Mitogenomics reveals phylogenetic relationships of Arcoida (Mollusca, Bivalvia) and multiple independent expansions and contractions in mitochondrial genome size

Lingfeng Kong, Yuanning Li, Kevin M. Kocot, Yi Yang, Lu Qi, Qi Li, Kenneth M. Halanych

PII: DOI: Reference:	S1055-7903(20)30129-9 https://doi.org/10.1016/j.ympev.2020.106857 YMPEV 106857
To appear in:	Molecular Phylogenetics and Evolution
Received Date:	7 September 2019
Revised Date:	9 April 2020
Accepted Date:	21 May 2020



Please cite this article as: Kong, L., Li, Y., Kocot, K.M., Yang, Y., Qi, L., Li, Q., Halanych, K.M., Mitogenomics reveals phylogenetic relationships of Arcoida (Mollusca, Bivalvia) and multiple independent expansions and contractions in mitochondrial genome size, *Molecular Phylogenetics and Evolution* (2020), doi: https://doi.org/10.1016/j.ympev.2020.106857

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Elsevier Inc. All rights reserved.

1	Mitogenomics reveals phylogenetic relationships of Arcoida (Mollusca, Bivalvia)
2	and multiple independent expansions and contractions in mitochondrial genome
3	size
4	
5	Lingfeng Kong ^{a,b,1*} , Yuanning Li ^{c,d,1} , Kevin M. Kocot ^e , Yi Yang ^a , Lu Qi ^a , Qi Li ^{a,b} ,
6	Kenneth M. Halanych ^{c*}
7	
8	^a Key Laboratory of Mariculture, Ministry of Education, Ocean University of China,
9	Qingdao 266003, China
10	^b Laboratory for Marine Fisheries Science and Food Production Processes, Qingdao
11	National Laboratory for Marine Science and Technology, Qingdao 266237, China
12	^c Department of Biological Sciences & Molette Biology Laboratory for Environmental
13	and Climate Change Studies, Auburn University, AL 36849, USA
14	^d Department of Ecology and Evolutionary Biology, Yale University, New Haven, CT
15	06520, USA ²
16	^e Department of Biological Sciences and Alabama Museum of Natural History, The
17	University of Alabama, Tuscaloosa, AL 35487, USA
18	
19	*Corresponding author at: Key Laboratory of Mariculture, Ministry of Education,
20	Ocean University of China, Qingdao 266003, China; Department of Biological
21	Sciences, Auburn University, AL 36849, USA.
22	E-mail address: klfaly@ouc.edu.cn (L. Kong), ken@auburn.edu (K.M. Halanych)
23	
24	¹ Contributed equally to this work.
25	² Current address
26	
27	
28	
29	

30 Abstract

31 Arcoida, comprising about 570 species of blood cockles, is an ecologically and 32 economically important lineage of bivalve molluscs. Current classification of arcoids 33 is largely based on morphology, which shows widespread homoplasy. Despite two 34 recent studies employing multi-locus analyses with broad sampling of Arcoida, 35 evolutionary relationships among major lineages remain controversial. Interestingly, 36 mitochondrial genomes of several ark shell species are 2-3 times larger than those found 37 in most bilaterians, and are among the largest bilaterian mitochondrial genomes 38 reported to date. These results highlight the need of detailed phylogenetic study to 39 explore evolutionary relationships within Arcoida so that the evolution of 40 mitochondrial genome size can be understood. To this end, we sequenced 17 41 mitochondrial genomes and compared them with publicly available data, including 42 those from other lineages of Arcoida with emphasis on the subclade Arcoidea species. 43 Our phylogenetic analyses indicate that Noetiidae, Cucullaeidae and Glycymerididae 44 are nested within a polyphyletic Arcidae. Moreover, we find multiple independent 45 expansions and potential contractions of mitochondrial genome size, suggesting that 46 the large mitochondrial genome is not a shared ancestral feature in Arcoida. We also 47 examined tandem repeats and inverted repeats in non-coding regions and investigated 48 the presence of such repeats with relation to genome size variation. Our results suggest 49 that tandem repeats might facilitate intraspecific mitochondrial genome size variation, 50 and that inverted repeats, which could be derived from transposons, might be 51 responsible for mitochondrial genome expansions and contractions. We show that 52 mitochondrial genome size in Arcoida is more dynamic than previously understood and 53 provide insights into evolution of mitochondrial genome size variation in metazoans.

Keywords: genome size, mitogenome, tandem repeats, inverted repeats

- 54
- 55

56

- 57
- 58

59 1. Introduction

60 Ark shells or blood cockles (order Arcoida Gray, 1854) are a well-known, economically important group of bivalves. The name blood cockle comes from the 61 62 presence of hemoglobin, a feature found in all Arcoidea but not exclusive to them 63 among bivalves (Oliver and Holmes, 2006). Ark shells have a long fossil record, dating 64 back to the Lower Ordovician (~450 Mya; Cope, 1997). Today, ark shells are globally 65 distributed with approximately 570 species (Huber, 2010). They occur predominantly 66 in shallow tropical waters and warm temperate seas, and are most diverse in the Indo-67 West Pacific, where they play important roles in the ecosystem (Nakamura, 2005). Ark 68 shell bivalves are perhaps best known for use as Akagai, red clam or surf clam sushi. 69 World fishery production of clams, cockles, and ark shells is estimated to be 591,000 70 tons per year, which is worth nearly USD \$600,000,000 (FAO, 2018). Production of 71 ark shells in China alone was more than 360,000 tons in 2016 (Chinese National 72 Department of Fisheries, 2017). Because of their economic and ecological importance, 73 ark shells have been the subject of considerable research efforts (Sturmer et al., 2009; 74 Morton and Puljas, 2016; Li et al., 2016; Hou et al., 2016).

75 The taxonomy of Arcoida has primarily been based on shell morphology and palaeontological records. Currently, two superfamilies are recognized: Arcoidea and 76 77 Limopsoidea. Arcoidea encompasses Arcidae, Noetiidae, Cucullaeidae and 78 Glycymerididae, whereas Limopsidae and Philobryidae are included in Limopsoidea 79 (Oliver and Holmes, 2006). However, current classification of families relies on a 80 limited number of morphological characters. In addition, morphological characters 81 such as ligament and byssus arrangement, which have been used within genera, are 82 homoplastic (reviewed in Oliver and Holmes, 2006). Although Arcoida is well-83 established as a clade (e.g., Steiner and Hammer, 2000; Giribet and Wheeler, 2005; 84 Matsumoto, 2003; Bieler et al., 2014; Combosch et al., 2017), its internal relationships 85 remain controversial and poorly understood. Recent multi-locus studies (Feng et al., 86 2015; Combosch and Giribet, 2016) have improved our understanding of evolutionary relationships within Arcoida, but placement of some important lineages remains 87

contentious. In particular, the precise phylogenetic position of Glycymerididae remains unresolved (Combosch and Giribet, 2016). Moreover, Arcidae, which is divided into two subfamilies Arcinae and Anadarinae, based on the strength of the byssus in the attached or free-living forms, e.g., epibyssate and endobyssate, respectively (Newell, 1969), remains a source of inconsistency in Arcoida classification and is in need of taxonomic revision (Feng et al., 2015; Combosch and Giribet, 2016).

94 Mitogenomics has proven useful in resolving phylogenetic relationships across a 95 wide range of metazoans (e.g. Miya et al., 2001; Osigus et al., 2013; Li et al., 2015; 96 Mikkelsen et al., 2018). Despite the evolutionary importance of Arcoida, genomic 97 resources from this clade are still scarce in comparison to other major bivalve clades. 98 Prior to this study, only seven complete mitochondrial genomes of ark shells were 99 available (Liu et al., 2013; Sun et al., 2015a, 2015b, 2015c; Sun et al., 2016; Feng et al., 100 2017) and the size of ark shell mitochondrial genomes are unusual. For example, the 101 mitochondrial genome of Scapharca broughtonii is 46,985 bp and Anadara vellicata's 102 mitochondrial genome is 34,147 bp in length, 2-3 times the 15-17 kb size typically 103 reported in other bilaterians (Boore, 1999). Much of this increase in size is due to the 104 presence of non-coding regions, which have the potential to influence mitochondrial 105 energetics and replication (Sun et al., 2016). Remarkably large mitochondrial genomes 106 are also found in other bivalves (e.g. sea scallop Placopecten magellanicus 31-41 kb, 107 Smith and Snyder, 2007; and Bryopa lata, Clavagellidae, which at least 31,969 bp long, 108 Williams et al., 2017). More than 50% of the mitochondrial genome of the scallop 109 *Placopecten magellanicus* is non-coding regions that contain tandemly organized, 110 perfectly repeated sequences and dispersed, imperfect members of repeat families. This 111 indicates that transposition may lead to its uncharacteristically large mitochondrial 112 genome (Smith and Snyder, 2007). However, the relationship between major types of 113 non-coding regions (e.g. tandem repeats, inverted repeats) to genome size variation has 114 not been investigated, and the mechanism of mitochondrial genome expansion of ark 115 shells as well as other bivalves is still not clear.

