



Swansea University
Prifysgol Abertawe



Cronfa - Swansea University Open Access Repository

This is an author produced version of a paper published in:

The Biological Bulletin

Cronfa URL for this paper:

<http://cronfa.swan.ac.uk/Record/cronfa43224>

Paper:

Costa-Paiva, E., Schrago, C., Coates, C. & Halanych, K. (in press). Discovery of novel hemocyanin-like genes in Metazoans. *The Biological Bulletin*

This item is brought to you by Swansea University. Any person downloading material is agreeing to abide by the terms of the repository licence. Copies of full text items may be used or reproduced in any format or medium, without prior permission for personal research or study, educational or non-commercial purposes only. The copyright for any work remains with the original author unless otherwise specified. The full-text must not be sold in any format or medium without the formal permission of the copyright holder.

Permission for multiple reproductions should be obtained from the original author.

Authors are personally responsible for adhering to copyright and publisher restrictions when uploading content to the repository.

<http://www.swansea.ac.uk/library/researchsupport/ris-support/>

1 **Discovery of novel hemocyanin-like genes in Metazoans**

2

3 Elisa M. Costa-Paiva^{1,2*}, Carlos G. Schrago^{1*}, Christopher J. Coates³ and Kenneth M. Halanych²

4

5 1. Departamento de Genética, Laboratório de Biologia Evolutiva Teórica e Aplicada,

6 Universidade Federal do Rio de Janeiro, Rio de Janeiro, RJ, Brazil

7 2. Department of Biological Sciences, Molette Biology Laboratory for Environmental and

8 Climate Change Studies, Auburn University, Auburn, AL, 36849, USA

9 3. Department of Biosciences, College of Science, Swansea University, Singleton Park,

10 Swansea, Wales SA2 8PP, UK

11

12 *Author for Correspondence: elisapolychaeta@gmail.com; carlos.schrago@gmail.com

13

14

15 Running head: Hemocyanin-like genes in Metazoans

16

17

18

19

20

21

22

23

24

25

26

27 **Abstract**

28 Among animals, two major groups of oxygen-binding proteins are found: proteins that use
29 iron to bind oxygen (hemoglobins and hemerythrins) and two non-homologous hemocyanins that
30 use copper. Although arthropod and mollusc hemocyanins (herein designated HcA and HcM,
31 respectively) bind oxygen in the same manner, they are distinct in their molecular structures. In
32 order to better understand the range of natural variation in Hcs, we searched for Hcs in a diverse
33 array of metazoan transcriptomes using bioinformatics tools to examine Hc evolutionary history and
34 consequently revive the discussion about whether all metazoan Hcs shared a common origin with
35 frequent losses, or, originated separately after the divergence of Lophotrochozoa and Ecdysozoa.
36 We confirm that the distribution of Hc-like genes is more widespread than previously reported,
37 including five putative novel HcM genes in two annelid species from Chaetopteridae. For HcA, 16
38 putative novel genes were retained, and the presence of HcA in 11 annelid species represent novel
39 observations. Interestingly, Annelida is the lineage that presents the greatest repertoire of oxygen
40 transport proteins reported to date, possessing all the main superfamily proteins, which could be
41 explained partially by the immense variability of life styles and habitats. Work presented here
42 contradicts the canonical view that hemocyanins are restricted to molluscs and arthropods,
43 suggesting the occurrence of copper-based blood pigments in metazoans has been underestimated.
44 Our results also support the idea of the presence of oxygen carrier Hcs being widespread across
45 metazoans with an evolutionary history characterized by frequent losses.

46

47 **Introduction**

48 Four families of oxygen-binding/transporting proteins have been characterised into two
49 major groups: iron-based including hemoglobins and hemerythrins, and copper-based including two
50 non-homologous families of hemocyanins (Hcs) (Terwilliger et al. 1976; Burmester 2002; Coates
51 and Decker, 2017). These proteins are inextricably linked to metazoan evolution as they are
52 constrained by oxygen requirements of tissues, and therefore selection has presumably favored

53 proteins that can reversibly bind and transport oxygen (Terwilliger, 1998; Schmidt-Rhaesa 2007).
54 Although these molecules can reversibly bind oxygen, their binding affinities and evolutionary
55 origins differ and it is generally regarded that the diversity of blood pigments in animals has been
56 underestimated (Martín-Durán et al, 2013; Koch et al, 2016, Costa-Paiva et al, 2017 A-B).

57 Hcs are extracellular, oligomeric proteins that bind molecular oxygen reversibly between
58 two copper ions (CuA and CuB; Markl and Decker, 1992) in a side-on bridging arrangement. These
59 megamolecules, which can be similar in size to viruses or ribosomes (up to 13 MDa), are found
60 predominantly in the hemolymph of gastropod molluscs and numerous arthropods (Burmester,
61 2002; Coates and Decker, 2017). They are also characterized by their affinity for oxygen, which can
62 vary from low to moderate, and are regulated by a variety of allosteric factors (e.g. urea, calcium)
63 related to specific ecophysiological adaptations (Decker et al, 2007). Arthropod and mollusc Hcs,
64 despite being given the same name due to the conserved dicupric active sites, are distinct in their
65 sequence and structural compositions (Terwilliger, 1998; van Holde et al, 2001).

66 Hcs present in Arthropoda (HcA) are organized into hexamers, where each subunit has
67 three separate domains: (I) 5 or 6 α -helices, (II) 4 α -helices grouped together wherein lies the
68 dicopper active site, and (III) 7 antiparallel strands forming a β -barrel (Magnus et al., 1994). HcAs,
69 are members of a protein superfamily that include (a) insect and crustacean phenoloxidases whose
70 functions include sclerotization of the cuticle, wound healing, and humoral immune defense
71 (Whitten and Coates, 2017); (b) hexamerins (HEX), metal-free proteins present in insects that do
72 not bind oxygen, but are considered storage proteins associated with molting cycles or nutritional
73 conditions (Burmester, 1999a); (c) copper-free pseudo-hemocyanins or cryptocyanins (pHc or
74 CRY) are similar to hemocyanin but appear to act as storage proteins in the haemolymph
75 (Burmester, 1999b); and (d) hexamerin receptors (reHEX), which are present in dipterans and are
76 related to their own ligands (Burmester and Scheller, 1996). Although these proteins form a
77 functionally diverse superfamily, primary structures have highly conserved core elements that allow
78 their evolutionary history to be traced (Burmester, 2001).

79 Mollusc Hcs (HcM) consist of paralogous functional units derived from successive gene
80 duplication events (Markl, 2013). They are large proteins (up to 13 MDa) with numerous
81 polypeptide channels consisting of 10-12 or more peptides, and 7 or 8 functional units connected to
82 each other by bridging peptides (Van Holde et al., 2001). Each functional unit consists of two
83 domains: the tyrosinase or α -helical core domain where the active site is located, and the β -
84 sandwich domain (Cuff et al., 1998). Functionally, the α -domain in HcM corresponds to domain II
85 in HcA and the β -domain corresponds to domain III, respectively (Decker et al, 2007). There is high
86 similarity between HcA and HcM near the ligands of the CuB site, comprising a segment of about
87 50 residues or 10% of the polypeptide chain length (Drexel et al, 1987; Preaux et al, 1988). HcA
88 and HcM are structurally heterogeneous and represent two different classes of proteins, however the
89 similar segment around CuB indicates that both Hcs have evolved independently from a common
90 ancestral mononuclear copper protein (Markl and Decker, 1992).

91 Despite extensive studies on hemocyanins (Hcs) over the years (Burmester 2002, 2015;
92 Coates and Nairn, 2014), knowledge remains limited for animals except in molluscs and arthropods
93 where Hcs are broadly distributed (Burmester, 2001; Lieb et al, 2006; Markl, 2013; Kato et al,
94 2017). Isolated records of Hc in single species of Porifera (*Amphimedon queenslandica*),
95 Hemichordata (*Saccoglossus kowalevski*), and Ctenophora (*Mnemiopsis leidyi*) (Martín-Durán et al.
96 2013) have been reported. The presence of Hc-like proteins in the tunicate *Ciona intestinalis* with
97 putative phenoloxidase-like activity suggested that respiratory Hcs evolved from a phenoloxidase
98 (Immesberger and Burmester, 2004), in addition to the well-characterized inducible phenoloxidase
99 activity of HcAs and HcMs (reviewed by Coates and Nairn, 2014). Thus, given the absence of
100 functional data for non-bilaterian animals (i.e., ctenophore and sponge) and recent evidence
101 concerning greater Hc distribution in animals, we revive the discussion about whether all metazoan
102 Hcs shared a common origin with frequent losses (Martín-Durán et al. 2013), or, if Hcs originated
103 separately after the Lophotrochozoa and the Ecdysozoa diverged (van Holde et al, 2001). In
104 addition to metazoan records, similarities of HcAs with sequences in fungi (*Aspergillus niger*) have

105 also been reported, suggesting that the origin of Hcs (or type-III copper proteins more generally)
106 occurred in Opisthokonta, followed by multiple independent loss events (Martín-Durán et al.,
107 2013). However, evidence is scant to establish what the function of this ancestral protein might have
108 been (Burmester, 2015), as well as large gaps in sampling across Metazoa. In order to understand
109 the range of natural variation in Hcs, we searched for Hcs in a diverse array of metazoan
110 transcriptomes using a phylogenetic approach to examine Hc evolutionary history in the context of
111 animal phylogeny.

112

113 **Methods**

114 **Sample collection**

115 Information on species employed herein is provided in Table 1 and Appendix 1.
116 Transcriptomes of these species were collected using a variety of techniques, which include
117 intertidal sampling, dredge and box cores. Samples collected were preserved in RNALater or frozen
118 at -80°C.

119

120 **Data collection & sequence assembly**

121 Methods for RNA extraction, cDNA preparation and high-throughput sequencing generally
122 followed Kocot et al. (2011) and Whelan et al. (2015). The total RNA was extracted from either
123 whole animals (for small specimens) or from the body wall and coelomic region (for larger
124 specimens). RNA purifications were performed after extraction using TRIzol (Invitrogen) or the
125 RNeasy kit (Qiagen) with on-column DNase digestion, respectively. Reverse transcription used
126 single stranded RNA template and the SMART cDNA Library Construction Kit (Clontech) with
127 double stranded cDNA synthesis employing The Advantage 2 PCR system (Clontech). cDNA
128 libraries were barcoded and sequenced using Illumina technology by The Genomic Services Lab at
129 the Hudson Alpha Institute (Huntsville, Alabama, USA). Since sequencing was performed from
130 2012-2015, Paired End (PE) runs were of 100bp or 125bp lengths, utilizing either v3 or v4

131 chemistry on Illumina HiSeq 2000 or 2500 platforms (San Diego, California). CD Hit was
132 employed to look for redundant sequences. In order to facilitate sequence assembly, paired-end
133 transcriptome data were digitally normalized to an average k-mer coverage of 30 using normalize-
134 by-median.py (Brown et al, 2012; McDonald and Brown, 2013) and assembled using Trinity r2013-
135 02-25 with default settings (Grabherr et al, 2011).

