717 *Stud*. 2014;80:508–17.
- 718 12. Wägele H, Martin WF. Endosymbioses in sacoglossan sea slugs: Plastid-bearing animals
719 that keep photosynthetic organelles without borrowing genes. In: Löffelhardt W, editor.
720 *Endosymbiosis*. Vienna: Springer Vienna; 2014. p. 291–324.
721 http://link.springer.com/10.1007/978-3-7091-1303-5_14. Accessed 21 Sep 2016.
- 722 13. Melo Clavijo J, Donath A, Serôdio J, Christa G. Polymorphic adaptations in metazoans to
723 establish and maintain photosymbioses: Evolution of photosymbiosis. *Biol Rev*.
724 2018;93:2006–20.
- 725 14. Laetz EMJ, Wägele H. Comparing amylose production in two solar-powered sea slugs:
726 the sister taxa *Elysia timida* and *E. cornigera* (Heterobranchia: Sacoglossa). *J Molluscan*
727 *Stud*. 2019;85:166–71.
- 728 15. Martin R, Heß M, Schrödl M, Tomaschko K-H. Cnidosac morphology in dendronotacean
729 and aeolidacean nudibranch molluscs: from expulsion of nematocysts to use in defense?
730 *Mar Biol*. 2009;156:261–8.
- 731 16. Martin R, Tomaschko K-H, Heß M, Schrödl M. Cnidosac-related structures in *Embletonia*
732 (Mollusca, Nudibranchia) compared with dendronotacean and aeolidacean species. *Open*
733 *Mar Biol J*. 2010;4:96–100.

- 734 17. Obermann D, Bickmeyer U, Wägele H. Incorporated nematocysts in *Aeolidiella*
735 *stephanieae* (Gastropoda, Opisthobranchia, Aeolidioidea) mature by acidification shown by
736 the pH sensitive fluorescing alkaloid Ageladine A. *Toxicon*. 2012;60:1108–16.
- 737 18. Goodheart JA, Bleidißel S, Schillo D, Strong EE, Ayres DL, Preisfeld A, et al.
738 Comparative morphology and evolution of the cnidosac in Cladobranchia (Gastropoda:
739 Heterobranchia: Nudibranchia). *Front Zool*. 2018;15:43.
- 740 19. Wägele H, Willan RC. Phylogeny of the Nudibranchia. *Zool J Linn Soc*. 2000;130:83–
741 181.
- 742 20. Pabst EA, Kocot KM. Phylogenomics confirms monophyly of Nudipleura (Gastropoda:
743 Heterobranchia). *J Molluscan Stud*. 2018;84:259–65.
- 744 21. Wägele H, Klusmann-Kolb A, Verbeek E, Schrödl M. Flashback and foreshadowing—a
745 review of the taxon Opisthobranchia. *Org Divers Evol*. 2014;14:133–49.
- 746 22. Korshunova T, Fletcher K, Picton B, Lundin K, Kashio S, Sanamyan N, et al. The
747 Emperor’s *Cadlina*, hidden diversity and gill cavity evolution: new insights for the taxonomy
748 and phylogeny of dorid nudibranchs (Mollusca: Gastropoda). *Zool J Linn Soc*.
749 2020;189:762–827.
- 750 23. Pola M, Gosliner TM. The first molecular phylogeny of cladobranchian opisthobranchs
751 (Mollusca, Gastropoda, Nudibranchia). *Mol Phylogenet Evol*. 2010;56:931–41.
- 752 24. Bleidissel S. Molekulare Untersuchungen zur Evolution der Aeolidida (Mollusca,
753 Gastropoda, Nudibranchia, Cladobranchia) und zur Evolution einer sekundären Symbiose
754 mit Symbiodinium (Dinoflagellata) in den Aeolidida. Dissertation. Bergische University of
755 Wuppertal; 2010. <https://d-nb.info/1012468550/34>.
- 756 25. Martynov A, Mehrotra R, Chavanich S, Nakano R, Kashio S, Lundin K, et al. The
757 extraordinary genus *Myja* is not a tergipedid, but related to the Facelinidae s. str. with the
758 addition of two new species from Japan (Mollusca, Nudibranchia). *ZooKeys*. 2019;818:89–
759 116.
- 760 26. Korshunova T, Martynov A, Bakken T, Evertsen J, Fletcher K, Mudianta IW, et al.
761 Polyphyly of the traditional family Flabellinidae affects a major group of Nudibranchia:

- 762 aeolidacean taxonomic reassessment with descriptions of several new families, genera, and
763 species (Mollusca, Gastropoda). *ZooKeys*. 2017;717:1–139.
- 764 27. Goodheart JA, Bazinet AL, Collins AG, Cummings MP. Relationships within
765 Cladobranchia (Gastropoda: Nudibranchia) based on RNA-Seq data: an initial investigation.
766 *R Soc Open Sci*. 2015;2:150196.
- 767 28. Goodheart JA, Bazinet AL, Valdés Á, Collins AG, Cummings MP. Prey preference
768 follows phylogeny: evolutionary dietary patterns within the marine gastropod group
769 Cladobranchia (Gastropoda: Heterobranchia: Nudibranchia). *BMC Evol Biol*. 2017;17:221.
- 770 29. Putz A, König GM, Wägele H. Defensive strategies of Cladobranchia (Gastropoda,
771 Opisthobranchia). *Nat Prod Rep*. 2010;27:1386–402.
- 772 30. Goodheart JA, Wägele H. Phylogenomic analysis and morphological data suggest left-
773 right swimming behavior evolved prior to the origin of the pelagic Phylliroidae (Gastropoda:
774 Nudibranchia). *Org Divers Evol*. 2020. doi:10.1007/s13127-020-00458-9.
- 775 31. Pruvot-Fol A. Faune De France n° 58, Mollusques Opisthobranches. Paris: Paul
776 Lechevalier; 1954. [https://www.abebooks.com/book-search/title/faune-de-france-58-](https://www.abebooks.com/book-search/title/faune-de-france-58-mollusques-opisthobranches/author/mme-alice-pruvot-fol/)
777 [mollusques-opisthobranches/author/mme-alice-pruvot-fol/](https://www.abebooks.com/book-search/title/faune-de-france-58-mollusques-opisthobranches/author/mme-alice-pruvot-fol/). Accessed 21 Aug 2020.
- 778 32. Schmekel L. Anatomie der Genitalorgane von Nudibranchiern (Gastropoda Euthyneura).
779 *Pubblicazioni Della Stazione Zool Napoli*. 1970;38:120–217.
- 780 33. Wägele JW, Letsch H, Klussmann-Kolb A, Mayer C, Misof B, Wägele H. Phylogenetic
781 support values are not necessarily informative: the case of the Serialia hypothesis (a mollusk
782 phylogeny). *Front Zool*. 2009;6:12.
- 783 34. Kück P, Wägele JW. Plesiomorphic character states cause systematic errors in
784 molecular phylogenetic analyses: a simulation study. *Cladistics*. 2016;32:461–78.
- 785 35. Wägele JW, Bartolomaeus T, editors. Deep Metazoan Phylogeny: The Backbone of the
786 Tree of Life: New insights from analyses of molecules, morphology, and theory of data
787 analysis. Berlin, Boston: DE GRUYTER; 2014. doi:10.1515/9783110277524.
- 788 36. Valdés Á, Lundsten L, Wilson NG. Five new deep-sea species of nudibranchs
789 (Gastropoda: Heterobranchia: Cladobranchia) from the Northeast Pacific. *Zootaxa*.

- 790 2018;4526:401.
- 791 37. Schillo D, Wipfler B, Undap N, Papu A, Böhringer N, Eisenbarth J-H, et al. Description of
792 a new *Moridilla* species from North Sulawesi, Indonesia (Mollusca: Nudibranchia:
793 Aeolidioidea)—based on MicroCT, histological and molecular analyses. *Zootaxa*.
794 2019;4652:265–95.
- 795 38. Shimodaira H. An approximately unbiased test of phylogenetic tree selection. *Syst Biol*.
796 2002;51:492–508.
- 797 39. Strimmer K, von Haeseler A. Likelihood-mapping: a simple method to visualize
798 phylogenetic content of a sequence alignment. *Proc Natl Acad Sci U S A*. 1997;94:6815–9.
- 799 40. Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, et al. Phylogenomics
800 resolves the timing and pattern of insect evolution. *Science*. 2014;346:763–7.
- 801 41. Pease JB, Brown JW, Walker JF, Hinchliff CE, Smith SA. Quartet Sampling distinguishes
802 lack of support from conflicting support in the green plant tree of life. *Am J Bot*.
803 2018;105:385–403.
- 804 42. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
805 assessing genome assembly and annotation completeness with single-copy orthologs.
806 *Bioinformatics*. 2015;31:3210–2.
- 807 43. Korshunova T, Bakken T, Grøtan VV, Johnson KB, Lundin K, Martynov A. A synoptic
808 review of the family Dendronotidae (Mollusca: Nudibranchia): a multilevel organismal
809 diversity approach. *Contrib Zool*. 2020;:1–61.
- 810 44. Misof B, Meyer B, von Reumont BM, Kück P, Misof K, Meusemann K. Selecting
811 informative subsets of sparse supermatrices increases the chance to find correct trees. *BMC*
812 *Bioinformatics*. 2013;14:348.
- 813 45. Wong TKF, Kalyanamoorthy S, Meusemann K, Yeates DK, Misof B, Jermiin LS. A
814 minimum reporting standard for multiple sequence alignments. *NAR Genomics Bioinforma*.
815 2020;2. doi:10.1093/nargab/lqaa024.
- 816 46. Jermiin LS, Ott M. SymTest. C++. 2017. <https://github.com/ottmi/symtest>. Accessed 28
817 May 2020.