116 In this study, we sequenced 17 mitochondrial genomes from representatives of 117 major Arcoida lineages. This information combined with data available from GenBank 118 can serve as a comparative framework for understanding mitochondrial genome 119 evolution of ark shells. We examine three questions: (1) if the phylogenetic 120 relationships inferred from mitochondrial genomes will be consistent with that of the 121 (Feng et al., 2015; Combosch and Giribet, 2016) multi-locus datasets; (2) if the 122 mitochondrial genome of arcoids experiences only expansion; and (3) if there are 123 correlation between the mitochondrial genome size and repeats. In addition, we 124 estimated divergence times of major clades of Arcoida to date the mitochondrial 125 genome expansion or contraction events in Arcoida molluscs.

126

127 **2. Materials and Methods**

128 **2.1.** Specimen collection and mitochondrial genome sequencing

129 Specimen collection data are shown in Table 1. All specimens were preserved by 130 freezing at -80 °C or in 80-100% non-denatured ethanol following collection. Arca 131 zebra, Anadara transversa, and Lunarca ovalis were ordered from Gulf Specimen 132 Supply in Panacea, Florida. These animals were collected locally by the company, but exact latitude and longitude are not available. Total genomic DNA was extracted using 133 134 the DNeasy Blood & Tissue Kit (Qiagen) following manufacture's protocols. Sequencing libraries were constructed using Illumina TruSeqTM DNA Sample Prep Kit 135 136 (Illumina, San Diego, CA, USA). For Anadara transversa and Lunarca ovalis, sequencing of genomic DNA was performed by the Genomic Services Lab at the 137 138 Hudson Alpha Institute in Huntsville, Alabama using 2 X 125 paired-end on an Illumina 139 HiSeq 2500 platform (San Diego, California). For Arcopsis adamsi and Arca zebra, 140 sequencing was conducted using 2 X 150 paired-end on the Illumina HiSeq X by 141 Novogene Corporation (USA). All other samples were sequenced on a single lane of 142 an Illumina HiSeq X using 2 X 150 paired-end by Beijing Novogene Technology Co., 143 Ltd (China).

144

145 **2.2. Mitochondrial genome assemblies and annotation**

146 Raw sequence data for all 17 samples were trimmed using Trimmomatic 0.36 147 (Bolger et al., 2014) with the parameters "ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 148 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:15 MINLEN:36". Resulting clean 149 reads were assembled *de novo* using Ray 2.3.1 (Boisvert et al., 2010) with k-mer = 31 150 and SPAdes v3.11.1 with k-mer of 21, 33, 55, 77 and the -careful flag (Bankevich et al., 2012). To identify putative mitochondrial contigs, BLAST (Altschul et al., 1997) 151 152 was performed on each assembly using the previously published Scapharca 153 broughtonii mitochondrial genome (GenBank Accession AB729113, Liu et al., 2013). 154 Assemblies of mitochondrial genomes from different assemblers were evaluated using Quast 4.5 (Gurevich et al., 2013) based on genome size, number of genes were 155 156 recovered, GC content and level of completeness. For Anadara crebricostata, Barbatia 157 virescens, Scapharca inaequivalvis, Scapharca kagoshimensis, Tegillarca sp., however, 158 two or three partial mtDNA contigs were recovered. In an attempt to join these partial 159 contigs, Price 1.2 (Ruby et al., 2013) was employed to extend existing contigs by 160 iteratively adding raw sequence reads to the contig ends as appropriate using default 161 settings. In all cases, partial contigs were successfully bridged together into a single contig. For additional quality control, raw paired-end reads were mapped to the 162 163 corresponding draft mitochondrial genome using Bowtie 2.3.4.1 (Langmead and 164 Salzberg, 2012) and alignments were visualized in Integrative Genomics Viewer, IGV 165 2.4.10 (Robinson et al., 2011). Annotation of the protein-coding genes (PCGs), 2 ribosomal RNAs (rRNAs) and transfer RNAs (tRNAs) was conducted initially with the 166 167 MITOS web server (revision 999, Bernt et al., 2013), with default settings and the 168 invertebrate genetic code for mitochondria, followed by manual genome annotation in 169 Artemis (Rutherford et al., 2000).

Transcriptomic data for *Arca noae*, *Anadara trapezia* and *Neocardia* sp. were
download from NCBI SRA database and raw reads were assembled using Trinity v.
2.6.6 (Grabherr et al., 2011) with the "trimmomatic" flag. Mitochondrial protein-coding
genes were identified by TBLASTX (Camacho et al., 2009) using the recovered Ark

174 shell mitochondrial genomes above as queries. The full-length sequences of all 17 175 newly sequenced mitochondrial genomes can be accessed through GenBank (Table 1). 176 Mitochondrial protein-coding genes derived from transcriptomic data acquired from 177 NCBI SRA Database deposited figshare were to 178 (https://doi.org/10.6084/m9.figshare.9746348.v1).

179

180 **2.3. Identification of Tandem repeats, inverted repeats and transposable elements**

181 Tandem repeats, inverted repeats and transposable elements were searched for in 182 the intergenic non-coding regions (NCRs) of new 17 mitochondrial genomes and four 183 mitochondrial genomes from GenBank. Only one specimen was analyzed when 184 multiple mitochondrial genomes were available in the same species. The intergenic 185 regions of each mitochondrial genome were extracted using Artemis (Carver et al., 186 2012). Tandem repeats were identified using the TRF (tandem repeats finder) V 4.09187 (Benson, 1999) with default parameters. Because the same repeat will be detected by 188 TRF at various period sizes, only those with the best alignment scoring for each 189 mitochondrial genome were chosen for further analyses. The Palindrome program from 190 the EMBOSS package (Rice et al., 2000) was employed to search for closely spaced 191 inverted repeats with identical sequences. The following parameters were used: 192 minimum length of palindromes = 7; maximum length of palindromes = 100; maximum 193 gap between elements = 10. Transposable elements are often the hallmark of nuclear-194 derived insertions (Alverson et al., 2010). With this in mind, each mitochondrial 195 genome was searched against the Repbase repetitive element database (version 23.05, 196 Kohany et al., 2006). Statistical tests were conducted in R v3.5.0 (R Development Core 197 Team 2018). As the data of mitochondrial genomic features are not necessarily 198 normally distributed, nonparametric Spearman's p values were calculated for the 199 correlation coefficient.

200

201 **2.4. Phylogenetic analyses**

202 Thirty-four taxa were included in the phylogenetic analyses of which data for eight 203 additional Ark shell and six outgroup species were retrieved from GenBank (Tables 1 204 and 2). The outgroup species Mizuhopecten vessoensis and Argopecten irradians from 205 Pectinidae, Pinctada maxima from Pteriidae, Ostrea denselamellosa, Crassostrea 206 hongkongensis and Crassostrea gigas from Ostreidae were used based on data 207 availability and current understanding of bivalve evolutionary history (Gonzalez et al. 208 2015; Lemer et al. 2016). Given that Ark shells date to the Ordovician (Cope, 1997), 209 relationships were reconstructed based on amino acid sequences from twelve 210 mitochondrial protein genes (atp6, cox1, cox2, cox3, cob, nad1, nad2, nad3, nad4, 211 nad4l, nad5 and nad6). Each gene was individually aligned using MAFFT (Katoh and 212 Standley, 2013). The selected genes were then trimmed using the default setting in 213 Gblocks (Talavera and Castresana, 2007) to remove ambiguously aligned regions with 214 default parameters. Resulting alignments were concatenated into final supermatrix 215 using FASconCAT (Kück and Meusemann, 2010) for downstream phylogenetic 216 analysis.

Maximum-likelihood (ML) analyses were conducted using IQ-TREE 1.6.1 (Nguyen et al., 2015). Partition schemes and best-fit models were selected by PartitionFinder (Lanfear et al., 2012). Nodal support was assessed with 1000 replicates of ultrafast bootstrapping (-bb 1000).