136

137 **Data mining and gene identification**

138 Methods employed here were similar to those described in Costa-Paiva et al. (2017b). Two
139 approaches were utilized to mine transcriptomic data from 179 metazoan species and two
140 choanoflagellate species for putative Hc genes *in silico* (Appendix 1).

141 The first approach employed BLASTX (Altschul et al, 1990) at an e-value cutoff of 10^{-6} in
142 order to compare each assembled transcriptome contig ('queries') to a protein database composed
143 of 22 Hcs sequences from the National Center for Biotechnology (NCBI) database (Appendix 2) of
144 at least 500 amino acid residues. The BLASTX approach assured that any transcriptome contig with
145 a significant 'hit' to an Hc would be further evaluated in the pipeline. Then, initial contigs recovered
146 from these BLASTX searches were utilized in a second set of BLASTX searches against the NCBI
147 protein database (minimum e-value of 10^{-10}) and only top hits longer than 600 nucleotides retained
148 and were considered as putative Hc genes.

149 A second approach processed the transcriptomic data from the same species (Appendix 1)
150 through the Trinotate annotation pipeline (Grabherr et al, 2011), which utilizes a BLAST-based
151 approach to provide, among others, GO annotation (The Gene Ontology Consortium, 2004).
152 Transcripts annotated as Hcs, using the 10^{-5} e-value cutoff obtained by using BLASTX, were also
153 considered putative Hc-like gene orthologs.

154 Contigs identified as putative Hc genes using the two approaches described above were
155 subsequently translated into amino acids using TransDecoder with default settings (Haas et al,
156 2013). All translations were additionally subject to a Pfam domain evaluation using the EMBL-EBI

157 database with an e-value cutoff of 10^{-5} . Translations which returned Hc domains N, M or C or Hc
158 beta associated to tyrosinases in Pfam and that were longer than 200 amino acids residues were
159 retained for subsequent analyses. Transcripts passing the criteria described above were considered
160 Hc genes (Table 1; Appendix 1).

161

162 **Sequence alignment**

163 Because Hcs have been treated as two distinct proteins (Terwilliger, 1998), two protein
164 datasets were formed based on the Pfam domain evaluation results: (a) an HcM dataset and (b) an
165 HcA dataset.

166

167 **a) HcM dataset sequence alignment**

168 The HcM dataset included eight mollusc sequences previously used as ‘queries’ (Appendix
169 2) and five new sequences from translated transcripts. As HcM consists of a series of functional
170 units contain α and β domains, we opted to used partial sequences consisting in two functional units
171 ($2[\alpha \text{ domain} + \beta \text{ domain}]$) for each sequence, with the exception of *Mesochaetopterus* sequences
172 which presented just one single functional unit each. The dataset was aligned with MAFFT using
173 default “FFT-NS-2” algorithm (Kato and Standley, 2013), followed by visual inspection and
174 manual curation in order to remove spuriously aligned sequences based on similarity to the protein
175 alignment as a whole. Subsequently, ends of aligned sequences were manually trimmed in Geneious
176 9.1.3 (Kearse et al, 2013) to exclude residues 5’ of the putative start of a Tyr domain and 3’ residues
177 following the amino acid subsequent to the end of the second Hc domain. The resulting alignment
178 was used for all subsequent analyses (Supplemental File 1 available online).

179

180 **b) HcA dataset sequence alignment**

181 The HcA dataset was formed using 40 arthropod Hc superfamily sequences (Burmester,
182 2001; Aguilera et al, 2013; Martín-Durán et al, 2013; Appendix 3) and a remaining 16 sequences

183 from translated transcripts presenting at least two of the three HcA domains I, II, and III. Dataset
184 were aligned with MAFFT using default “FFT-NS-2” algorithm (Katoh and Standley, 2013),
185 followed by visual inspection and manual curation in order to remove spuriously aligned sequences
186 based on similarity to the protein alignment as a whole. Subsequently, aligned sequences were
187 trimmed using trimAl (Capella-Gutiérrez et al, 2009) with a 90% gap threshold were also
188 performed in order to eliminate poorly aligned regions. The resulting alignment was used for all
189 subsequent analyses (Supplemental File 2 available online).

190

191 **Phylogenetic analysis**

192 For each dataset, ProtTest3.4 was applied to carry out statistical selection of best-fit models
193 of protein evolution for each dataset separately using the Akaike and Bayesian Information Criteria
194 (AIC and BIC, respectively) methods (Darriba et al, 2011).

195 Bayesian phylogenetic inference was performed with MrBayes 3.2.1 (Ronquist and
196 Huelsenbeck, 2003) for each database separately employing two independent MCMC runs. In each
197 run, four Metropolis-coupled chains were sampled every 500th cycle for 10⁷ generations. In order to
198 confirm if chains achieved stationary and determine an appropriate burn-in, we evaluated trace plots
199 of all MrBayes parameter output in Tracer v1.6 (Rambaut et al, 2014). The first 25% of samples
200 were discarded as burn-in and a majority rule consensus tree generated using MrBayes. Bayesian
201 posterior probabilities were used for assessing statistical support of each bipartition.

202

203 **Contamination screening**

204 We under took procedures to reduce the chance of false positives in our analysis, meaning
205 genes that are not homologous to Hcs or genes obtained from contamination. In most samples
206 included in our work, the total RNA was extracted from the body wall and coelomic region, which
207 excluded any possible contamination from food residues. Moreover, all the species included in our

208 analyses in which Hcs were found were prepared and sequenced separate from any mollusc or
209 arthropod species, which makes cross-contamination highly unlikely.

210 Furthermore, for *in silico* analyses, we opted for a more conservative approach and employed
211 very stringent e-values cutoff ($< 10^{-5}$). After translation, proteins from contigs identified as putative
212 Hc genes were subsequently subject to a Pfam Domain evaluation as described above and a
213 BLASTP (Altschul et al, 1990) search against NCBI protein database (minimum e-value of 10^{-10}).
214 The top hits were Hc sequences from either mollusc (for new HcM sequences) or arthropod (for
215 new HcA sequences) and this could be easily explained since sequences of Hcs from other
216 metazoan species are still rare in NCBI.

217 Therefore, as we opted for a very conservative approach, we decided to not include two
218 sequences identified as HcAs in the dataset, one from an annelid *Cossura longocirrata* and one
219 from a cycliophor *Symbion americanus*. Although these sequences had an Hc Pfam domain and
220 presented the six conservative histidine residues, they also showed high identity values ($> 80\%$) to
221 crustacean Hcs when submitted to BLASTP and MEGABLAST searches.

222

223 **Results**

224 **a) HcM**

225 Our bioinformatic pipeline (Figure 1) recovered a total of 18 unique protein sequences of
226 mollusc Hc-like genes from 181 transcriptomes representing 15 metazoan phyla and two
227 choanoflagellate species (Appendix 1). Following translation, Pfam domain evaluation, and
228 removal of sequences with fewer than 200 amino acid residues, five putative novel HcM genes were
229 retained from all taxa examined here, representing two annelid species from Chaetopteridae (Table
230 1, Supplemental File 1 available online). For both choanoflagellate species and all other metazoan
231 phyla, we did not find any HcM gene (Appendix 1). Alignment of translated transcripts possessed
232 956 residue positions and partial sequences consisted of two functional units (2[α domain + β
233 domain]) for each sequence, with the exception of *Mesochaetopterus* sequences which had one

234 functional unit each. New sequences were combined with eight publically available HcMs
235 previously used as ‘queries’ (Appendix 2) to produce a final dataset of 13 HcM sequences (Figure
236 2; Supplemental File 1 available online; Costa-Paiva, 2018a).

237 The Bayesian inference analysis (Figure 3) revealed two highly supported clades (Figure 3,
238 p.p. = 0.99) even though other internal nodes were less resolved, which are often observed in gene
239 genealogies (DeSalle, 2015). Aside from previous mollusc records (Markl, 2013), we found novel
240 HcM genes in two annelid species (Table 1) both within Chaetopteridae. All five sequences
241 included canonical functional units composed of α - and β -domains. The topology of the HcM gene
242 tree did not mirror recent phylogenies of Mollusca and the relationship between Mollusca and
243 Annelida (Kocot et al, 2011; Halanych, 2016; Kocot et al, 2017). We found a strongly supported
244 sub-group (Figure 3, pink clade B, p.p. = 0.99) with two *Phyllochaetopterus* Hcs and HcM
245 sequences from caudofoveates, cephalopods, and bivalves. However, the remaining
246 *Phyllochaetopterus* sequence formed a clade with other *Mesochaetopteus* sequences (Figure 3, blue
247 clade A, p.p. = 0.76), which is in a clade with other polyplacophorans, gastropods, and cephalopods.

248 Our data were based on available transcriptomes, and the absence of a certain gene when
249 searching in transcriptomes does not necessarily mean that it is absent from the genome. Thus, we
250 carefully addressed that we did not find any HcM genes in any other metazoan and
251 choanoflagellates as expected (Appendix 1).

252

253 **b) HcA**

254 Our bioinformatic analyses for HcA (Figure 1) recovered a total of 137 unique protein
255 sequences of arthropod Hc-like genes from 181 transcriptomes from 15 metazoan phyla and two
256 choanoflagellate species (Table 1). After nucleic acid translation, Pfam domain evaluation, and
257 removal of sequences with fewer than 200 amino acid residues, 16 putative novel genes
258 representing at least two of the three HcA domains were retained, representing 16 individual species
259 distributed across four phyla (Table 1). Concerning these 16 new putative HcA genes, just one

260 possessed all three canonical arthropod domains (I, II, and III), the annelid *Streblosoma hartmanae*
261 (Table I). The remaining 15 sequences contained two domains (I and II or II and III). Domain II, the
262 location of the dicopper active site with six histidine residues, was found in each species, therefore
263 we included these records as putative HcA genes (Figure 4). In order to understand the relationship
264 of these putative HcA genes with other members of the Hc superfamily, the alignment of translated
265 transcripts possessing 359 residue positions were performed for a dataset containing these 16
266 putative HcA genes and 40 publically available sequences from arthropod Hc superfamily
267 representatives (Burmester, 2001; Aguilera et al., 2013; Martín-Durán et al., 2013) (Supplemental
268 File 2 available online; Costa-Paiva, 2018b).

269 The Bayesian inference analysis of HcA superfamily revealed five supported clades: A) a
270 green clade (Figure 5, p.p. = 1) formed by hexapod hexamerin sequences; B) a blue clade (Figure 5,
271 p.p. = 1) formed exclusively by crustaceans HcAs, cryptocyanins and pseudo-hemocyanins; C) an
272 orange clade (Figure 5, p.p. = 1) formed by myriapod and chelicerate HcAs; D) a pink clade (Figure
273 5, p.p. = 1) formed by hexapod and crustacean prophenoloxidase sequences; and E) a yellow clade
274 (Figure 5, p.p. = 1) formed by non-arthropod HcA sequences, including the 16 novel sequences
275 from annelids, hemichordates, sponges, and ctenophores. The circled clade inside the yellow clade
276 was formed exclusively by ctenophores HcAs sequences.