- 818 47. Ho SYW, Jermiin LS. Tracing the decay of the historical signal in biological sequence
819 data. *Syst Biol.* 2004;53:623–37.
- 820 48. Ababneh F, Jermiin LS, Ma C, Robinson J. Matched-pairs tests of homogeneity with
821 applications to homologous nucleotide sequences. *Bioinformatics.* 2006;22:1225–31.
- 822 49. Bowker AH. A test for symmetry in contingency tables. *J Am Stat Assoc.* 1948;43:572–4.
- 823 50. Bouchet P, Rocroi J-P, Hausdorf B, Kaim A, Kano Y, Nützel A, et al. Revised
824 classification, nomenclator and typification of gastropod and monoplacophoran families.
825 *Malacologia.* 2017;61:1–526.
- 826 51. WoRMS Editorial Board. World Register of Marine Species. 2020. doi:10.14284/170.
- 827 52. MolluscaBase eds. MolluscaBase. MolluscaBase. 2020. Accessed at
828 <http://www.molluscabase.org> on 2020-08-25.
- 829 53. Martynov A, Lundin K, Picton B, Fletcher K, Malmberg K, Korshunova T. Multiple
830 paedomorphic lineages of soft-substrate burrowing invertebrates: parallels in the origin of
831 *Xenocratena* and *Xenoturbella*. *PLOS ONE.* 2020;15:1–24.
- 832 54. Alder J, Hancock A. A monograph of the British nudibranchiate Mollusca: with figures of
833 all the species. Part 5. London: The Ray Society; 1851.
- 834 55. Alder J, Hancock A. Descriptions of *Pterochilus*, a new genus of nudibranchiate
835 mollusca, and two new species of *Doris*. *Ann Mag Nat Hist.* 1844;14:329–31.
- 836 56. Edmunds M. Opisthobranchiate mollusca from Ghana: Flabellinidae, Piseinotecidae,
837 Eubranthidae & Embletoniidae. *J Conchol.* 2015;42:105–24.
- 838 57. Miller MC, Willan RC. Redescription of *Embletonia gracile* Risbec, 1928 (Nudibranchia:
839 Embletoniidae): relocation to suborder Dendronotacea with taxonomic and phylogenetic
840 implications. *J Molluscan Stud.* 1992;58:1–11.
- 841 58. Dell’Ampio E, Meusemann K, Szucsich NU, Peters RS, Meyer B, Borner J, et al.
842 Decisive data sets in phylogenomics: lessons from studies on the phylogenetic relationships
843 of primarily wingless insects. *Mol Biol Evol.* 2014;31:239–49.
- 844 59. Wägele H, Vonnemann V, Wägele J-W. Toward a phylogeny of the Opisthobranchia. In:
845 Lydeard C, Lindberg D, editors. *Molecular systematics and phylogeography of mollusks.*