221 Bayesian inference (BI) with the site-heterogeneous CAT-GTR substitution model 222 (Lartillot and Philippe, 2004) was employed using PhyloBayes MPI 1.6j (Lartillot et 223 al., 2013). The concatenated AA dataset was run with four parallel Markov chain Monte 224 Carlo (MCMC) chains for 20,000-30000 generations each. Burn-in was determined 225 based on trace plots of log likelihood scores as viewed with plot phylobayes traces.R 226 script (https://github.com/wrf/graphphylo). Chains were considered to have reached 227 convergence when the maxdiff value was less than 0.1 as measured by bpcomp (Lartillot et al., 2013), the rel_diff value was less than 0.3, and the effective sample size 228 229 was greater than 50 as measured by tracecomp (Lartillot et al., 2013). A 50% majority

rule consensus tree was computed with bpcomp and nodal support was estimated byposterior probability.

232

233 2.5. Ancestral mitogenome size reconstruction

Ancestral mitogenome sizes (with 95% confidence interval) were reconstructed using maximum likelihood under a Brownian motion model by the anc.ML function of the R package Phytools (Revell 2012) on the phylogeny inferred with the molecular clock analyses (see below). The species *Arca noae*, *Anadara trapezia* and *Neocardia* sp. were removed from the phylogeny for this analysis using the drop.tip function in the R package Ape (Paradis et al., 2004) since the mitogenome size is not available.

240

241 **2.6. Molecular clock analyses**

242 A Bayesian relaxed-clock approach was employed to estimate divergence times of 243 major clades of Arcoida using BEAST 2.5.0 (Bouckaert et al., 2014). The molecular 244 clock was estimated using the uncorrelated lognormal distribution clock model, random 245 starting trees and a Calibrated Yule model. Divergence times were calibrated with three 246 timepoints from the fossil record. The age of the last common ancestor of Bivalvia was 247 constrained with a uniform distribution prior between 520.5 and 530 Mya (Bieler et al., 248 2014). Arcida was constrained using a normal distribution prior spanning 471.8–488.6 249 Ma, based on Glyptarca serrata Cope, 1996 (Arenigian; Cope, 1997). The age of 250 Northwestern Pacific Anadara species was constrained using a normal distribution 251 prior with the means 138.3 Mya and standard deviations of 5 Mya, based on Anadara 252 ferruginea (Jaccard, 1869; Huber, 2010). We performed two independent BEAST runs 253 (different seeds) for 20 million generations and with a sampling frequency of 1000 254 generations. We combined both runs using LogCombiner and after checking for 255 convergence of the runs with Tracer, the first 30% of generations were discarded as the burn-in period. The maximum clade credibility tree with median ages and the 95% 256 257 credibility interval (CI) at each node was computed using TreeAnnotator.

258

259 **3. Results**

260 **3.1.** Mitochondrial genome assemblies and general features

Of the 17 new taxa, complete mitochondrial genomes were recovered from 15 (Table 1). The remaining two were nearly complete: *Arca zebra* was missing rrnS, *Barbatia virescens* was missing trnA, trnE and trnV. Mitochondrial genome sequences varied in size from 17,479 bp (*B. lima*) to 56,170 bp (*S. kagoshimensis*). Four mitochondrial genomes had sequence lengths more than 40 kb and three ranged from 30 kb to 40 kb (Table 1). The average coverage of mitochondrial genomes varied from 164-fold to >1,000-fold coverage.

Nucleotide composition of the mitochondrial genome was biased toward A and T for all species (Table S1). The A + T content ranged from 54.06% in *B. virescens* to 67.84% in *S. kagoshimensis*. GC-skew and AT-skew for a given strand were calculated as (G - C)/(G + C) and (A - T)/(A + T), respectively, with negative values in skewness meaning the coding strand is enriched for C or T. We found the values of the AT-skew were mostly negative, while values of the GC-skew were all positive (Table S1).

Six different start codons were observed (Table S2 & S3) but most protein-coding genes (50.98%) start with the codon ATG. Similar results were observed in mitochondrial genomes of anomalodesmatan bivalves (Williams et al., 2017). In most genes, complete stop codons were identified (TAG or TAA), but a truncated stop codon T, which may be modified to a complete TAA stop codon via posttranscriptional polyadenylation (Ojala et al., 1981), was found in *cox1-3, cob* and *nad3* (Table S2).

280 All genes are encoded on the same strand, as typically found in other bivalves 281 (Williams et al., 2017). All complete mitochondrial genomes sequenced herein are 282 composed of 12 protein-coding genes (all taxa lacked *atp8*), two ribosomal (rRNA) 283 genes, and 21-29 tRNAs. Gene order of protein-coding genes and rRNA genes (Fig. 2) 284 differed among taxa, but is generally conserved in closely related species. 285 Arrangements of tRNAs are highly positional variable across the arcoid mitochondrial genomes, but there are several conserved gene blocks within different lineages. The 286 287 tRNA complex trnP-trnI-trnG-trnE-trnV is present in all species of Scapharca,

Anadara, Potiarca, Tegillarca and Lunarca. There are two tRNA complexes (trnN trnD-trnQ-trnH and trnC-trnM-trnE-trnL1-trnW) shared by *A. zebra* and *A. navicularis*.

290 Of the protein-coding genes and rRNA genes, six possible cases of gene 291 duplication were observed: cox2 is duplicated in Anadara crebricostata, S. 292 inaequivalvis, S. kagoshimensis, and Tegillarca sp. and rrnS is duplicated in A. 293 transversa and Lunarca ovalis (Fig. 2). Two copies of cox2 in A. crebricostata, S. 294 inaequivalvis and S. kagoshimensis were sequentially organized and located 295 downstream of *cob*, while one copy of *cox2* in *Tegillarca* sp. is located on downstream 296 of *cob* and the other is downstream of *cox3*. Overall identity between the two *cox2* 297 sequences in each species where it is duplicated is low. However, for A. transversa and 298 L. ovalis, the sequences of the two copies of rrnS in each species are nearly identical. 299 Consistent with a previous study (Feng et al., 2017), the newly sequenced mitochondrial 300 genome of *Cucullaea labiata* revealed the presence of an intron within *cox1*. The intron 301 was 606 bp long and no open reading frames were found. Only two nucleotide bp were 302 found different of the overlap 606 bp when we aligned the intron sequences of C. 303 *labiata* from this and previous study.

304

305 **3.2.** Non-coding regions

306 The total length of intergenic non-coding region ranged from 3,013 bp (17.24% of 307 the genome) in B. lima to 41,101 bp (73.17% of the genome) in S. kagoshimensis (Table 308 S4). All sequences are AT-rich and the A + T content of the intergenic non-coding 309 regions in each mitochondrial genome is higher than that of the whole mitochondrial 310 genome except B. virescens, while the AT-skew for non-coding region had a negative 311 value in all species and showed a large variation (Table S1). A significant and strong 312 positive correlation is observed between mitochondrial genome size and proportion of 313 non-coding region ($\rho = 0.96$, P = 5.12e-6; Fig. 3A). Tandem repeats and inverted 314 repeats in intergenic non-coding regions show significant variation in number and 315 sequence lengths (Table S4). We found significantly positive correlations between

mitochondrial genome size and proportion of tandem repeats and proportion of inverted repeats ($\rho = 0.64$, P = 0.002, Fig. 3B; $\rho = 0.76$, P = 9.84e-5, Fig. 3C).

318

319 **3.3. Phylogenetic analyses**

320 The concatenated alignment of amino acid sequences from 34 taxa had a total 321 length of 4,866 positions. Both ML and BI analyses of the two concatenated datasets 322 yielded identical branching orders with high posterior probabilities (PP) and bootstrap 323 support values (BS, Fig. 4). Arcidae was found to be polyphyletic with three well-324 supported lineages. Three Arca species formed a lineage; the second lineage included 325 Barbatia virescens and Trisidos species; and the third enclosed the subfamily 326 Anadarinae and the sister taxon Barbatia lima. Northwestern Pacific Anadarinae 327 species which included Anadara, Scapharca, Potiarca and Tegillarca were sister to the 328 Gulf of Mexico Anadarinae species A. transversa and Lunarca ovalis. The only 329 representative of Noetiidae, Arcopsis adamsi, was found to nest within the polyphyletic 330 Arcidae as the sister taxon of the B. virescens/Trisidos clade. Arcoidea, which 331 encompasses Arcidae, Noetiidae, Cucullaeidae and Glycymerididae, forms a clade that 332 is not well-supported in the ML analysis (PP = 1; BS = 56%) and *Neocardia* sp. (Limopsoiea) is the sister taxon of Arcoidea (PP = 1; BS = 100%). C. labiata, which is 333 the only extant species of Cucullaeidae, was strongly supported as the sister taxon of 334 Glycymerididae (PP = 1; BS = 92%). 335

336

337 **3**.