277

278 **Discussion**

279 Herein, we confirm that the distribution of Hc genes is more widespread than previously
280 reported. Our results describe actively transcribed HcMs genes in two chaetopterid annelid species,
281 and HcAs in 16 species distributed across four metazoan phyla (Table 1, Figure 6). Of the four
282 phyla, HcAs in Ctenophora, Porifera, and Hemichordata were reported previously (Aguilera et al.
283 2013; Martín-Durán et al. 2013). Importantly, the presence of HcAs in Annelida represent new
284 records. Our work is contrary to the traditional view that Hcs are restricted to Molluscs and
285 Arthropods (Burmester, 2001; Markl, 2013), corroborating the recent findings of Martin-Duran et al

286 (2013) that the presence of Hcs, as well as other oxygen carrier molecules such as hemerythrins
287 (Costa-Paiva et al., 2017 B) and globins (Koch et al., 2016), are underestimated in animals.
288 Although our results cannot empirically prove that these newly discovered genes effectively
289 transport oxygen, we present evidence here that this function is entirely possible (e.g., orthology to
290 Hcs, PFAM structure, and GO ontology). Additionally, these genes were identified based on
291 sequence similarity to previously well-characterized Hcs under the assumption that sequence
292 similarity is generally indicative of function (Gabaldón and Huynen, 2004).

293 Previous studies (e.g., Martín-Durán et al., 2013) have demonstrated that the α domain
294 (tyrosinase domain) has a wide distribution across metazoan lineages, with the exception of
295 arthropods, which can be explained by the expansion and diversification of the Hc domain II in this
296 group of animals. However, the tyrosinase domain itself can play several roles in addition to
297 respiratory function, such as melanin biosynthesis (Sugumaran, 2002). Although respiratory
298 function requires the presence of both domains α and β , as we found in molluscs, the presence of
299 few functional units in chaetopterids could indicate another function besides oxygen transport for
300 these molecules in agreement with recent studies on the functional versatility of Hcs. Examples
301 include antimicrobial peptide production, host-symbiont dynamics, and Hc-derived phenoloxidase
302 activity (Zhuang et al., 2015; Coates and Nairn, 2014; Kremer et al., 2014).

303 In relation to HcAs, our data corroborate previous findings of HcAs in hemichordates,
304 sponges, and ctenophores (Martín-Durán et al, 2013). Moreover, not all of our newly discovered
305 genes possessed the three domains, but all of them were considered to be an HcA because they
306 possessed domain II, where the active site is located (Decker et al, 2007). The presence of domain II
307 spread across metazoan lineages and the presence of the same domain in amebozoans and in the
308 filamentous fungus *Aspergillus niger* (Martín-Durán et al, 2013) are likely to be an unikont
309 synapomorphy as suggested before. Despite previous suggestions that the domain I (*N*-domain) of
310 HcA can be used as a specific molecular signature of the Pan-arthropoda (Martín-Durán et al,

311 2013), our findings concerning the presence of domain I in annelids and ctenophores contradict this
312 idea (Figure 6).

313 Interestingly, Annelida is the lineage that presents the greatest repertoire of oxygen transport
314 proteins reported to date, possessing all the main superfamily proteins: Hrs, Hbs, HcAs, and HcMs
315 (Rouse and Pleijel, 2001; Costa-Paiva et al, 2017a). This fact could be explained by their ancient
316 origin and the early radiation of this group, and also by the great diversification of life forms and
317 habitats which leads to a great variety of oxygen absorption and transport strategies inside the body
318 (Schumway, 1979; Rouse & Pleijel, 2001). Furthermore, the high diversity of oxygen-binding
319 proteins can be found even in the same lineage within annelids, as observed in chaetopterids. Based
320 on our data, a single individual of *Mesochaetopterus alipes* can actively transcribe Hrs (Costa-Paiva
321 et al, 2017a) and mollusc-like Hcs. Other annelids also presented more than one family of oxygen-
322 binding protein, for example some species of Terebellidae, Opheliidae and Sipuncula (Bailly et al,
323 2008; Liu et al, 2013). Such organisms may simultaneously express more than one protein or may
324 have different protein expression in different parts of the body (Bailly et al, 2008).

325 Although this unexpected large repertoire of oxygen-binding proteins in an annelid could be
326 partially explained by the need to transport oxygen molecules, secondary functional specializations
327 could also have been an important trigger for the diversification of these molecules throughout their
328 evolutionary history in this group. There are records of Hrs molecules involved in other functions
329 besides oxygen loading, such as in the storage of iron atoms, detoxification of heavy metals, and
330 metabolic pathways related to the innate immunity of some species of annelids (e.g., *Theromyzon*
331 *tessulatum*, *Hirudo medicinalis*, and *Neanthes diversicolor*; Baert et al, 1992; Demuyne et al,
332 1993; Vergote et al, 2004; reviewed by Coates and Decker, 2017). Moreover, Hcs are known for
333 contributing in many ways to immune defenses, such as the inhibition of viral replication, precursor
334 of antimicrobial peptides, and the conformational switch to a phenoloxidase-like enzyme (Coates
335 and Nairn, 2014; Coates and Talbot, 2018).

336 Our results support the idea of the presence of Hcs being widespread across metazoans (Figure
337 6) with an evolutionary history characterized by frequent losses (Aguilera et al, 2013; Martín-Durán
338 et al, 2013). Such losses are observed within several lineages of molluscs and arthropods. For
339 example, the gastropod family Planorbidae, which lack HcM in their hemolymph and utilize
340 extracellular hemoglobin for oxygen transport (Ochiai et al, 1989; Arndt and Santoro, 1998),
341 probably due to hemoglobin's higher affinity for oxygen than the ancestral Hc (Lieb et al, 2001).
342 The same can be found in crustacean lineages, as branchiopods, ostracods, copepods, cirripeds, and
343 decapods that lost HcA and presented hemoglobin for handling oxygen (Terwilliger and Ryan,
344 2001).

345 The revised distribution in expression of HcM and HcA genes across metazoans could be
346 explained by the differences in physio-chemical properties of the oxygen binding domains and the
347 life histories of disparate animal lineages. Obtaining functional data on these newly discovered Hc
348 genes is needed to evaluate the significance of their widespread occurrence in metazoans, and
349 oxygen-binding/transport proteins in general.

350

351 **Acknowledgements**

352 We thank Fernando Avila Queiroz for valuable help with the figures. EMCP was supported by
353 CAPES (Coordenação de Aperfeiçoamento de Pessoal de Nível Superior, Brazil). Use of SkyNet
354 computational resources at Auburn University is acknowledged. This work was funded by National
355 Science Foundation grant DEB-1036537 to K. M. Halanych and Scott R. Santos and OCE-1155188
356 to K. M. Halanych. This is Molette Biology Laboratory contribution ## and Auburn University
357 Marine Biology Program contribution ###. We are grateful to the three anonymous reviewers and
358 editorial board for their candor.

359

360

361

362 **References**

- 363 **Aguilera F., C. McDougall, and B. M. Degnan. 2013.** Origin, evolution and classification of type-
364 3 copper proteins: lineage-specific gene expansions and losses across the Metazoa. *BMC Evol. Biol.*
365 **13**(1): 96.
- 366 **Altschul S. F., W. Gish, W. Miller, E. W. Myers, D. J. Lipman. 1990.** Basic local alignment
367 search tool. *J. Mol. Biol.***215**: 403–410.
- 368 **Arndt M. H., and M. M. Santoro. 1998.** Structure of the extracellular hemoglobin of
369 *Biomphalaria glabrata*. *Comp. Biochem. Physiol. B: Biochem. Mol. Biol.***119**(4): 667-675.
- 370 **Baert J.L., M. Britel, P. Sautière, and J. Malécha. 1992.** Ovohemerythrin, a major 14-kDa yolk
371 protein distinct from vitellogenin in leech. *Eur. J. Biochem.***209**: 563-569.
- 372 **Bailly X., S. Vanin, C. Chabasse, K. Mizuguchi, and S. N. Vinogradov. 2008.** A phylogenomic
373 profile of hemerythrins, the nonheme diiron binding respiratory proteins. *BMC Evol. Biol.***8**: 244.
- 374 **Brown C. T., A. Howe, Q. Zhang, A. B. Pyrkosz, and T. H. Brom. 2012.** A reference-free
375 algorithm for computational normalization of shotgun sequencing data. arXiv:1203.4802 [q-
376 bio.GN], available: <https://arxiv.org/pdf/1203.4802.pdf>. [2018, August 05].
- 377 **Burmester T. 1999.** Evolution and function of the insect hexamerins. *Eur. J. Entomol.***96**: 213-226.
378 **(A)**
- 379 **Burmester T. 1999.** Identification, molecular cloning, and phylogenetic analysis of a non-
380 respiratory pseudo-hemocyanin of *Homarus americanus*. *J. Biol. Chem.* **274**(19): 13217-13222.
381 **(B)**
- 382 **Burmester T. 2001.** Molecular evolution of the arthropod hemocyanin superfamily. *Mol. Biol.*
383 *Evol.* **18**(2): 184-195.

384 **Burmester T. 2002.** Origin and evolution of arthropod hemocyanins and related proteins. *J. Comp.*
385 *Physiol. B.***172**(2): 95–107.

386 **Burmester T. 2015.** Evolution of respiratory proteins across the Pancrustacea. *Integr. Comp.*
387 *Biol.***55**(5): 765-770.

388 **Burmester T., and K. Schellen. 1996.** Common origin of arthropod tyrosinase, arthropod
389 hemocyanin, insect hexamerin, and dipteran arylphorin receptor. *J. Mol. Evol.***42**(6): 713-728.

390 **Capella-Gutiérrez S., J. M. Silla-Martínez, and T. Gabaldón. 2009.** TrimAl: a tool for automated
391 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*, **25**(15): 1972-1973.

392 **Coates C. J. and H. Decker. 2017.** Immunological properties of oxygen-transport proteins:
393 hemoglobin, hemocyanin and hemerythrin. *Cell. Mol. Life Sci.***74**(2): 293-317.

394 **Coates C. J. and J. Nairn. 2014.** Diverse immune functions of hemocyanins. *Dev. Comp.*
395 *Immunol.***45**(1): 43-55.

396 **Coates C. J. and J. Talbot. 2018.** Hemocyanin-derived phenoloxidase reaction products display
397 anti-infective properties. *Dev. Comp. Immunol.***86**: 47-51.