- 846 Smithsonian Books; 2003. p. 185–228. [https://www.researchgate.net/profile/Heike_Waegele/](https://www.researchgate.net/profile/Heike_Waegele/publication/284477221_Toward_a_phylogeny_of_the_Opisthobranchia/links/566039f808ae1ef929857b4d.pdf)
847 [publication/284477221_Toward_a_phylogeny_of_the_Opisthobranchia/links/](https://www.researchgate.net/profile/Heike_Waegele/publication/284477221_Toward_a_phylogeny_of_the_Opisthobranchia/links/566039f808ae1ef929857b4d.pdf)
848 [566039f808ae1ef929857b4d.pdf](https://www.researchgate.net/profile/Heike_Waegele/publication/284477221_Toward_a_phylogeny_of_the_Opisthobranchia/links/566039f808ae1ef929857b4d.pdf). Accessed 18 Sep 2016.
- 849 60. Wägele H, Klussmann-Kolb A. Opisthobranchia (Mollusca, Gastropoda) – more than just
850 slimy slugs. Shell reduction and its implications on defence and foraging. *Front Zool.*
851 2005;2:3.
- 852 61. Simon S, Blanke A, Meusemann K. Reanalyzing the Palaeoptera problem – The origin of
853 insect flight remains obscure. *Arthropod Struct Dev.* 2018;47:328–38.
- 854 62. Vasilikopoulos A, Balke M, Beutel RG, Donath A, Podsiadlowski L, Pflug JM, et al.
855 Phylogenomics of the superfamily Dytiscoidea (Coleoptera: Adephaga) with an evaluation of
856 phylogenetic conflict and systematic error. *Mol Phylogenet Evol.* 2019;135:270–85.
- 857 63. Vasilikopoulos A, Misof B, Meusemann K, Lieberz D, Flouri T, Beutel RG, et al. An
858 integrative phylogenomic approach to elucidate the evolutionary history and divergence
859 times of Neuropterida (Insecta: Holometabola). *BMC Evol Biol.* 2020;20:64.
- 860 64. Szucsich NU, Bartel D, Blanke A, Böhm A, Donath A, Fukui M, et al. Four myriapod
861 relatives - but who are sisters? No end to debates on relationships among the four major
862 myriapod subgroups. *BMC Evol Biol.* 2020;accepted.
- 863 65. Zapata F, Wilson NG, Howison M, Andrade SCS, Jörger KM, Schrödl M, et al.
864 Phylogenomic analyses of deep gastropod relationships reject Orthogastropoda. *Proc R Soc*
865 *Lond B Biol Sci.* 2014;281:20141739.
- 866 66. Senatore A, Edirisinghe N, Katz PS. Deep mRNA sequencing of the *Tritonia diomedea*
867 brain transcriptome provides access to gene homologues for neuronal excitability, synaptic
868 transmission and peptidergic signalling. *PLOS ONE.* 2015;10:e0118321.
- 869 67. Andrews S. FastQC: a quality control tool for high throughput sequence data. 2010.
870 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>.
- 871 68. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
872 data. *Bioinformatics.* 2014;30:2114–20.
- 873 69. Liu J, Li G, Chang Z, Yu T, Liu B, McMullen R, et al. BinPacker: Packing-based *de novo*

- 874 transcriptome assembly from RNA-seq data. PLOS Comput Biol. 2016;12:e1004772.
- 875 70. Peng Y, Leung HCM, Yiu S-M, Lv M-J, Zhu X-G, Chin FYL. IDBA-tran: a more robust de
876 novo de Bruijn graph assembler for transcriptomes with uneven expression levels.
877 Bioinformatics. 2013;29:i326–34.
- 878 71. Kannan S, Hui J, Mazooji K, Pachter L, Tse D. Shannon: An information-optimal de novo
879 RNA-Seq assembler. bioRxiv. 2016;:039230.
- 880 72. Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, et al. SOAPdenovo-Trans: de novo
881 transcriptome assembly with short RNA-Seq reads. Bioinformatics. 2014;30:1660–6.
- 882 73. Robertson G, Schein J, Chiu R, Corbett R, Field M, Jackman SD, et al. De novo
883 assembly and analysis of RNA-seq data. Nat Methods. 2010;7:909–12.
- 884 74. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length
885 transcriptome assembly from RNA-Seq data without a reference genome. Nat Biotechnol.
886 2011;29:644–52.
- 887 75. Kriventseva EV, Tegenfeldt F, Petty TJ, Waterhouse RM, Simão FA, Pozdnyakov IA, et
888 al. OrthoDB v8: update of the hierarchical catalog of orthologs and the underlying free
889 software. Nucleic Acids Res. 2015;43:D250–6.
- 890 76. DeJong RJ, Emery AM, Adema CM. The mitochondrial genome of *Biomphalaria glabrata*
891 (Gastropoda: Basommatophora), intermediate host of *Schistosoma mansoni*. J Parasitol.
892 2004;90:991–7.
- 893 77. Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al. The oyster genome reveals stress
894 adaptation and complexity of shell formation. Nature. 2012;490:49–54.
- 895 78. Simakov O, Marletaz F, Cho S-J, Edsinger-Gonzales E, Havlak P, Hellsten U, et al.
896 Insights into bilaterian evolution from three spiralian genomes. Nature. 2013;493:526–31.
- 897 79. Petersen M, Meusemann K, Donath A, Dowling D, Liu S, Peters RS, et al. Orthograph: a
898 versatile tool for mapping coding nucleotide sequences to clusters of orthologous genes.
899 BMC Bioinformatics. 2017;18:111.
- 900 80. Subramanian AR, Kaufmann M, Morgenstern B. DIALIGN-TX: greedy and progressive
901 approaches for segment-based multiple sequence alignment. Algorithms Mol Biol. 2008;3:6.