3.4. Ancestral state reconstruction

Ancestral state reconstruction indicates that the evolution of the mitochondrial genome size experiences different trends across different arcoid lineages. When considering a difference of >4kb from their respective recent ancestral mitogenomes, mitochondrial genome expansion is apparent in Northwestern Pacific Anadarinae species and *Arca zebra*, but the opposite trend of genome contraction occurs in *Arca navicularis*, *Barbatia lima*, and the Gulf of Mexico Anadarinae species (Fig. 5). Notably, both mitochondrial genome expansion and contraction were found within the

Northwestern Pacific Anadarinae linage, e.g. *Tegillarca* sp. and a clade comprising *S*. *broughtonii*, *S. kagoshimensis* and *S. inaequivalvis* occur expansions, and *T. granosa*, *Potiarca pilula* experience contraction. A medium-sized mitochondrial genome (23.4 kb) was estimated to be the ancestral state of Arcoidea.

349

350 3.5. Diversification times

Bayesian inference with a relaxed molecular clock recovered the same topology 351 352 as ML and BI analyses with strong support (Fig. 6), with the exception of the placement of A. noae and A. navicularis. The dated topology shows that A. noae is the sister taxon 353 354 of the clade containing A. navicularis and A. zebra. Our analyses suggest size of 355 mitochondrial genome expansions in Northwestern Pacific Anadarinae species (34.2 356 kb) primarily originated approximately 158 MYA (95% highest posterior density 357 interval [HPD] 139.9-184.4), while mitochondrial genome contractions in the Gulf of 358 Mexico Anadarinae species (20.5 kb) happened at 58.9 MYA (95% HPD 15.6-132.3). 359 Within the Northwestern Pacific Anadarinae linage, a clade comprising S. broughtonii, 360 S. kagoshimensis and S. inaequivalvis experienced expansions (48.1 kb) at 61 MYA 361 (95% HPD 36.9-84.1). Moreover, mitochondrial genome expansion of A. zebra (44.7 kb) and contractions of A. navicularis (18.0 kb) occurred at 54.8 MYA (95% HPD 11.0-362 363 91.4).

364

365 **4. Discussion**

366 4.1. Arcoid phylogeny

Arcidae was recovered as a polyphyletic group, consistent with earlier studies (e.g. Marko, 2002; Matsumoto, 2003; Feng et al., 2015; Combosch and Giribet, 2016), whereas Anadarinae was recovered as monophyletic in both ML and BI analyses (Fig. 4), consistent with previous findings (Marko, 2002; Matsumoto, 2003; Feng et al., 2015; Combosch and Giribet, 2016). Within Anadarinae, *Anadara* and *Scapharca* were found to be polyphyletic which suggest that extensive taxonomy revise in the genus level are needed. Interestingly, Anadarinae was further split into two subclades based on

374 geographic locality, northwestern Pacific species and Mexico Gulf species. By 375 comparison, Arcinae was found to be non-monophyletic, with three well-supported 376 lineages. The three Arca species were recovered in one clade, sister to all other 377 Arcoidea species with high posterior probability (PP = 1). Morphologically, this clade 378 corresponds to the "noae" morphotype within the genus Arca of Oliver and Holmes 379 (2006). The second clade included Barbatia virescens and Trisidos species, shows 380 Trisidos and Barbatia shared a common ancestor. The other Barbatia lineage, which 381 was only represented by B. lima, was sister to the clade Anadarinae with high nodal 382 support, consistent with Combosch and Giribet (2016). Results from the present study 383 support raising these three lineages to the familial rank. However, confirmation from 384 future studies including missing genera such as Bathyarca and Bentharca is needed. 385 Moreover, reconsideration of morphological characters used in the traditional 386 taxonomy of the group in light of these results is needed as it is evident that many of 387 these characters exhibit homoplasy (Oliver and Holmes, 2006).

388 The phylogenetic position and taxonomic status of Glycymerididae has been 389 contentious and nesting Glycymerididae within Arcoidea or Limopsoidea remains 390 unresolved (see Feng et al., 2015; Combosch and Giribet, 2016 for details). In the 391 present study, both ML and BI phylogenetic analyses support placement of 392 Glycimerididae within Arcidea (Fig. 4), in agreement with previous studies based on 393 combinations of nuclear 18S rRNA, 28S rRNA, histone H3, mitochondrial 12S and 394 COI data (Feng et al., 2015), and corroborated the morphological analyses of Amler 395 (1999) and Oliver and Holmes (2006). Cucullaea labiata, the only extant species within 396 Cucullaeidae, was recovered as the sister taxon of a clade containing Bathyarca 397 glomerula and Arca boucardi (Combosch and Giribet, 2016), however, the 398 phylogenetic position of the clade containing the three species was not congruent 399 among parsimony, maximum likelihood and Bayesian inference analyses (Combosch 400 and Giribet, 2016). Although B. glomerula and A. boucardi were not included in our 401 study, our analyses support the phylogenetic position of C. labiata as the sister lineage 402 of Glycimerididae both ML and BI analysis. Arcopsis adamsi (Noetiidae) was

recovered as the sister taxon of the *B. virescens/Trisidos* clade, consistent with that of
Feng et al. (2015) and Combosch and Giribet (2016).

405

406 **4.2. Genome size**

407 We found evidence for multiple expansions and contractions of the mitochondrial 408 genome (Fig. 5), which implies that mitochondrial genome size evolution could have 409 been involved in the diversification process. Based on ancestral state reconstruction, 410 the large mitochondrial genome found in previous studies (Liu et al., 2013; Sun et al., 411 2015c) is not a shared ancestral feature of arcoids. Molecular clock analyses suggest 412 mitochondrial genome expansions in Northwestern Pacific Anadarinae species 413 primarily originated approximately 158 million years ago, but mitochondrial genome 414 occurred contraction in the Gulf of Mexico Anadarinae species about 59 million years 415 ago. Interestingly, Arca zebra, a Gulf of Mexico Arcinae species has undergone an 416 expansion event while A. navicularis, a Northwestern Pacific Arcinae species 417 underwent contraction about 55 million years ago (Fig. 5, Fig. 6), which indicates that 418 mitochondrial genome expansion and contraction are independent evolutionary events. 419 Despite the fact that the mitochondrial genome does not have recombination and therefore only represents the history of one marker, we stress that conclusions from 420 421 mitochondrial sequences regarding divergence times are comparable with those from 422 concatenated nuclear or nuclear and mitochondrial gene sequences data (Bieler et al., 423 2014; Combosch and Giribet, 2016). The profound climatic and ecological changes that 424 occurred during this period (e.g., Cretaceous-Paleogene mass extinction and 425 Paleocene–Eocene Thermal Maximum) may have contributed to the adaptive radiation 426 of these arcoids to new niches, which then influence the mitochondrial genome 427 evolution process. Animal mitochondrial genomes have purportedly been under 428 selection for small genome size and are generally devoid of noncoding sequences to 429 improve replication and translation efficiency (Rand, 1993). However, several studies 430 show that low metabolic rates and limited locomotive ability are correlated with weak 431 selective constraints on mitochondrial genes (e.g., Chong and Mueller, 2013; Sun et al.,

432 2017). Up to now, all remarkably large mitochondrial genomes found in molluscs are 433 from bivalves, a group that is primarily sedentary, except the limpet *Lottia digitalis*, of 434 which the mitochondrial genome size is 25.4 kb (Simisiom et al., 2006). We propose 435 that the large mitochondrial genome size in these bivalves may be correlated with their 436 metabolic rates (Strotz et al., 2018). If the metabolic rate is low, bivalves perhaps 437 experienced weak purifying selection on their long intergenic regions.