398 **Costa-Paiva E. M. 2018.** Mollusc hemocyanin dataset [online]. Figshare file DOI:
399 10.6084/m9.figshare.5301730 [2018, May 08]. (A)

400 **Costa-Paiva E. M. 2018.** Arthropod hemocyanin superfamily dataset [online]. Figshare file DOI:
401 10.6084/m9.figshare.5301691) [2018, May 08]. (B)

402 **Costa-Paiva E. M., N. V. Whelan, D. S. Waits, S. R. Santos, C. G. Schrago, and K. M.**
403 **Halanych. 2017.** Discovery and evolution of novel hemerythrin genes in annelid worms. *BMC*
404 *Evol. Biol.* **17**: 85-96. (A)

405 **Costa-Paiva E. M., C. G. Schrago, and K. M. Halanych. 2017.** Broad phylogenetic occurrence of
406 the oxygen-binding hemerythrins in bilaterians. *Genome Biol. Evol.*,9(10): 2580-2591. (B)

407 **Cuff M., K. Miller, K. van Holde, and W. Hendrickson. 1998.** Crystal structure of a functional
408 unit from Octopus hemocyanin. *J. Mol. Biol.* **278**: 855–870.

409 **Darriba D., G. L. Taboada, R. Doallo, and D. Posada D. 2011.** ProtTest 3: fast selection of best-
410 fit models of protein evolution. *Bioinformatics.* **27**(8): 1164-1165.

411 **Decker H., N. Hellmann, E. Jaenicke, B. Lieb, U. Meissner, and J. Markl. 2007.** Recent
412 progress in hemocyanin research. *Integr. Comp. Biol.* **47**(4): 631–644.

413 **Demuyneck S., K. W. Li, R. van der Schors, and N. Dhainaut-Courtois. 1993.** Amino acid
414 sequence of the small cadmium-binding protein (MP-II) from *Nereis diversicolor* (Annelida,
415 Polychaeta) – evidence for a myohemerythrin structure. *Eur. J. Biochem.*217: 151–156.

416 **DeSalle R. 2015.** Can single protein and protein family phylogenies be resolved better? *J.*
417 *Phylogenetics Evol. Biol.* **3**: e116.

418 **Drexel R., S. Siegmund, H. J. Schneider, B. Linzen, C. Gielens, G. Préaux, R. Lontie, J.**
419 **Kellermann, and F. Lottspeich. 1987.** Complete amino-acid sequence of a functional unit from a
420 molluscan hemocyanin (*Helix pomatia*). *Biol. Chem.* **368**(1): 617-636.

421 **Gene Ontology Consortium. 2004.** The Gene Ontology (GO) database and informatics resource.
422 *Nucleic Acids Res.* **32**(1): D258-D261.

423 **Grabherr M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson, I. Amit, X. Adiconis, L.**
424 **Fan, R. Raychowdhury, Q. Zeng, et al. 2011.** Full-length transcriptome assembly from RNA-Seq
425 data without a reference genome. *Nature Biotechnol.***29**: 644–652.

426 **Haas B. J., A. Papanicolaou, M. Yassour, M. Grabherr, P.D. Blood, J. Bowden, M. B. Couger,**
427 **D. Eccles, B. Li, M. Lieber, et al. 2013.** De novo transcript sequence reconstruction from RNA-seq
428 using the Trinity platform for reference generation and analysis. *Nature Protocols*.**8**(8): 1494-1512.

429 **Halanych K. M. 2016.** How our view of animal phylogeny was reshaped by molecular approaches:
430 lessons learned. *Org. Divers. Evol.***16**(2): 319-328.

431 **Immesberger A., and T. Burmester. 2004.** Putative phenoloxidases in the tunicate *Ciona*
432 *intestinalis* and the origin of the arthropod hemocyanin superfamily. *J. Comp. Physiol. B.***174**(2):
433 169-180.

434 **Kato S., T. Matsui, C. Gatsogiannis, and Y. Tanaka. 2018.** Molluscan hemocyanin: structure,
435 evolution, and physiology. *Biophys. Rev.* **10**(2):191-202.

436 **Katoh K., and D. M. Standley. 2013.** MAFFT multiple sequence alignment software version 7:
437 improvements in performance and usability. *Mol. Biol. Evol.* **30**(4): 772–780.

438 **Kearse M., R. Moir, A. Wilson, S. Stones-Havas, M. Cheung, S. Sturrock, S. Buxton, A.**
439 **Cooper, S. Markowitz, C. Duran, et al. 2012.** Geneious Basic: an integrated and extendable
440 desktop software platform for the organization and analysis of sequence data.
441 *Bioinformatics*;**28**(12): 1647–1649.

442 **Koch J., J. Lüdemann, R. Spies, M. Last, C. T. Amemiya, and T. Burmester. 2016.** Unusual
443 diversity of myoglobins genes in the lungfish. *Mol. Biol. Evol.***33**(12): 3033-3041.

444 **Kocot K. M., Cannon J.T., Todt C., Citarella M.R., Kohn A.B., Meyer A., Santos S.R.,**
445 **Schander C., Moroz L.L., Lieb B., et al. 2011.** Phylogenomics reveals deep molluscan
446 relationships. *Nature*; **477**: 452–456.

447 **Kocot K. M., T. H. Struck, J. Merkel, D. S. Waits, C. Todt, P. M. Brannock, D. A. Weese, J. T.**
448 **Cannon, L. L. Moroz, B. Lieb, et al. 2017.** Phylogenomics of Lophotrochozoa with consideration
449 of systematic error. *Syst. Biol.***66**(2): 256-282.

450 **Kremer N., J. Schwartzman, R. Augustin, L. Zhou, E. G. Ruby, S. Hourdez, and M. J. McFall-**
451 **Ngai. 2014.** The dual nature of haemocyanin in the establishment and persistence of the squid–
452 vibrio symbiosis. *Proc. R. Soc. Lond. B: Biol. Sci.***281**(1785): 20140504.

453 **Lieb B., B. Altenhein, J. Markl, A. Vincent, E. van Olden, K. E. van Holde, and K. I. Miller.**
454 **2001.** Structures of two molluscan hemocyanin genes: significance for gene evolution. *Proc. Natl.*
455 *Acad. Sci.***98**(8): 4546-4551.

456 **Liu Y., C. Li, X. Su, M. Wang, Y. Li, Y. Li, and T. Li. 2013.** Cloning and characterization of
457 hemerythrin gene from Sipuncula Phascolosoma esculenta. *Genes Genom.***35**(1), 95-100.

458 **Magnus K., B. Hazes, H. Ton-That, C. Bonaventura, J. Bonaventura, and W. Hol. 1994.**
459 Crystallographic analysis of oxygenated and deoxygenated states of arthropod hemocyanin shows
460 unusual differences. *Proteins*;**19**: 302–309.

461 **Markl J. 2013.** Evolution of molluscan hemocyanin structures. *Biochim. Biophys. Acta (BBA)-*
462 *Proteins and Proteomics*;**1834**(9): 1840-1852.

463 **Markl J., and H. Decker. 1992.** Molecular structure of the arthropod hemocyanins. Pp. 325-376 in
464 *Advances in Comparative & Environmental Physiology Vol 13 – Blood and Tissue Oxygen Carriers,*
465 Mangum ChP, ed. Springer, Berlin, Heidelberg.

466 **Martín-Durán J.M., A. De Mendoza, A. Sebé-Pedrós, I. Ruiz-Trillo, and A. Hejzol. 2013.** A
467 broad genomic survey reveals multiple origins and frequent losses in the evolution of respiratory
468 hemerythrins and hemocyanins. *Genome Biol. Evol.***5**: 1435–1442.

469 **McDonald E. and C. T. Brown. 2013.** Khmer: working with big data in bioinformatics. *CoRR*,
470 abs/1303.2223.

471 **Ochiai T., Y. Enoki, and I. Usuki. 1989.** Physicochemical properties of the extracellular
472 hemoglobin from the planorbid snail, *Indoplanorbis exustus*. *Comp. Biochem. Physiol. B: Comp.*
473 *Biochem.***93**(4): 935-940.

474 **Préaux G., C. Gielens, R. Witters, and R. Lontie. 1988.** The structure of molluscan haemocyanins
475 and their homology with tyrosinases. *B. Soc. Chim. Belg.***97**(11-12): 1037-1044.

476 **Rambaut A., M. A. Suchard, D. Xie, and A. J. Drummond. 2014.** Tracer v1.6. Available from
477 <http://beast.bio.ed.ac.uk/Tracer>. 2014. [2017, May 15]

478 **Ronquist F., and J. P. Huelsenbeck. 2003.** MrBayes 3: Bayesian phylogenetic inference under
479 mixed models. *Bioinformatics*;**19**(12): 1572-1574.

480 **Rouse G. and F. Pleijel. 2001.** *Polychaetes*. Oxford university press. 354 pp.

481 **Schmidt-Rhaesa A. 2007.** *The Evolution of Organs Systems*. Oxford University Press. 385 pp.

482 **Sugumaran M. 2002.** Comparative biochemistry of eumelanogenesis and the protective roles of
483 phenoloxidase and melanin in insects. *Pigment Cell & Melanoma Research*;**15**(1): 2-9.

484 **Terwilliger N. B. 1998.** Functional adaptations of oxygen-transport proteins. *J. Exp. Biol.***201**:
485 1085-1098.

486 **Terwilliger N. B., and M. Ryan. 2001.** Ontogeny of crustacean respiratory proteins. *Am.*
487 *Zool.***41**(5): 1057-1067.

488 **Terwilliger R. C., N. B. Terwilliger, and E. Schabtach. 1976.** Comparison of chlorocruorin and
489 annelid hemoglobin quaternary structures. *Comp. Biochem. Physiol. A.***55**(1): 51-55.

490 **van Holde K. E., K. I. Miller, and H. Decker. 2001.** Hemocyanins and Invertebrate Evolution. *J.*
491 *Biol. Chem.***276**(19): 15563–66.

492 **Vergote D., P. E. Sautière, F. Vandenbulcke, D. Vieau, G. Mitta, E. R. Macagno, and M. Salzet.**
493 **2004.** Up-regulation of neurohemerythrin expression in the central nervous system of the medicinal
494 leech, *Hirudo medicinalis*, following septic injury. *J. Biol. Chem.***279** (42): 43828–37.

495 **Weigert A., and C. Bleidorn. 2016.** Current status of annelid phylogeny. *Org. Divers. Evol.***16**(2):
496 345-362.

497 **Whelan N. V., K. M. Kocot, L. L. Moroz, and K. M. Halanych. 2015.** Error, signal, and the
498 placement of Ctenophora sister to all other animals. *Proc. Natl. Acad. Sci. USA.***112**(18): 5773-5778.

499 **Whitten M., and C. J. Coates. 2017.** Re-evaluation of insect melanogenesis research: Views from
500 the dark side. *Pigment Cell & Melanoma Research*;**30**(4): 386-401.

501 **Zhuang J., C. J. Coates, H. Zhu, P. Zhu, Z. Wu, and L. Xie. 2015.** Identification of candidate
502 antimicrobial peptides derived from abalone hemocyanin. *Dev. Comp. Immunol.***49**(1), 96-102.

503

504

505

506 **Table 1.** List of taxa analyzed with novel genes and total number of contigs after
507 assembly. Number and type of putative Hc genes and accession numbers are also
508 provided.