- 902 81. Misof B, Misof K. A Monte Carlo approach successfully identifies randomness in multiple
903 sequence alignments: a more objective means of data exclusion. *Syst Biol.* 2009;58:21–34.
- 904 82. Kück P, Meusemann K, Dambach J, Thormann B, von Reumont BM, Wägele JW, et al.
905 Parametric and non-parametric masking of randomness in sequence alignments can be
906 improved and leads to better resolved trees. *Front Zool.* 2010;7:10.
- 907 83. Kück P. AliCUT. 2019. <https://github.com/PatrickKueck/AliCUT>.
- 908 84. Kück P, Longo GC. FASconCAT-G: extensive functions for multiple sequence alignment
909 preparations concerning phylogenetic studies. *Front Zool.* 2014;11:81.
- 910 85. Nguyen L-T, Schmidt HA, Haeseler A von, Minh BQ. IQ-TREE: A fast and effective
911 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.*
912 2015;32:268–74.
- 913 86. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. ModelFinder:
914 fast model selection for accurate phylogenetic estimates. *Nat Methods.* 2017;14:587–9.
- 915 87. Chernomor O, Haeseler A von, Minh BQ. Terrace aware data structure for phylogenomic
916 inference from supermatrices. *Syst Biol.* 2016;65:997–1008.
- 917 88. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms
918 and methods to estimate maximum-likelihood phylogenies: Assessing the performance of
919 PhyML 3.0. *Syst Biol.* 2010;59:307–21.
- 920 89. Anisimova M, Gil M, Dufayard J-F, Dessimoz C, Gascuel O. Survey of branch support
921 methods demonstrates accuracy, power, and robustness of fast likelihood-based
922 approximation schemes. *Syst Biol.* 2011;60:685–99.
- 923 90. Aberer AJ, Krompass D, Stamatakis A. Pruning rogue taxa improves phylogenetic
924 accuracy: an efficient algorithm and webservice. *Syst Biol.* 2013;62:162–6.

925

926 **Figures**

927 **Figure 1: Best ML tree (phylogram) from the strict data set.** Maximum likelihood (ML)
928 tree with bootstrap (BS) support values calculated on the strict data set. Black dots (●)

929 indicate a BS support value of 100. The numbers represent splits that are discussed in the
930 main text and the surrounding coloured circles represent Quartet Sampling (QS) scores for
931 the corresponding split. QFreq. = Quartet frequencies. QC = Quartet concordance. QD =
932 Quartet differential. QI = Quartet informativeness.

933 **Figure 2: Best ML tree (cladogram): AU tests + FcLM.** Cladogram with summarized
934 major families/clades and images of representative species. Splits for which additional
935 testing was performed are marked with Roman numerals (i-vii) in a coloured circle (AU test)
936 and a triangle (FcLM, splits i-iii). The original position of *E. pulchra* as obtained from the strict
937 data set is marked by a blue branch (T1). Alternative positions of *E. pulchra* are indicated by
938 a red (T2) and yellow branch (T3), respectively. We thank Craig A. Hoover for providing the
939 picture of *Flabellinopsis iodinea* and Karen Cheney for permissions to use the picture of
940 *Unidentia angelvaldesi*.

941 **Figure 3: Analysis workflow.** Schematic workflow representing all steps from NGS data to
942 the testing of alternative topologies with major steps being highlighted in shades of gray.

943

944 **Additional Files**

945 [Additional file 1](#) - Supplementary Text (pdf)

946

947 [Additional file 2](#) - Supplementary Tables S1 - S14 (xlsx)

948 **Table S1:** Sampling information for the species collected for this study.

949 **Table S2:** NCBI accession numbers for all species used in this study.

950 **Table S3:** Statistics of raw sequence reads before and after trimming.

951 **Table S4:** Assembly statistics.

952 **Table S5:** BUSCO results.

953 **Table S6:** Results of the Quality Checker script and selection of the best assembly.

954 **Table S7:** Information on sequences removed during contamination filtering.

955 **Table S8:** Number of removed outlier sequences per species.

956 **Table S9:** Group definitions to compile the intermediate data set.

957 **Table S10:** Supermatrix diagnostics of data sets used in this study.

958 **Table S11:** Results of the Quartet Sampling analysis.

959 **Table S12:** Group definitions used for Four-cluster Likelihood Mapping (FcLM) analyses.

960 **Table S13:** FcLM results testing the position of *Embletonia*.

961 **Table S14:** AU test results on the strict and intermediate data set.