438 As seen in Fig. 3A, the size variation of mitochondrial genomes is mainly due to 439 the different lengths of the NCRs. Several characteristics were revealed when the NCRs 440 were examined in detail. Firstly, the distribution of the NCRs was variable among 441 different linages (Fig. 2). The two large non-coding regions in A. zebra located between 442 cob and rrnL, rrnL and nad3, respectively. In the mitochondrial genomes of S. 443 broughtonii, S. kagoshimensis, S. inaequivalvis and Anadara crebricostata, the two 444 largest NCRs were observed between cox2 and nad6, nad2 and cox1. Tegillarca sp. has 445 three main non-coding regions which located separately between nad2 and cox1, cox3 446 and cox2, nad4l and rrnS. However, S. globosa, A. vellicata, T. nodifera and T. granosa 447 has only one large concentrated non-coding region, located between *nad2* and *cox1*. 448 Although less remarkable in size, several cases of major NCRs have previously been 449 reported in other bivalves mitogenomes. However, the intergenic location of these 450 NCRs does not coincide with that found in ark shells. Major NCRs were described 451 between nad4 and nad1 in three scallops (Pectinidae; Wu et al., 2009), and between 452 atp6 and nad2 in Crassostrea iredalei (Ostreidae; Wu et al., 2010). The different 453 distribution of NCRs in ark shells and close related taxa strongly suggests multiple 454 independent insertions instead of a single ancestral event with subsequent loss in other 455 lineages.

456 Secondly, among the 21 known mitochondrial genomes from ark shells, one to ten 457 tandem repeats were detected in sixteen species (Table S4). Significant correlations 458 were observed between mitochondrial genome size and proportion of each 459 mitochondrial genome made up of tandem repeats ($\rho = 0.64$, P = 0.002; Fig. 3B). Indeed, 460 the tandem repeats are potentially a main factor leading to the variation of intraspecific

461 mitochondrial genome size. Comparing our newly sequenced mitochondrial genome of 462 Cucullaea labiata to that of previous study (Feng et al., 2017), we found that the 463 difference of mitochondrial genome size between two individuals was due to the 464 number of copies of a 659 bp repeat (one copy vs. 8.9 copies). Similar results have been 465 found in other bivalves, including the sea scallop Placopecten magellanicus whose mitochondrial genome ranges in size from 30.7 kb to 41 kb among different individuals, 466 467 depending on the number of copies of a 1.4 kb repeat (Smith and Snyder, 2007). The 468 mitochondrial genomes of S. broughtonii ranges in size from about 47 kb to ~50 kb 469 because of variation in the number of tandem repeat arrays (Liu et al., 2013).

470 Thirdly, a systematic search detected inverted repeats (ranging from 1 to 83 copies) 471 in all samples arcoid species. Strong correlations between presence of inverted repeats 472 and the genome size is observed ($\rho = 0.76$, P = 9.84e-5; Fig. 3C), which indicates the 473 inverted repeats might be the major elements responsible for mitochondrial genome 474 size variation, especially accounting for the interspecific mitochondrial genome size. 475 Inverted repeats are not common in bilaterian mitochondrial genomes, but they have 476 been found previously in fungi, algae, plants and demosponges (see Lavrov, 2010 and 477 references therein). This rapid rate of proliferation combined with the broad phylogenetic distribution of hairpin elements was suggested as a defining force in the 478 479 evolution of mitochondrial genomes of demosponges (Lavrov, 2010). However, the 480 origin and propagation mechanism is still largely unknown. Lavrov (2010) suggested 481 that the occurrence of stem-loop elements in the mitochondrial genome of the sponge 482 Lubomirskia baicalensis may be linked to the activity of the Baikalum-1 transposon in 483 its nuclear genome. Notably, we found numerous fragments of transposable elements 484 in the non-coding region of arcoid mitochondrial genomes when they were blasted 485 against the Repbase repetitive element database (Table S4). Previously, two 486 transposon-like noncanonical open reading frames were also found to integrate into the 487 mitochondrial genome of a deep-sea anemone Bolocera sp. (Zhang et al., 2017). Further 488 studies are needed to test this association and the molecular mechanisms involved.

489

490 **4.3. Genome composition and structure**

491 Similar to other bivalves, arcoids show highly variable mitochondrial gene order, 492 even within the same genus. All arcoid mitochondrial genomes sampled to date possess 493 substantially different mitochondrial gene orders (Fig. 2). When tRNAs were not 494 considered, there are still twelve different gene orders for 21 species. Coding on both 495 strands has been assumed to restrain the gene arrangement of mitochondrial genomes, 496 because rearranging a genome with dual-strand coding may be more complicated than 497 if all genes are on the same strand (Ren et al., 2010). Perhaps coding on only one strand 498 may be one of the factors behind the extensive rearrangement of bivalve mitochondrial 499 genomes. Although variability in gene arrangement is high, gene order of protein-500 coding genes and ribosomal genes is generally conserved within each phylogenetic 501 clade. Interestingly, Tegillarca granosa/T. nodifera clade share same mitochondrial 502 gene order with A. vellicata/S. globosa/A. crebricostata/S. inaequivalvis/S. 503 kagoshimensis/S. broughtonii clade, which might imply synapomorphy of identical 504 PCG/rRNA clusters.

505 Six putative cases of gene duplication, especially in *cox2* and *rrn*, were observed 506 in the 17 mitochondrial genome orders found in ark shells. By revisiting the annotation 507 of mtDNA genomes of Scapharca broughtonii (AB729113, Liu et al., 2013), A. sativa 508 (= S. broughtonii, KF667521, Hou et al., 2016) and S. kagoshimensis (KF750628, Sun 509 et al., 2015a), we found the existence of a putative duplication of the cox2 gene. 510 Duplication of *cox2* was detected in other bivalve taxa such as *Meretrix*, *Loripes* and 511 Venerupis within Heterodonta, hence duplicated cox2 was considered to be 512 synapomorphic of Heterodonta (Stöger and Schrödl, 2013). In addition, cox2 is 513 duplicated in tandem in the mitochondrial genomes of *Ruditapes philippinarum* and 514 Musculista senhousia, two bivalves with doubly uniparental inheritance (DUI) of 515 mtDNA (Breton et al., 2014). However, to our knowledge, the blood cockles do not 516 have DUI, hence their duplicate cox2 cannot be at all related to this mechanism and its 517 functionality is still under study. Duplicated rrnS has been also reported as a

synapomorphy for all *Crassostrea* species except *C. virginica*, which is the sister taxon
to other *Crassostrea* species (Wu et al., 2010).

520 To the best of our knowledge, C. labiata is the only mollusk species to contain an 521 intron in a mitochondrial protein-coding gene. Introns rarely occur in animal 522 mitochondrial genomes. However, Group I introns have been found in Placozoa, 523 Anthozoa (corals, soft corals, sea anemones) and Porifera (sponges), while Group II 524 introns have been found in Placozoa and Annelida (Vallès et al., 2008; Huchon et al., 525 2015; Bernardino et al., 2017). Although ORFs were not found in the cox1 introns of 526 C. labiata, the 606 bp-intron region could possibly be derived from ancient transposable 527 elements that have since lost any function (Bernardino et al., 2017). Notably, based on 528 the result of Palindrome program from the EMBOSS package, 18 bp palindromic motifs 529 were found in the intron. The palindrome is thought to provide the site for DNA-binding 530 proteins involved in the transcriptional machinery (Arunkumar and Nagaraju, 2006). In 531 addition, the conventional splicing signals (GT-AG) of a nuclear intron (Breathnach et 532 al., 1978) was found at the *cox1* exon-intron boundaries by aligning mitochondrial 533 genome of *C. labiata* with the GenBank published one (KP091889.1, Feng et al., 2017). 534 Thus, we speculate that the GT-AG motif in the 5' and 3' end of the intron may 535 represent a splicing signal involved the transcription process of the *cox1*.

536 Here, we show that mitochondrial genomes of some ark shells have large 537 insertions of non-coding regions that are unusual for bilaterian mitochondrial genomes, 538 leading several ark shell species to have some of the largest animal mitochondrial 539 genomes observed to date (~50kb). Moreover, our results suggest that the large 540 variation of mitochondrial genome size are the result of multiple independent events. 541 Tandem repeats likely facilitate mitochondrial genome size variation, and that inverted 542 repeats which could related to the transposition events might be responsible elements 543 for interspecific mitochondrial genome expansions and contraction in ark shells.