Taxon	Total contigs number	Number and type of putative Hc genes	Accession number
METAZOA			
Porifera			
<i>Kirkpatrickia variolosa</i> (Kirkpatrick, 1907)	100,231	1 partial HcA Domains II + III	MF998096
<i>Latrunculia apicalis</i> Ridley & Dendy, 1886	76,210	1 partial HcA Domains II + III	MF998097
Ctenophora			
<i>Coeloplana astericola</i> Mortensen, 1927	222,614	1 partial HcA Domains I + II	MF998091
<i>Mnemiopsis leidyi</i> A. Agassiz, 1865	385,798	1 partial HcA Domains II + III	MF998101
<i>Pleurobrachia bachei</i> A. Agassiz, 1860	38,856	2 partial HcA Domains II + III	MF998107 MF998108
Hemichordata			
Harrimaniidae gen sp. (from Iceland)	230,054	1 partial HcA Domains II + III	MF998095
Annelida			
<i>Dorydrilus michaelsoni</i> Piguet, 1913	136,096	1 partial HcA Domains II + III	MF998093
<i>Eupolymnia nebulosa</i> (Montagu, 1819)	139,021	1 partial HcA Domains II + III	MF998094
<i>Lumbrineris perkinsi</i> Carrera-Parra, 2001	144,648	1 partial HcA Domains II + III	MF998098
<i>Mesochaetopterus alipes</i> Monro, 1928	83,209	2 HcM	MF998099 MF998100
<i>Paramphinome jeffreysii</i> (McIntosh, 1868)	165,337	1 partial HcA Domains II + III	MF998102
<i>Phyllochaetopterus prolifica</i> Potts, 1914	193,836	3 HcM	MF998103 MF998104 MF998105
<i>Pista macrolobata</i> Hesse, 1917	126,764	1 partial HcA Domains II + III	MF998106
<i>Streblosoma hartmanae</i> Kritzler, 1971	108,080	1 HcA Domains I + II + III	MF998109
<i>Stylodrilus heringianus</i> Claparede, 1862	239,935	1 partial HcA Domains II + III	MF998110
<i>Terebellides stroemii</i> Sars, 1835	169,760	1 partial HcA Domains II + III	MF998112
<i>Thelepus crispus</i> Johnson, 1901	67,478	1 partial HcA Domains II + III	MF998113

509 **Figure Legends**

510 **Figure 1** – Bioinformatics pipeline. Rounded rectangles represent input / output files, ovals
511 represent software or scripts, and the hexagon represents a step which involving manual evaluation.
512 Eight mollusc Hc sequences previously used as query sequences from Genbank (Appendix 2) were
513 also included in the HcM dataset and 42 arthropod HcA superfamily protein sequences (Appendix
514 3) were also included in the HcA dataset for posterior analyses.

515 **Figure 2** – HcM dataset partial alignment evidencing the active-site of six conserved histidine
516 residues (arrows).

517 **Figure 3** – Unrooted bayesian tree for HcM using MrBayes 3.2.1. A) Blue clade represents the
518 remaining annelid sequences with polyplacophoran, gastropods, and cephalopod sequences. The
519 circled blue clade represents the *Phyllochaetopterus* sequence with other *Mesochaetopteus* novel
520 sequences. B) Pink clade represents two *Phyllochaetopterus* Hcs (circled clade) and HcM
521 sequences from caudofoveats, cephalopods, and bivalves. The number after the name of each
522 sequence indicates the GenBank accession numbers for each previously identified HcM gene and it
523 is indicated in Appendix 2.

524 **Figure 4** – HcA dataset partial alignment evidencing the active-site of six conserved histidine
525 residues (arrows).

526 **Figure 5** – Unrooted bayesian tree for HcA using MrBayes 3.2.1. A) green clade is formed by
527 hexapod hexamerin sequences; B) blue clade by crustacean HcAs, cryptocyanins and pseudo-
528 hemocyanins; C) orange clade included myriapod and chelicerate HcAs; D) pink clade
529 prophenoloxidase sequences; and E) yellow clade is formed by non-arthropod HcA sequences,
530 which include the 16 novel sequences from annelids, hemichordates, sponges, and ctenophores. The
531 circled yellow clade represented a supported clade formed by ctenophores HcAs sequences. The
532 number after the name of each sequence indicates the GenBank accession numbers for each
533 previously identified HcA superfamily gene and it is indicated in Appendix 3.

534 **Figure 6** - Hypothesized relationships among metazoan phyla derived from recent phylogenomic
535 studies. Red rectangles represent HcM records and blue rectangles represent HcA records. The
536 domain structure of the Hc protein found is demonstrated alongside each respective taxon.

537

538 **Appendix**

539 **Appendix 1** - List of all taxa analyzed and total number of contigs after assembly. For underlined
540 taxa, number and type of putative Hc genes and accession numbers are also provided.

541 **Appendix 2** - Queries used to search the assembled translated transcriptomes. All HcM sequences
542 were also included in the dataset previous to the alignment.

543 **Appendix 3** - HcA superfamily protein sequences used in Burmester (2001), Aguilera et al. (2013),
544 and Martín-Durán et al. (2013) with genes accession numbers for each species.

545

546 **Supplemental Files**

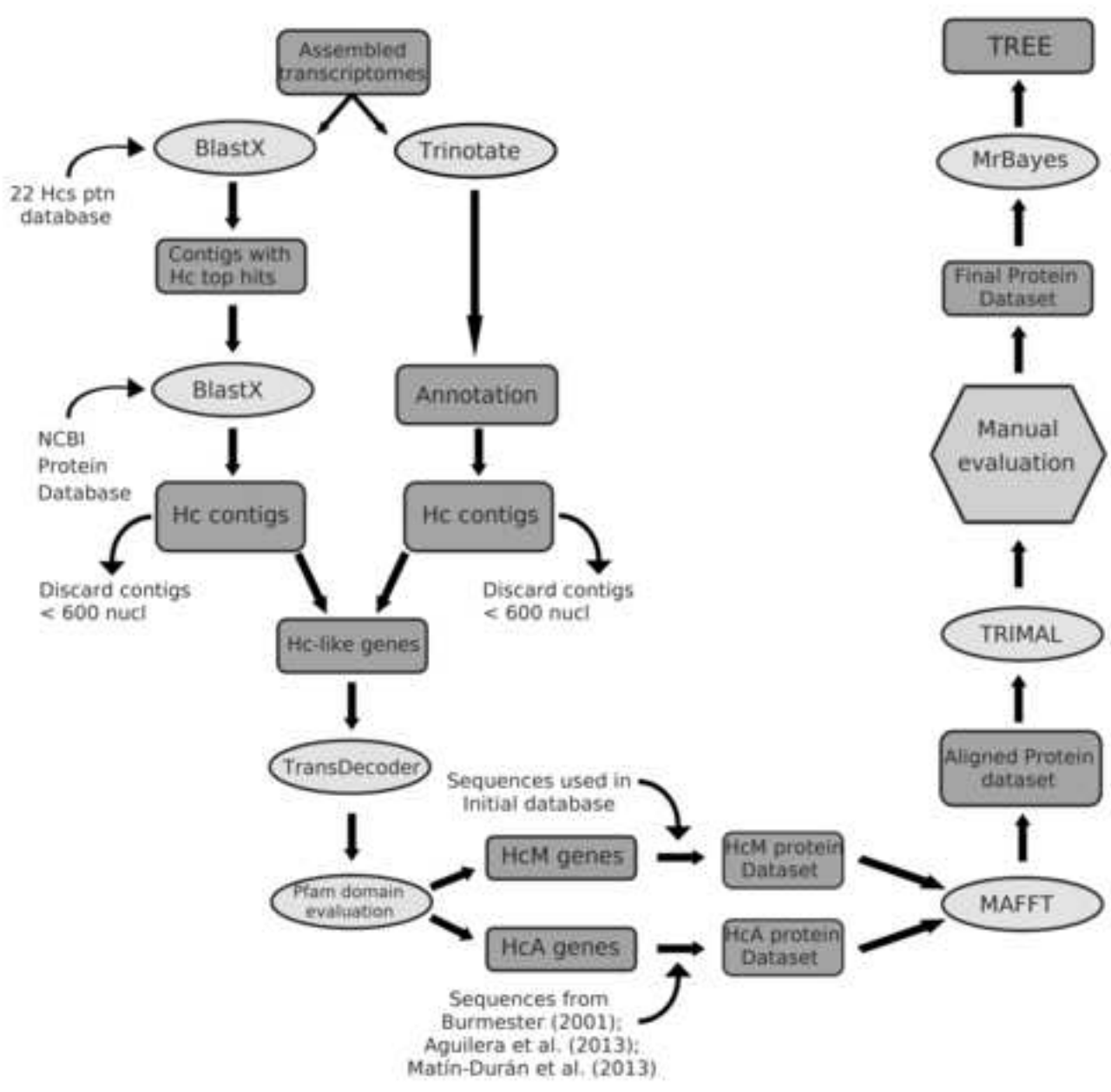
547 **Supplemental File 1** - The amino acid alignment for HcM used in analyses.

548 **Supplemental File 2** - The amino acid alignment for HcA used in analyses.