962

963 [Additional file 3](#) - Supplementary Figures S1 - S12 (pdf)

964 **Figure S1: Species-pairwise site-coverage of the original unreduced and reduced data**
965 **sets.**

966 Heat maps indicate species-pairwise amino acid site-coverage of the sequences of 58
967 species in the original data sets inferred with AliStat. Low shared site-coverage is in shades
968 of red and high shared site-coverage is in shades of green. AliStat scores are given in
969 Supplementary Table S10, Additional File 2. **a)** original unreduced data set. **b)** original
970 reduced data set.

971

972 **Figure S2: Heat maps calculated with SymTest applying the Bowker's test on the**
973 **original unreduced and reduced data sets.**

974 Heat maps show the results of pairwise Bowker's test as implemented in SymTest 2.0.47
975 analysing the original data sets unreduced and reduced. The percentage of pairwise p-
976 values < 0.05 rejecting SRH conditions are given in parentheses. **a)** original unreduced data
977 set (p-values < 0.05: 83.36%). **b)** original reduced data set (p-values < 0.05: 42.65%). Note
978 that especially *Calmella* and *Doris* are obvious with respect to violating SRH conditions.

979

980 **Figure S3: Heat map visualising the information content of the final unreduced data**
981 **set calculated with MARE.**

982 The information content (IC) is colour-coded in shades of blue, with darker shades
983 representing higher IC and white squares indicating missing data. Red squares indicate

984 gene partitions with an IC = 0. Species are displayed in rows (x-axis) and gene partitions are
985 displayed in columns (y-axis). Supermatrix diagnostics of MARE are provided in
986 Supplementary Table S10, Additional File 2.

987

988 **Figure S4: Heat map visualising the information content of the final intermediate data**
989 **set calculated with MARE.**

990 The information content (IC) is colour-coded in shades of blue, with darker shades
991 representing higher IC and white squares indicating missing data. Red squares indicate
992 gene partitions with an IC = 0. Species are displayed in rows (x-axis) and gene partitions are
993 displayed in columns (y-axis). Supermatrix diagnostics of MARE are provided in
994 Supplementary Table S10, Additional File 2.

995

996 **Figure S5: Heat map visualising the information content of the final strict data set**
997 **calculated with MARE.**

998 The information content (IC) is colour-coded in shades of blue, with darker shades
999 representing higher IC and white squares indicating missing data. Red squares indicate
1000 gene partitions with an IC = 0. Species are displayed in rows (x-axis) and gene partitions are
1001 displayed in columns (y-axis). Supermatrix diagnostics of MARE are provided in
1002 Supplementary Table S10, Additional File 2.

1003

1004 **Figure S6: Species-pairwise site-coverage of the final unreduced, intermediate, and**
1005 **strict data set.**

1006 Heat maps indicate species-pairwise amino acid site-coverage of the sequences of 56
1007 species in the final data sets inferred with AliStat. Low shared site-coverage is in shades of
1008 red and high shared site-coverage is in shades of green. AliStat scores are given in
1009 Supplementary Table S10, Additional File 2. **a)** unreduced data set. **b)** intermediate data set.
1010 **c)** strict data set.

1011

1012 **Figure S7: Heat maps calculated with SymTest applying the Bowker's test on the final**
1013 **unreduced, intermediate, and strict data sets.**

1014 Heat maps show the results of pairwise Bowker's test as implemented in SymTest 2.0.47
1015 analysing the final data sets unreduced, intermediate, and strict. The percentage of pairwise
1016 p-values < 0.05 rejecting SRH conditions are given in parentheses. **a)** unreduced data set
1017 (p-values < 0.05: 82.14%). **b)** intermediate data set (p-values < 0.05: 63.96%). **c)** strict data
1018 set (p-values < 0.05: 46.17%).

1019

1020 **Figure S8: Best ML tree of the strict data set with aLRT and aBayes support.**

1021 The phylogram is identical to the phylogram in Fig. 1. The first value displays branch support
1022 based on 10,000 SH-aLRT replicates, the second value displays support derived from the
1023 approximate Bayesian support.

1024

1025 **Figure S9: Best ML tree of the intermediate data set with non-parametric bootstrap**
1026 **support.**

1027 Statistical support was inferred from 300 non-parametric bootstrap replicates.

1028

1029 **Figure S10: Best ML tree of the intermediate data set with aLRT and aBayes support.**

1030 The first value displays branch support based on 10,000 SH-aLRT replicates, the second
1031 value displays support derived from the approximate Bayesian support.

1032

1033 **Figure S11: Best ML tree of the unreduced data set with non-parametric bootstrap**
1034 **support.**

1035 Statistical support was inferred from 100 non-parametric bootstrap replicates.

1036

1037 **Figure S12: Best ML tree of the unreduced data set with aLRT and aBayes support.**

1038 The first value displays branch support based on 10,000 SH-aLRT replicates, the second
1039 value displays support derived from the approximate Bayesian support.