544

545 Acknowledgements

546 We thank Dr. Rudiger Bieler from Field Museum of Natural History for kindly 547 the sample of Arcopsis adamsi. Two anonymous reviewers providing 548 provided feedback that helped to improve this article. This study was supported by 549 awards from the National Natural Science Foundation of China (31772414, 41276138) 550 and the Fundamental Research Funds for the Central Universities (201964001) to L.K., 551 the Ocean University of China and Auburn University foundation/integrated grant program to K.M.H. and L.K. Lingfeng Kong was also supported by a one-year 552 553 scholarship from the China Scholarship Council (CSC) as a visiting scholar at Auburn 554 University. This is Molette Biology Laboratory contribution # and Auburn University 555 Marine Biology Program contribution #.

556

557 **References**

- Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W., Lipman,
 D.J., 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database
 search programs. Nucleic Acids Res. 25, 3389–402.
 https://doi.org/10.1093/nar/25.17.3389.
- Alverson, A.J., Wei, X., Rice, D.W., Stern, D.B., Barry, K., Palmer, J.D., 2010. Insights
 into the evolution of mitochondrial genome size from complete sequences of *Citrullus lanatus* and *Cucurbita pepo* (Cucurbitaceae). Mol. Biol. Evol. 27, 1436–
 1448. https://doi.org/10.1093/molbev/msq029.
- 566 Amler, M.R.W., 1999. Synoptical classification of fossil and Recent Bivalvia. Geol.
 567 Palaeontol. 33, 237–248.
- Arunkumar, K.P., Nagaraju, J., 2006. Unusually long palindromes are abundant in
 mitochondrial control regions of insects and nematodes. PLoS One. 1, e110.
 https://doi.org/10.1371/journal.pone.0000110.
- Bankevich, A., Nurk, S., Antipov, D., Gurevich, A.A., Dvorkin, M., Kulikov, A.S.,
 Lesin, V.M., Nikolenko, S.I., Pham, S., Prjibelski, A.D., Pyshkin, A. V, Sirotkin,
- 573 A. V, Vyahhi, N., Tesler, G., Alekseyev, M.A., Pevzner, P.A., 2012. SPAdes: A
- 574 New Genome Assembly Algorithm and Its Applications to Single-Cell

- 575
 Sequencing.
 J.
 Comput.
 Biol.
 19,
 455–477.

 576
 https://doi.org/10.1089/cmb.2012.0021.
 Image: Comput.
 Im
- 577 Benson, G., 1999. Tandem repeats finder: a program to analyze DNA sequences.
 578 Nucleic Acids Res 27, 573–580. https://doi.org/10.1093/nar/27.2.573.
- Bernardino, A.F., Li, Y., Smith, C.R., Halanych, K.M., 2017. Multiple introns in a
 deep-sea Annelid (Decemunciger: Ampharetidae) mitochondrial genome. Sci.
 Rep. 7, 4295. https://doi.org/10.1038/s41598-017-04094-w.
- Bernt, M., Donath, A., Jühling, F., Externbrink, F., Florentz, C., Fritzsch, G., Pütz, J.,
 Middendorf, M., Stadler, P.F., 2013. MITOS: Improved de novo metazoan
 mitochondrial genome annotation. Mol. Phylogenet. Evol. 69, 313–319.
 https://doi.org/10.1016/j.ympev.2012.08.023.
- 586 Bieler, R., Mikkelsen, P.M., Collins, T.M., Glover, E.A., González, V.L., Graf, D.L.,
- 587 Harper, E.M., Healy, J., Kawauchi, G.Y., Sharma, P.P., Staubach, S., Strong, E.E.,
- 588 Taylor, J.D., Tëmkin, I., Zardus, J.D., Clark, S., Guzmán, A., McIntyre, E., Sharp,
- P., Giribet, G., 2014. Investigating the Bivalve Tree of Life An exemplar-based
 approach combining molecular and novel morphological characters. Invertebr.
 Syst. 28, 32–115. https://doi.org/10.1071/IS13010.
- Boisvert, S., Laviolette, F., Corbeil, J., 2010. Ray: simultaneous assembly of reads from
 a mix of high-throughput sequencing technologies. J. Comput. Biol. 17, 1519–
 1533. https://doi.org/10.1089/cmb.2009.0238.
- Bolger, A.M., Lohse, M., Usadel, B., 2014. Trimmomatic: A flexible trimmer for
 Illumina sequence data. Bioinformatics 30, 2114–2120.
 https://doi.org/10.1093/bioinformatics/btu170.
- Boore, J.L., 1999. Animal mitochondrial genomes. Nucleic Acids Res. 27, 1767–1780.
 https://doi.org/10.1093/nar/27.8.1767.
- Bouckaert, R., Heled, J., Kühnert, D., Vaughan, T., Wu, C.H., Xie, D., Suchard, M.A.,
 Rambaut, A., Drummond, A.J., 2014. BEAST 2: A Software Platform for
 Bayesian Evolutionary Analysis. PLoS Comput. Biol. 10, 1–6.
 https://doi.org/10.1371/journal.pcbi.1003537.

604	Breathnach, R., Benoist, C., O'Hare, K., Gannon, F., Chambon, P., 1978. Ovalbumin						
605	gene: evidence for a leader sequence in mRNA and DNA sequences at the exon-						
606	intron boundaries. Proc. Natl. Acad. Sci. 75, 4853-4857.						
607	https://doi.org/10.1073/pnas.75.10.4853						
608	Breton, S., Milani, L., Ghiselli, F., Guerra, D., Stewart, D.T., Passamonti, M., 2014. A						
609	resourceful genome: updating the functional repertoire and evolutionary role of						
610	animal mitochondrial DNAs. Trends Genet. 30, 555–564.						
611	https://doi.org/10.1016/j.tig.2014.09.002.						
612	Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K.,						
613	Madden, T.L., 2009. BLAST+: Architecture and applications. BMC						
614	Bioinformatics 10, 1-9. https://doi.org/10.1186/1471-2105-10-421.						
615	Carver, T., Harris, S.R., Berriman, M., Parkhill, J., McQuillan, J.A., 2012. Artemis: an						
616	integrated platform for visualization and analysis of high-throughput sequence-						

- 617based experimental data.Bioinformatics28,464–469.618https://doi.org/10.1093/bioinformatics/btr703.
- 619 Chinese National Department of Fisheries, 2017. China fisheries statistic yearbook
 620 2016 (in Chinese). China Agriculture Press, Beijing.
- Chong, R.A., Mueller, R.L., 2013. Low metabolic rates in salamanders are correlated
 with weak selective constraints on mitochondrial genes. Evolution (N. Y). 67,
 894–899. https://doi.org/10.1111/j.1558-5646.2012.01830.x.
- Combosch, D.J., Collins, T.M., Glover, E.A., Graf, D.L., Harper, E.M., Healy, J.M.,
 Kawauchi, G.Y., Lemer, S., McIntyre, E., Strong, E.E., Taylor, J.D., Zardus, J.D.,
 Mikkelsen, P.M., Giribet, G., Bieler, R., 2017. A family-level Tree of Life for
 bivalves based on a Sanger-sequencing approach. Mol. Phylogenet. Evol. 107,
 191–208. https://doi.org/10.1016/j.ympev.2016.11.003.
- Combosch, D.J., Giribet, G., 2016. Clarifying phylogenetic relationships and the
 evolutionary history of the bivalve order Arcida (Mollusca: Bivalvia:
 Pteriomorphia). Mol. Phylogenet. Evol. 94, 298–312.
 https://doi.org/10.1016/j.ympev.2015.09.016.