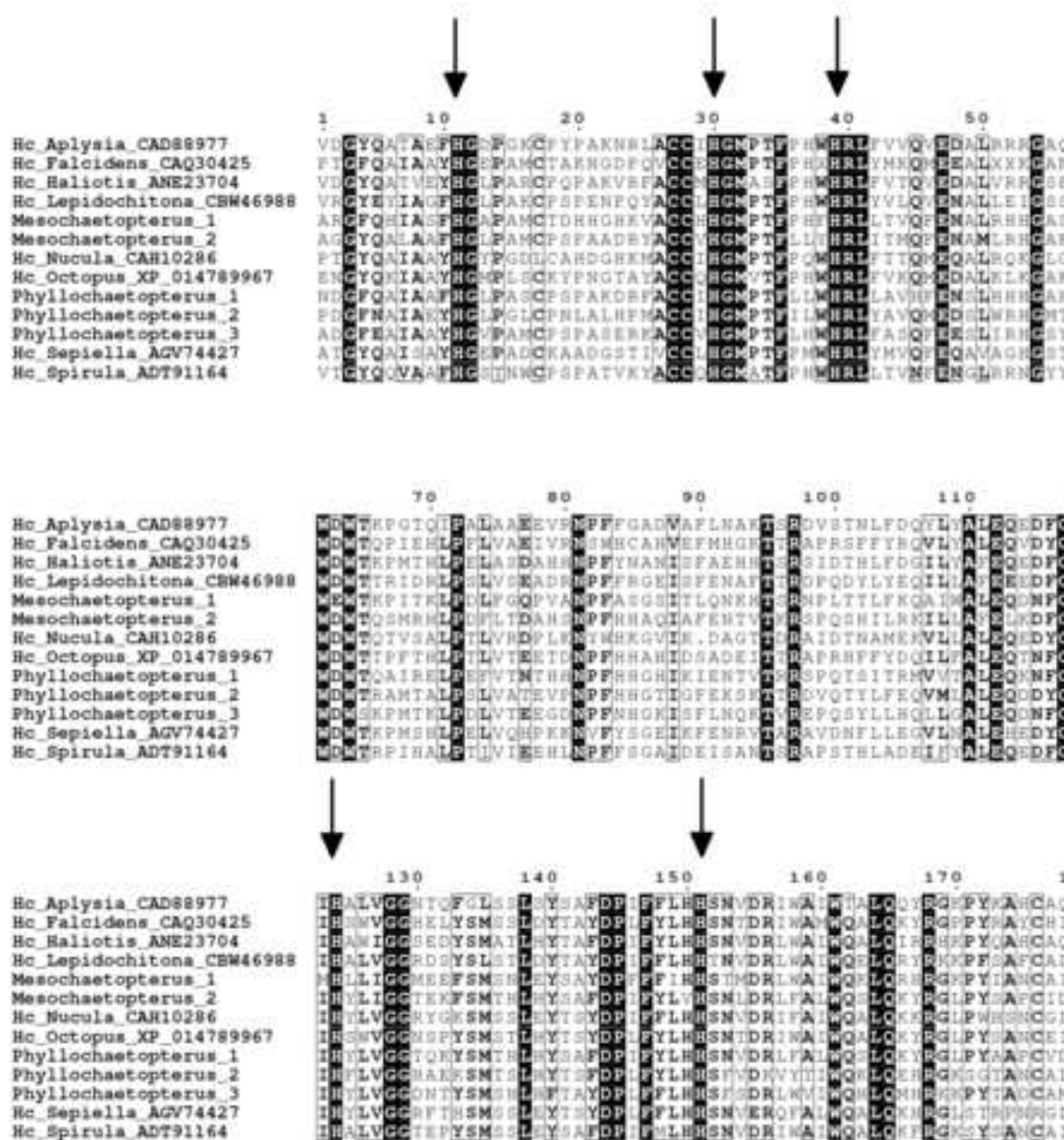
549

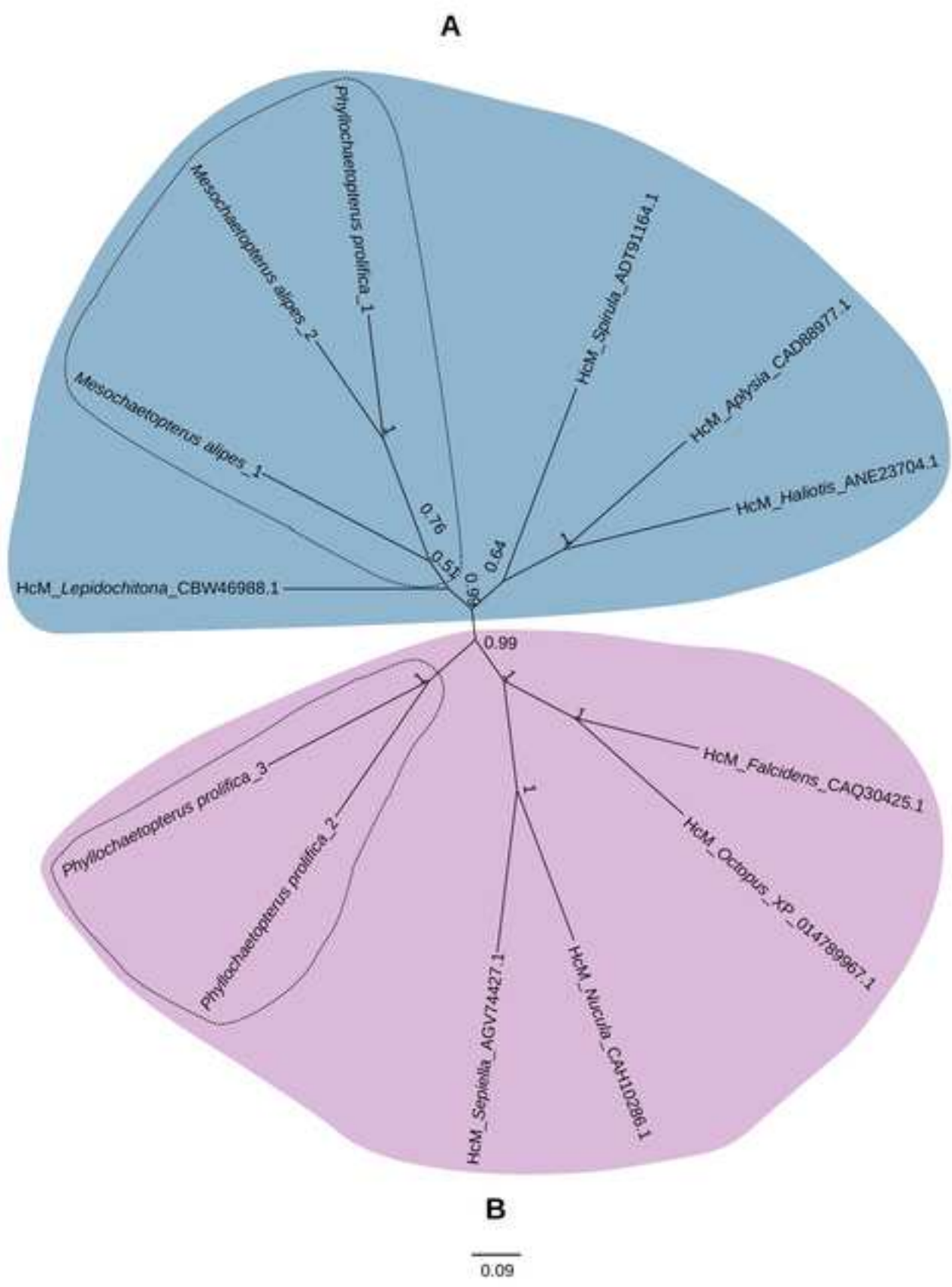
550

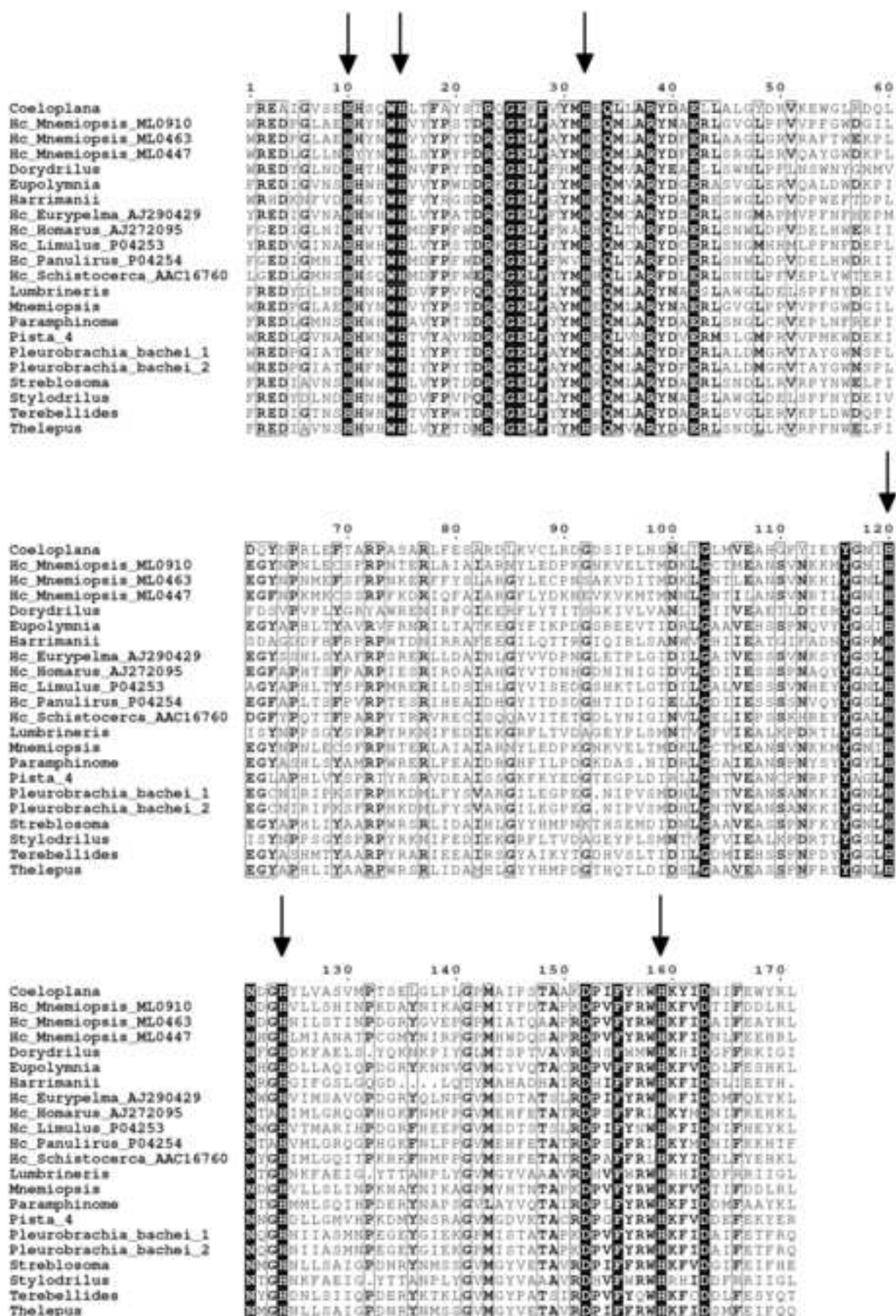
Figure



Figure

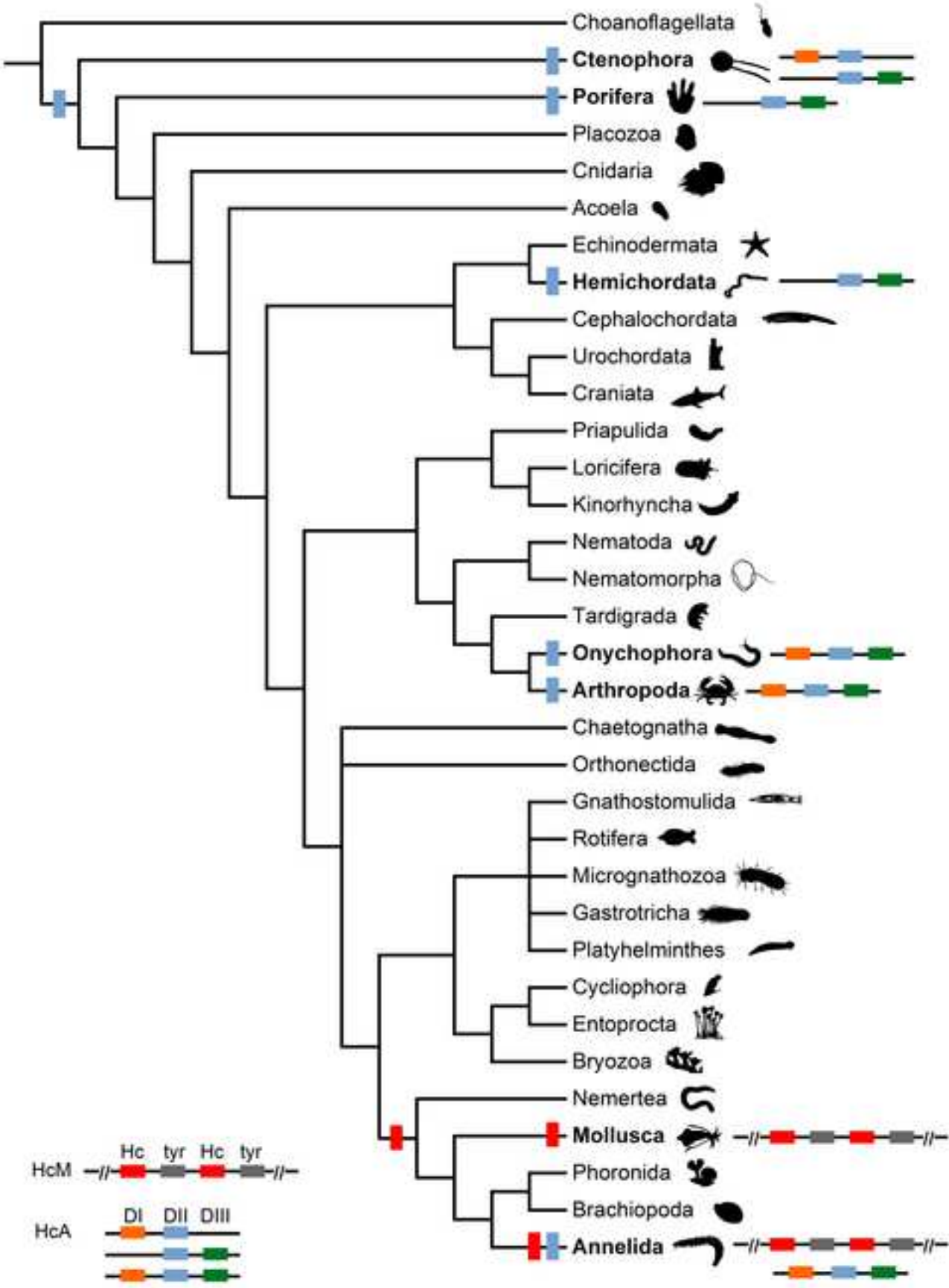






Figure

[Click here to download Figure Figure6.jpg](#)



Appendix 1. List of all taxa analyzed and total number of contigs after assembly. For undelined taxa, number and type of putative Hc genes and accession numbers are also provided.

Taxon	Total contigs number	Number and type of putative Hc genes	Accession number
CHOANOFLAGELATA			
<i>Acanthoecca spectabilis</i> W.Ellis, 1930	198,922		
<i>Salpingoeca pyxidium</i> Kent	202,399		
METAZOA			
Porifera			
<i>Hyalonema populiferum</i> Schulze, 1899	58,839		
<i>Kirkpatrickia variolosa</i> (Kirkpatrick, 1907)	100,231	1 partial HcA Domains II + III	MF998096
<i>Latrunculia apicalis</i> Ridley & Dendy, 1886	76,210	1 partial HcA Domains II + III	MF998097
<i>Rossella fibulata</i> Schulze & Kirkpatrick, 1910	40,103		
<i>Sympagella nux</i> Schmidt, 1870	85,237		
Ctenophora			
<i>Beroe abyssicola</i> Mortensen, 1927	83,798		
<i>Coeloplana astericola</i> Mortensen, 1927	222,614	1 partial HcA Domains I + II	MF998091
<i>Dryodora glandiformis</i> (Mertens, 1833)	101,598		
<i>Euplokamis dunlapae</i> Mills, 1987	321,550		
<i>Mnemiopsis leidyi</i> A. Agassiz, 1865	385,798	1 partial HcA Domains II + III	MF998101
<i>Pleurobrachia bachei</i> A. Agassiz, 1860	38,856	2 partial HcA Domains II + III	MF998107 MF998108
<i>Vallicula multiformis</i> Rankin, 1956	339,814		
Cnidaria			
<i>Gersemia antarctica</i> (Kukenthal, 1902)	20,023		
<i>Periphylla periphylla</i> (Peron & Lesueur, 1810)	212,658		
Staurozoa gen. sp.	45,023		
Echinodermata			
<i>Apostichopus californicus</i> (Stimpson, 1857)	134,640		
<i>Astrotopmma agassizii</i> Lyman, 1875	156,062		
<i>Labidiaster annulatus</i> Sladen, 1889	108,871		

<i>Labidiaster</i> sp.	168,720		
<i>Leptosynapta clarki</i> Heding, 1928	242,126		
Hemichordata			
<i>Balanoglossus aurantiaca</i> Girard, 1853	143,815		
<i>Cephalodiscus gracilis</i> Harmer, 1905	57,139		
<i>Cephalodiscus hodgsoni</i> Ridewood, 1907	200,052		
<i>Cephalodiscus nigrescens</i> Lankester, 1905	11,565		
<u>Harrimaniidae gen sp. (from Iceland)</u>	230,054	1 partial HcA Domains II + III	MF998095
Harrimaniidae gen sp. (from Norway)	274,434		
<i>Ptychodera bahamensis</i> Spengel, 1893	115,310		
<i>Rhabdopleura</i> sp.	4,790		
<i>Saccoglossus mereschkowskii</i> Wagner, 1885	145,937		
<i>Schizocardium brasiliense</i> Spengel, 1893	101,493		
<i>Stereobalanus canadensis</i> Spengel, 1893	12,741		
Torquaratoridae gen. sp.	102,971		
Annelida			
<i>Abarenicola pacifica</i> Healy & Wells, 1959	94,376		
<i>Aeolosoma</i> sp.	190,647		
<i>Aglaophamus verrilli</i> (McIntosh, 1885)	118,343		
<i>Alciopa</i> sp.	233,051		
<i>Amphisamytha galapagensis</i> Zottoli, 1983	14,313		
<i>Ancistrosyllis groenlandica</i> McIntosh, 1878	94,924		
<i>Andiorrhinus</i> sp.	139,858		
<i>Ankyrodrilus legaeus</i> Holt, 1965	54,246		