1040

1041 [Additional file 4](#) - FASTA file in zip archive

1042 **Archive S1:** Illumina adapters used for adapter trimming.

1043

1044 [Additional file 5](#) - FASTA files in zip archive

1045 **Archive S2:** Included in this archive are the five alternative assemblies for each sample that
1046 has been sequenced in the frame of this study (FASTA format). Note that the best selected
1047 assembly has been deposited at the NCBI TSA database. *Pseudobornella orientalis* (HW08)
1048 has been removed from the NCBI TSA database due to exceptionally low sequence quality.
1049 Its alternative assemblies are therefore also not part of this archive.

1050

1051 [Additional file 6](#) - FASTA/txt files in zip archive

1052 **Archive S3:** This archive includes official gene sets of the three reference species
1053 *Biomphalaria glabrata*, *Crassostrea gigas*, and *Lottia gigantea* on translational and
1054 transcriptional level, the list of all orthologous sequence clusters (OGs) as required for
1055 Orthograph, and an exemplary Orthograph config file.

1056

1057 [Additional file 7](#) – Python script/txt files in zip archive

1058 **Archive S4:** Included in this archive is the *Orthograph_Quality_Checker.py* script, a manual,
1059 an example configuration file, and an example output file.

1060

1061 [Additional file 8](#) – Alignment files in zip archive

1062 **Archive S5:** Unmasked multiple sequence alignments on amino acid level including *Doris*
1063 *kerгуelenensis* and *Calmella cavolini* prior to the removal of outliers.

1064

1065 [Additional file 9](#) – Python scripts in zip archive

1066 **Archive S6:** This archive contains two custom Python scripts. The *remove_outliers.py* script
1067 removes all identified outlier sequences from a given alignment. The

1068 *remove_reference_sequences.py* script removes all sequences from the reference species
1069 *Biomphalaria glabrata*, *Crassostrea gigas*, and *Lottia gigantea* from the alignments.

1070

1071 Additional file 10 – Alignment files in zip archive

1072 **Archive S7:** 1,767 Multiple sequence alignments (FASTA format) on amino acid level, from
1073 which sequences belonging to *Doris kerguelenensis* and *Calmella cavolini* as well as
1074 ambiguously aligned sections and gap-only sites were removed. These served as the basis
1075 for compiling the final unreduced supermatrix.

1076

1077 Additional file 11 – Alignment/txt files in zip archive

1078 **Archive S8:** The unreduced, intermediate, and strict supermatrix (FASTA format) plus
1079 respective gene partition information including the selected substitution model used in the
1080 phylogenetic analyses.

1081

1082 Additional file 12 – Tree files (NEWICK format) in zip archive

1083 **Archive S9:** Seven tree topologies (NEWICK format) displaying differing positions of
1084 *Embletonia pulchra* that were tested using the approximately unbiased (AU) test with IQ-
1085 TREE.