- Cope, J.C.W., 1997. The early phylogeny of the class Bivalvia. Palaeontology 40, 713–
 746.
- FAO, 2018. The State of World Fisheries and Aquaculture. Food and AgricultureOrganization of the United Nations, Rome, Italy.
- Feng, Y., Li, Q., Kong, L., 2015. Molecular phylogeny of Arcoidea with emphasis on
 Arcidae species (Bivalvia: Pteriomorphia) along the coast of China: Challenges to
 current classification of arcoids. Mol. Phylogenet. Evol. 85, 189–196.
 https://doi.org/10.1016/j.ympev.2015.02.006.
- Feng, Y., Li, Q., Yu, H., Kong, L., 2017. Complete mitochondrial genome sequence of *Cucullaea labiata* (Arcoida: Cucullaeidae) and phylogenetic implications. Genes
 Genomics 39, 867–875. https://doi.org/10.1007/s13258-017-0548-1.
- Giribet, G., Wheeler, W., 2005. On bivalve phylogeny: a high-level analysis of the
 Bivalvia (Mollusca) based on combined morphology and DNA sequence data.
 Invertebr. Biol. 121, 271–324. https://doi.org/10.1111/j.17447410.2002.tb00132.x.
- Gonzalez, V.L., Andrade, S.C.S., Bieler, R., Collins, T.M., Dunn, C.W., Mikkelsen,
 P.M., Taylor, J.D., Giribet, G., 2015. A phylogenetic backbone for Bivalvia: an
 RNA-seq approach. Proc. R. Soc. B Biol. Sci. 282, 20142332.
 https://doi.org/10.1098/rspb.2014.2332.
- Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I.,
 Adiconis, X., Fan, L., Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E.,
 Hacohen, N., Gnirke, A., Rhind, N., Di Palma, F., Birren, B.W., Nusbaum, C.,
 Lindblad-Toh, K., Friedman, N., Regev, A., 2011. Full-length transcriptome
 assembly from RNA-Seq data without a reference genome. Nat. Biotechnol. 29,
 644–652. https://doi.org/10.1038/nbt.1883.
- Gurevich, A., Saveliev, V., Vyahhi, N., Tesler, G., 2013. QUAST: Quality assessment
 tool for genome assemblies. Bioinformatics. 29, 1072–1075.
 https://doi.org/10.1093/bioinformatics/btt086.

- 661 Hou, Y., Wu, B., Liu, Z.H., Yang, A.G., Ren, J.F., Zhou, L.Q., Dong, C.G., Tian, J.T.,
- 662 2016. Complete mitochondrial genome of Ark shell *Scapharca subcrenata*.
 663 Mitochondrial DNA 27, 939–940. https://doi.org/10.3109/19401736.2014.926495.

Huber, M., 2010. Compendium of bivalves. A full-color guide to 3,300 of the world

- 665 marine bivalves. A status on Bivalvia after 250 years of research, Hackenheim:666 ConchBooks.
- Huchon, D., Szitenberg, A., Shefer, S., Ilan, M., Feldstein, T., 2015. Mitochondrial
 group I and group II introns in the sponge orders Agelasida and Axinellida. BMC
 Evol. Biol. 15, 1–14. https://doi.org/10.1186/s12862-015-0556-1.
- Jaccard, A., 1869. Description géologique du Jura Vaudois et Neuchatelois. Matériaux
 pour la Cart. géologique la Suisse 6, 1–340.
- Katoh, K., Standley, D.M., 2013. MAFFT multiple sequence alignment software
 version 7: Improvements in performance and usability. Mol. Biol. Evol. 30, 772–
 780. https://doi.org/10.1093/molbev/mst010.
- Kohany, O., Gentles, A.J., Hankus, L., Jurka, J., 2006. Annotation, submission and
 screening of repetitive elements in Repbase: RepbaseSubmitter and Censor. BMC
 Bioinformatics 7, 1–7. https://doi.org/10.1186/1471-2105-7-474.
- Kück, P., Meusemann, K., 2010. FASconCAT: Convenient handling of data matrices.
 Mol. Phylogenet. Evol. 56, 1115–1118.
 https://doi.org/10.1016/j.ympev.2010.04.024.
- Lanfear, R., Calcott, B., Ho, S.Y.W., Guindon, S., 2012. PartitionFinder: Combined
 selection of partitioning schemes and substitution models for phylogenetic
 analyses. Mol. Biol. Evol. 29, 1695–1701. https://doi.org/10.1093/molbev/mss020.
- Langmead, B., Salzberg, S.L., 2012. Fast gapped-read alignment with Bowtie 2. Nat.
 Methods 9, 357–359. https://doi.org/10.1038/nmeth.1923.
- Lartillot, N., Philippe, H., 2004. A Bayesian mixture model for across-site
 heterogeneities in the amino-acid replacement process. Mol. Biol. Evol. 21, 1095–
 1109. https://doi.org/10.1093/molbev/msh112.

- Lartillot, N., Rodrigue, N., Stubbs, D., Richer, J., 2013. PhyloBayes MPI: Phylogenetic
 Reconstruction with Infinite Mixtures of Profiles in a Parallel Environment. Syst.
 Biol. 62, 611–615. https://doi.org/10.1093/sysbio/syt022.
- Lavrov, D. V., 2010. Rapid proliferation of repetitive palindromic elements in mtDNA
 of the endemic Baikalian sponge *Lubomirskia baicalensis*. Mol. Biol. Evol. 27,
 757–60. https://doi.org/10.1093/molbev/msp317.
- Lemer, S., González, V.L., Bieler, R., Giribet, G., 2016. Cementing mussels to oysters
 in the pteriomorphian tree: a phylogenomic approach. Proc. R. Soc. B Biol. Sci.
 283, 20160857. https://doi.org/10.1098/rspb.2016.0857.
- Li, G. yao, Zhang, L., Liu, J.Z., Chen, S.G., Xiao, T.W., Liu, G.Z., Wang, J.X., Wang,
 L.X., Hou, M., 2016. Marine drug Haishengsu increases chemosensitivity to
 conventional chemotherapy and improves quality of life in patients with acute
 leukemia. Biomed. Pharmacother. 81, 160–165.
 https://doi.org/10.1016/j.biopha.2016.04.005.
- Li, Y., Kocot, K.M., Schander, C., Santos, S.R., Thornhill, D.J., Halanych, K.M., 2015.
 Mitogenomics reveals phylogeny and repeated motifs in control regions of the
 deep-sea family Siboglinidae (Annelida). Mol. Phylogenet. Evol. 85, 221–229.
 https://doi.org/10.1016/j.ympev.2015.02.008.
- Liu, Y.G., Kurokawa, T., Sekino, M., Tanabe, T., Watanabe, K., 2013. Complete
 mitochondrial DNA sequence of the ark shell *Scapharca broughtonii*: An ultralarge metazoan mitochondrial genome. Comp. Biochem. Physiol. Part D
 Genomics Proteomics 8, 72–81. https://doi.org/10.1016/j.cbd.2012.12.003.
- Marko, P.B., 2002. Fossil calibration of molecular clocks and the divergence times of
 geminate species pairs separated by the Isthmus of Panama. Mol. Biol. Evol. 19,
 2005–2021. https://doi.org/10.1093/oxfordjournals.molbev.a004024.
- Matsumoto, M., 2003. Phylogenetic analysis of the subclass Pteriomorphia (Bivalvia)
 from mtDNA COI sequences. Mol. Phylogenet. Evol. 27, 429–440.
 https://doi.org/10.1016/S1055-7903(03)00013-7.

717	Mikkelsen, N.T.,	Kocot, F	К.М., Н	alanych,	K.M.,	2018.	Mitog	enomics	reveals
718	phylogenetic	relationsh	nips of	caudofov	veate	aplacopl	noran	molluscs.	Mol.
719	Phylogenet. E	vol. 127, 4	129–436.	https://do	oi.org/1	0.1016/j	.ympev	v.2018.04.	.031

- Miya, M., Kawaguchi, A., Nishida, M., 2001. Mitogenomic exploration of higher
 teleostean phylogenies: A case study for moderate-scale evolutionary genomics
 with 38 newly determined complete mitochondrial DNA sequences. Mol. Biol.
 Evol. 18, 1993–2009. https://doi.org/10.1093/oxfordjournals.molbev.a003741.
- Morton, B., Puljas, S., 2016. The ectopic compound ommatidium-like pallial eyes of
 three species of Mediterranean (Adriatic Sea) *Glycymeris* (Bivalvia: Arcoida).
 Decreasing visual acuity with increasing depth? Acta Zool. 97, 464–474.
 https://doi.org/10.1111/azo.12140.
- Nakamura, Y., 2005. Suspension feeding of the ark shell *Scapharca subcrenata* as a
 function of environmental and biological variables. Fish. Sci. 71, 875–883.
 https://doi.org/10.1111/j.1444-2906.2005.01040.x.
- Newell, N.D., 1969. Classification of Bivalvia. In: Moore, R. (Ed.), Treatise on
 Invertebrate Paleontology, Part N, Mollusca 6. vol. 1. Bivalvia. Geological Society
 of America & University of Kansas, Boulder-Lawrence. pp. N205–N244.
- Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., Minh, B.Q., 2015. IQ-TREE: A Fast
 and Effective Stochastic Algorithm for Estimating Maximum-Likelihood
 Phylogenies. Mol. Biol. Evol. 32, 268–274.
 https://doi.org/10.1093/molbev/msu300.
- Ojala, D., Montoya, J., Attardi, G., 1981. TRNA punctuation model of RNA processing
 in human mitochondria. Nature 290, 470–474. https://doi.org/10.1038/290470a0.
- Oliver, P.G., Holmes, A.M., 2006. The Arcoidea (Mollusca: Bivalvia): A review of the
 current phenetic-based systematics. Zool. J. Linn. Soc. 148, 237–251.
 https://doi.org/10.1111/j.1096-3642.2006.00256.x.
- Osigus, H.J., Eitel, M., Bernt, M., Donath, A., Schierwater, B., 2013. Mitogenomics at
 the base of Metazoa. Mol. Phylogenet. Evol. 69, 339–351.
 https://doi.org/10.1016/j.ympev.2013.07.016.