<i>Antarctodrilus proboscidea</i> (Brinkhurst & Fulton, 1979)	49,656		
<i>Aphelochaeta</i> sp.	165,566		
<i>Aphrodita japonica</i> Marenzeller, 1879	120,025		
<i>Arabella</i> sp.	217,183		
<i>Areco reco</i> Righi, Ayres & Bittencourt, 1978	170,510		
<i>Arenicola loveni</i> Kinberg, 1866	27,028		
<i>Arhynchite pugettensis</i> Fisher, 1949	20,724		
<i>Arichlidon gathofi</i> Watson Russell, 2000	140,980		
<i>Aricidea quadrilobata</i> Webster & Benedict, 1887	81,139		

<i>Armandia</i> sp.	137,440		
<i>Aspidosiphon laevis</i> Quatrefages, 1865	168,072		
<i>Auchenoplax crinita</i> Ehlers, 1887	144,974		
<i>Aulodrilus japonicus</i> Yamaguchi, 1953	109,361		
<i>Autolytus tuberculatus</i> (Schmarda, 1861)	137,934		
<i>Axiothella rubrocincta</i> (Johnson, 1901)	107,215		
<i>Bathydrilus rohdei</i> (Jamieson, 1977)	226,538		
<i>Bdellodrilus illuminatus</i> (Moore, 1894)	67,562		
<i>Bhawania goodei</i> Webster, 1884	70,615		
<i>Boccardia proboscidea</i> Hartman, 1940	117,570		
<i>Cambarincola holti</i> Hoffman, 1963	46,015		
<i>Capilloventer</i> sp.	221,627		
<i>Chaetogaster diaphanus</i> (Gruithuisen, 1828)	128,034		
<i>Chaetopterus variopedatus</i> (Renier, 1804)	147,132		
<i>Chaetacanthus magnificus</i> (Grube, 1876)	95,443		
<i>Chaetozone</i> sp.	143,597		
<i>Chloeia pinnata</i> Moore, 1911	130,037		
<i>Chone</i> sp.	106,577		
<i>Cirratulus spectabilis</i> (Kinberg, 1866)	120,244		
<i>Clymenella torquata</i> (Leidy, 1855)	111,567		
<i>Cossura longocirrata</i> Webster & Benedict, 1887	75,079		
<i>Crucigera zygophora</i> (Johnson, 1901)	116,092		
<i>Dichogaster saliens</i> (Beddard 1893)	98,665		

<i>Diopatra cuprea</i> (Bosc, 1802)	138,779		
<i>Dodecaceria pulchra</i> Day, 1955	229,501		
<u><i>Dorydrilus michaelsoni</i> Piguet, 1913</u>	136,096	1 partial HcA Domains II + III	MF998093
<i>Eteone</i> sp.	41,912		
<i>Enchytraeus crypticus</i> Westheide & Graefe, 1992	161,842		
<i>Eulalia myriacyclum</i> (Schmarda, 1861)	110,762		
<i>Eunice norvegica</i> (Linnaeus, 1767)	122,784		
<i>Euphrosine capensis</i> Kinberg, 1857	72,220		
<u><i>Eupolymnia nebulosa</i> (Montagu, 1819)</u>	139,021	1 partial HcA Domains II + III	MF998094
<i>Galathowenia oculata</i> (Zachs, 1923)	179,612		
<i>Galeolaria caespitosa</i> Lamarck, 1818	143,655		
<i>Gatesona chaetophora</i> (Bouché, 1972)	104,334		
<i>Glossodrilus</i> sp.	122,993		
<i>Glycera americana</i> Leidy, 1855	126,229		
<i>Glycera dibranchiata</i> Ehlers, 1868	101,455		
<i>Glycinde armigera</i> Moore, 1911	79,528		
<i>Glyptonotobdella antarctica</i> (Sawyer & White, 1969)	64,208		
<i>Goniada brunnea</i> Treadwell, 1906	89,398		
<i>Harmothoe oculinarum</i> (Storm, 1879)	94,991		
<i>Hemipodia simplex</i> (Grube, 1857)	55,653		
<i>Hermenia verruculosa</i> Grube, 1856	111,026		
<i>Hermodice carunculata</i> (Pallas, 1766)	110,813		
<i>Heteromastus filiformis</i> (Claparède, 1864)	148,196		
<i>Histriobdella homari</i> Beneden, 1858	143,130		
<i>Idanthursus</i> sp.	201,049		
<i>Laetmonice producta</i> Grube, 1876	73,530		