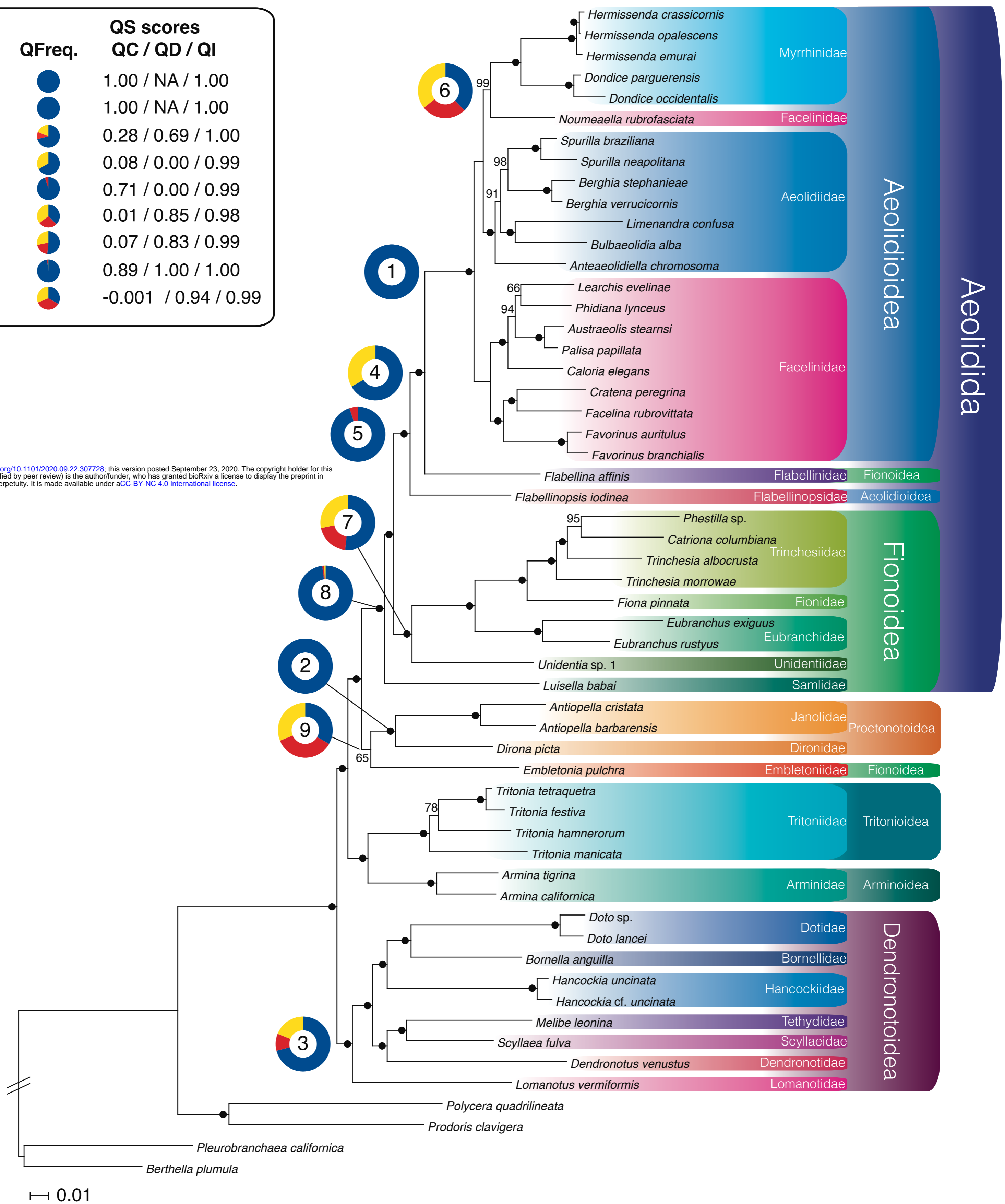
1086

1087 Additional file 13 - Alignment/NEXUS/txt files in zip archive

1088 **Archive S10:** Data used for Four-cluster Likelihood Mapping (FcLM). This archive includes
1089 two directories (one per approach, with a) 19 species included in Group 4 and b) 15 species
1090 included in Group 4; see section 17). Each directory includes four subdirectories: original,
1091 permutationI, permutationII, and permutationIII. In each subdirectory, the following files that
1092 served as input for the FcLM with IQ-TREE are provided: superalignment (FASTA format),
1093 partition file with gene boundaries and respective models, and the group file (NEXUS format)
1094 listing the species included in the defined groups (see Supplementary Table 13).

Splits	QFreq.	QS scores QC / QD / QI
Split 1		1.00 / NA / 1.00
Split 2		1.00 / NA / 1.00
Split 3		0.28 / 0.69 / 1.00
Split 4		0.08 / 0.00 / 0.99
Split 5		0.71 / 0.00 / 0.99
Split 6		0.01 / 0.85 / 0.98
Split 7		0.07 / 0.83 / 0.99
Split 8		0.89 / 1.00 / 1.00
Split 9		-0.001 / 0.94 / 0.99

bioRxiv preprint doi: <https://doi.org/10.1101/2020.09.22.307728>; this version posted September 23, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



Aeolidida

Aeolidioidea

Fionoidea

Dendronotoidea

Myrrhinae

Facelinidae

Aeolidiidae

Facelinidae

Flabellinidae

Flabellinopsidae

Trinchesiidae

Fionidae

Eubranchiidae

Unidentiidae

Samlidae

Janolidae

Dironidae

Embletoniidae

Tritoniidae

Arminidae

Dotidae

Bornellidae

Hancockiidae

Tethyidae

Scyllaeidae

Dendronotidae

Lomanotidae

Fionoidea

Aeolidioidea

Proctonotoidea

Fionoidea

Fionoidea

Tritonioidea

Tritonioidea

Arminoidea

Arminoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

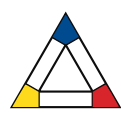
Dendronotoidea

Dendronotoidea

Dendronotoidea

Dendronotoidea

X Topologies investigated with AU test

 Investigated with FcLM

bioRxiv preprint doi: <https://doi.org/10.1101/2020.09.22.307728>; this version posted September 23, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.