<i>Leanira</i> sp.	115,908		
<i>Lumbrineris crassicephala</i> Hartman, 1965	196,426		
<u><i>Lumbrineris perkinsi</i> Carrera-Parra, 2001</u>	144,648	1 partial HcA Domains II + III	MF998098
<i>Lysilla</i> sp.	104,324		
<i>Magelona berkeleyi</i> Jones, 1971	50,123		
<i>Marphysa sanguinea</i> (Montagu, 1813)	110,924		
<i>Melinna maculata</i> Webster, 1879	135,712		
<u><i>Mesochaetopterus alipes</i> Monroe, 1928</u>	83,209	2 HcM	MF998099 MF998100
<i>Microphthalmus similis</i> Bobretzky, 1870	169,427		
<i>Myxicola infundibulum</i> (Montagu, 1808)	217,996		
<i>Naineris laevigata</i> (Grube, 1855)	218,272		
<i>Neosabellaria cementarium</i> (Moore, 1906)	82,479		
<i>Nephasoma flagriferum</i> (Selenka, 1885)	170,216		
<i>Nephtys incisa</i> Malmgren, 1865	188,338		
<i>Nicolea macrobranchia</i> (Schmarda, 1861)	53,572		
<i>Nicomache venticola</i> Blake & Hilbig, 1990	124,708		
<i>Notomastus tenuis</i> Moore, 1909	129,745		
<i>Odontosyllis gibba</i> Claparède, 1863	131,487		
<i>Oenone fulgida</i> (Savigny in Lamarck, 1818)	144,726		
<i>Ophelina acuminata</i> Örsted, 1843	81,846		
<i>Ophiodromus pugettensis</i> (Johnson, 1901)	92,341		
<i>Ophryotrocha globopalpata</i> Blake & Hilbig, 1990	129,450		
<i>Owenia fusiformis</i> Delle Chiaje, 1844	106,476		
<i>Palola</i> sp.	211,279		
<i>Paralvinella palmiformis</i> Desbruyères & Laubier, 1986	85,363		
<u><i>Paramphinome jeffreysii</i> (McIntosh, 1868)</u>	165,337	1 partial HcA Domains II + III	MF998102
<i>Pectinaria gouldii</i> (Verrill, 1874)	92,091		
<i>Phascolosoma agassizii</i> Keferstein, 1866	87,403		
<i>Pherecardia striata</i> (Kinberg, 1857)	216,722		
<u><i>Phyllochaetopterus prolifica</i> Potts, 1914</u>	193,836	3 HcM	MF998103 MF998104 MF998105

<u><i>Pista macrolobata</i> Hessle, 1917</u>	126,764	1 partial HcA Domains II + III	MF998106
<i>Prionospio dubia</i> Day, 1961	119,949		
<i>Sabaco elongatus</i> (Verrill, 1873)	84,082		
<i>Schizobranchia insignis</i> Bush, 1905	102,002		
<i>Sclerolinum brattstromi</i> Webb, 1964	149,694		
<i>Scolelepis squamata</i> (Müller, 1806)	147,343		
<i>Serpula vermicularis</i> Linnaeus, 1767	151,097		
<i>Siboglinum ekmani</i> Jägersten, 1956	270,658		
<i>Siboglinum fiordicum</i> Webb, 1963	75,226		
<i>Sphaerodorum papillifer</i> Moore, 1909	52,411		
<i>Spirobranchus kraussii</i> (Baird, 1865)	167,761		
<i>Sternaspis scutata</i> (Ranzani, 1817)	10,634		
<i>Sternaspis</i> sp.	10,878		
<u><i>Streblosoma hartmanae</i> Kritzler, 1971</u>	108,080	1 HcA Domains I + II + III	MF998109
<u><i>Stylodrilus heringianus</i> Claparede, 1862</u>	239,935	1 partial HcA Domains II + III	MF998110
<i>Stygocapitella subterranea</i> Knöllner, 1934	74,556		
<i>Syllis</i> cf. <i>hyalina</i> Grube, 1863	106,283		
<u><i>Terebellides stroemii</i> Sars, 1835</u>	169,760	1 partial HcA Domains II + III	MF998112
<i>Tharyx kirkegaardii</i> Blake, 1991	114,157		
<u><i>Thelepus crispus</i> Johnson, 1901</u>	67,478	1 partial HcA Domains II + III	MF998113
<i>Themiste pyroides</i> (Chamberlin, 1919)	88,157		
<i>Travisia brevis</i> Moore, 1923	69,827		
<i>Trypanosyllis</i> sp.	167,501		
Brachiopoda			

<i>Glottidia pyramidata</i> (Stimpson, 1860)	131,562		
<i>Hemithiris psittacea</i> (Gmelin, 1791)	103,581		
<i>Laqueus californicus</i> (Koch, 1848)	133,086		
<i>Macandrevia cranium</i> (O. F. Müller, 1776)	9,695		
<i>Novocrania anomala</i> (O. F. Müller, 1776)	117,369		
Phoronida			
<i>Phoronis psammophila</i> Cori, 1889	193,702		
<i>Phoronopsis harmeri</i> Pixell, 1912	283,821		
Nemertea			
<i>Malacobdella grossa</i> (Müller, 1779)	79,313		
<i>Paranemertes peregrina</i> Coe, 1901	99,203		
<i>Parborlasia corrugatus</i> (McIntosh, 1876)	911,662		
<i>Tubulanus polymorphus</i> Renier, 1804	109,120		
Bryozoa			
<i>Pectinatella magnifica</i> (Leidy, 1851)	191,465		
Cycliophora			
<i>Symbion americanus</i> Obst, Funch & Kristensen, 2006	135,725		
Entoprocta			
<i>Barentsia gracilis</i> M. Sars, 1835	146,310		
<i>Loxosoma pectinaricola</i> Franzen, 1962	144,339		
Platyhelminthes			
<i>Acipensericola petersoni</i> Bullard, Snyder, Jensen & Overstreet, 2008	152,140		
<i>Cardicola currani</i> Bullard & Overstreet, 2004	86,962		
<i>Cardicola palmeri</i> Bullard & Overstreet, 2004	52,837		
<i>Elaphrobates euzeti</i> Bullard & Overstreet, 2003	118,013		
<i>Elopicola</i> sp.	64,384		
<i>Hapalorhynchus</i> sp.	42,863		
<i>Myliobaticola richardheardi</i> Bullard & Jensen, 2008	15,147		
<i>Myliobaticola</i> sp.	73,883		
<i>Psettarium anthicum</i> Bullard & Overstreet, 2006	39,616		
<i>Sanguinicola</i> sp.	145,041		

Selachohemecus olsoni Short, 1954

135,169

Orthonectida

Orthonectida gen. sp.

231,032

Priapulida

Priapulid sp.

50,034

Appendix 2: Queries used to search the assembled translated transcriptomes. All HcM sequences were also included in the dataset previous to the alignment.

Taxon	Protein	GenBank accession number
<u>Arthropoda</u>		
<i>Archispirosreptus gigas</i>	Hc subunit type I	CCC55877.1
<i>Cherax quadricarinatus</i>	Hc	AFP23115.1
<i>Cupiennius salei</i>	Hc subunit 1	CAC44749.1
<i>Cyamus scammoni</i>	Hc	ABB59715.1
<i>Limulus polyphemus</i>	Hc II	NP_001301072.1
<i>Macrobrachium nipponense</i>	Hc	AHJ90473.1
<i>Nebalia kensleyi</i>	Hc	ACV33306.1
<i>Penaeus monodon</i>	Hc	AEB77775.1
<i>Periplaneta americana</i>	Hc subunit 1 precursor	CAR85701.1
<i>Scutigera coleoptrata</i>	Hc subunit A	CAC69246.1
<i>Zootermopsis nevadensis</i>	Hc	KDR21641.1
<u>Mollusca</u>		
<i>Aplysia californica</i>	Hc	CAD88977.1
<i>Falciidens crossotus</i>	Hc fgh, partial	CAQ30425.1
<i>Haliotis rubra</i>	Hc type 1	ANE23704.1
<i>Lepidochitona cinerea</i>	Hc, partial	CBW46988.1
<i>Nucula nucleus</i>	Hc isoform 1	CAH10286.1
<i>Octopus bimaculoides</i>	Hc units G and H-like, partial	XP_014789967
<i>Sepiella maindroni</i>	Hc	AGV74427.1
<i>Spirula spirula</i>	Hc, partial	ADT91164.1

Taxon	Protein	GenBank accession number
<u>Porifera</u>		
<i>Amphimedon queenslandica</i>	Phenoloxidase subunit 2-like	XP_003390261.1
<u>Ctenophora</u>		
<i>Mnemiopsis leidyi</i>	Hc	Contig ML0910 (Supplementary Data in Martín-Durán et al., 2013)
<u>Hemichordata</u>		
<i>Saccoglossus kowalevski</i>	Hc-like, partial	ACY92544

Appendix 3. HcA superfamily protein sequences used in Burmester (2001), Aguilera et al. (2013), and Martín-Durán et al. (2013) with genes accession numbers for each species.

Protein	Species	Accession number
Prophenoloxidase	<i>Penaeus monodon</i>	AAD45201
	<i>Pacifastacus leniusculus</i>	X83494
	<i>Tenebrio molitor</i>	AB020738
	<i>Bombyx mori</i>	BBA08368
	<i>Manduca sexta</i>	AAC05796
	<i>Drosophila melanogaster</i>	NP476812
	<i>Neobellieria bullata</i>	AAD45526
	<i>Galleria mellonella</i>	AAK64363
	<i>Anopheles gambiae</i>	AF004915
Hemocyanin Arthropoda	<i>Eurypelma californicum</i>	AJ290429
	<i>Limulus polyphemus</i>	P04253
	<i>Callinectes sapidus</i>	AAF64305
	<i>Cupiennius salei</i>	CAC44749
	<i>Penaeus semisulcatus</i>	AAM77690
	<i>Pacifastacus leniusculus</i>	AAM81357
	<i>Panaeus vannamei</i>	CAA57880
	<i>Panulirus interruptus</i>	P04254
	<i>Homarus americanus</i>	AJ272095

Non-Arthropoda	<i>Palinurus vulgaris</i>	CAC69243
	<i>Scutigera coletrata</i>	CAC69246
	<i>Amphimedon queenslandica</i>	XP003390261
	<i>Mnemiopsis leidyi</i>	Contig ML0910 (Supplementary Data in Martín-Durán et al., 2013)
	<i>Mnemiopsis leidyi</i>	Contig ML0463 (Supplementary Data in Martín-Durán et al., 2013)
	<i>Mnemiopsis leidyi</i>	Contig ML0447 (Supplementary Data in Martín-Durán et al., 2013)
	<i>Saccoglossus kowalevskii</i>	ACY92544
Cryptocyanin		
	<i>Cancer magister</i>	AF091261
Pseudo-hemocyanin		
	<i>Homarus americanus</i>	CAB38042
	<i>Homarus americanus</i>	CAB38043
	<i>Metacarcinus magister</i>	AAD09762
Hexamerin		
	<i>Locusta migratoria</i>	U74469
	<i>Drosophila melanogaster</i>	NP476624
	<i>Periplaneta americana</i>	AAB09629
	<i>Blaberus discoidalis</i>	AAA74579
	<i>Spodoptera litura</i>	CAB55603
	<i>Camponotus festinatus</i>	AJ251271
	<i>Apriona germari</i>	AAM44045
	<i>Plodia interpunctella</i>	AAK71136

<i>Bracon hebetor</i>	I25974
<i>Anopheles gambiae</i>	AF020870
<i>Anopheles merus</i>	AF020875