- Paradis, E., Claude, J., Strimmer, K., 2004. APE: Analyses of phylogenetics and
 evolution in R language. Bioinformatics 20, 289–290.
 https://doi.org/10.1093/bioinformatics/btg412
- R Development Core Team, 2018. R: a language and environment for statistical
 computing. Vienna (Austria): R Foundation for Statistical Computing.
- Rand, D.M., 1993. Endotherms, ectotherms, and mitochondrial genome-size variation.
 J. Mol. Evol. 37, 281–295. https://doi.org/10.1007/BF00175505.
- Ren, J., Liu, X., Jiang, F., Guo, X., Liu, B., 2010. Unusual conservation of
 mitochondrial gene order in *Crassostrea* oysters: evidence for recent speciation in
 Asia. BMC Evol. Biol. 10, 394. https://doi.org/10.1186/1471-2148-10-394.
- Revell, L.J., 2012. phytools: an R package for phylogenetic comparative biology (and other things). Methods Ecol. Evol. 3, 217–223. https://doi.org/10.1111/j.2041-210X.2011.00169.x.
- Rice, P., Longden, I., Bleasby, A., 2000. EMBOSS: the European Molecular Biology
 Open Software Suite. Trends Genet. 16, 276–277. https://doi.org/10.1016/S01689525(00)02024-2.
- Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz, G.,
 Mesirov, J.P., 2011. Integrative genomics viewer. Nat. Biotechnol. 29, 24–26.
 https://doi.org/10.1038/nbt.1754.
- Ruby, J.G., Bellare, P., Derisi, J.L., 2013. PRICE: software for the targeted assembly
 of components of (Meta) genomic sequence data. G3 (Bethesda). 3, 865–880.
 https://doi.org/10.1534/g3.113.005967.
- Rutherford, K., Parkhill, J., Crook, J., Horsnell, T., Rice, P., Rajandream, M.A., Barrell,
 B., 2000. Artemis: sequence visualization and annotation. Bioinformatics 16, 944–
 945. https://doi.org/10.1093/bioinformatics/16.10.944.
- Smith, D.R., Snyder, M., 2007. Complete mitochondrial DNA sequence of the scallop *Placopecten magellanicus*: Evidence of transposition leading to an
 uncharacteristically large mitochondrial genome. J. Mol. Evol. 65, 380–391.
 https://doi.org/10.1007/s00239-007-9016-x.

- Steiner, G., Hammer, S., 2000. Molecular phylogeny of the Bivalvia inferred from 18S
 rDNA sequences with particular reference to the Pteriomorphia. Geol. Soc.
 London, Spec. Publ. 177, 11–29. https://doi.org/10.1144/GSL.SP.2000.177.01.02.
- Stöger, I., Schrödl, M., 2013. Mitogenomics does not resolve deep molluscan
 relationships (yet?). Mol. Phylogenet. Evol. 69, 376–392.
 https://doi.org/10.1016/j.ympev.2012.11.017.
- Strotz, L.C., Saupe, E.E., Kimmig, J., Lieberman, B.S., 2018. Metabolic rates, climate
 and macroevolution: a case study using Neogene molluscs. Proc. R. Soc. B Biol.
 Sci. 285, 20181292. https://doi.org/10.1098/rspb.2018.1292.
- Sturmer, L., Nunez, J., Creswell, R., Baker, S., 2009. The potential of blood ark and
 ponderous ark aquaculture in Flordia: results of spawning, larval rearing, nursery
 and growout trials. SeaGrant, Florida TP-169. Shellfish Aquaculture Extension
 Program, University of Florida.
- Sun, S, Kong, L., Yu, H., Li, Q., 2015a. The complete mitochondrial genome of *Scapharca kagoshimensis* (Bivalvia: Arcidae). Mitochondrial DNA 26, 957–958.
 https://doi.org/10.3109/19401736.2013.865174.
- Sun, S, Kong, L., Yu, H., Li, Q., 2015b. The complete mitochondrial DNA of
 Tegillarca granosa and comparative mitogenomic analyses of three Arcidae
 species. Gene 557, 61–70. https://doi.org/10.1016/j.gene.2014.12.011.
- Sun, S, Kong, L., Yu, H., Li, Q., 2015c. Complete mitochondrial genome of *Anadara vellicata* (Bivalvia: Arcidae): A unique gene order and large atypical non-coding
 region. Comp. Biochem. Physiol. Part D Genomics Proteomics 16, 73–82.
 https://doi.org/10.1016/j.cbd.2015.08.001.
- 798 Sun, S., Li, Q., Kong, L., Yu, H., 2016. Complete mitochondrial genomes of Trisidos 799 kiyoni and Potiarca pilula: Varied mitochondrial genome size and highly 800 Arcidae. 33794. rearranged gene order in Sci. Rep. 6, 801 https://doi.org/10.1038/srep33794.

- Sun, S., Li, Q., Kong, L., Yu, H., 2017. Limited locomotive ability relaxed selective
 constraints on molluscs mitochondrial genomes. Sci. Rep. 7, 10628.
 https://doi.org/10.1038/s41598-017-11117-z.
- Talavera, G., Castresana, J., 2007. Improvement of phylogenies after removing
 divergent and ambiguously aligned blocks from protein sequence alignments. Syst.
 Biol. 56, 564–577. https://doi.org/10.1080/10635150701472164.
- Vallès, Y., Halanych, K.M., Boore, J.L., 2008. Group II introns break new boundaries:
 Presence in a bilaterian's genome. PLoS One 3, e1488.
 https://doi.org/10.1371/journal.pone.0001488.
- 811 Williams, S.T., Foster, P.G., Hughes, C., Harper, E.M., Taylor, J.D., Littlewood, D.T.J.,
- B12 Dyal, P., Hopkins, K.P., Briscoe, A.G., 2017. Curious bivalves: Systematic utility
 and unusual properties of anomalodesmatan mitochondrial genomes. Mol.
 Phylogenet. Evol. 110, 60–72. https://doi.org/10.1016/j.ympev.2017.03.004.
- Wu, X., Xu, X., Yu, Z., Kong, X., 2009. Comparative mitogenomic analyses of three
 scallops (Bivalvia: Pectinidae) reveal high level variation of genomic organization
 and a diversity of transfer RNA gene sets. BMC Res. Notes 2, 1–7.
 https://doi.org/10.1186/1756-0500-2-69.
- Wu, X., Xu, X., Yu, Z., Wei, Z., Xia, J., 2010. Comparison of seven *Crassostrea*mitogenomes and phylogenetic analyses. Mol. Phylogenet. Evol. 57, 448–454.
 https://doi.org/10.1016/j.ympev.2010.05.029.
- Zhang, B., Zhang, Y.H., Wang, X., Zhang, H.X., Lin, Q., 2017. The mitochondrial
 genome of a sea anemone *Bolocera* sp. exhibits novel genetic structures
 potentially involved in adaptation to the deep-sea environment. Ecol. Evol. 7,
 4951–4962. https://doi.org/10.1002/ece3.3067.
- 826

827 Highlights:

• 17 newly sequenced mtDNA genomes of Arcoida were used to explore phylogeny

829 and mitochondrial genome size evolution.

- Noetiidae, Cucullaeidae and Glycymerididae are nested within a polyphyletic
 Arcidae.
- Multiple independent expansions and contractions of mitochondrial genome size
 were found in Arcoida.
- Mitochondrial genome size variation appears correlated to tandem repeats and
- 835 inverted repeats.

836

	Journal Pre-proofs
837	CRediT author statement
838	
839	L.K., Y.L. Q.L. and K.M.H. conceived and designed this study. L.K.,
840	Y.Y., L.Q., Y.L., K.M.K. and K.M.H. collected specimens. L.K., Y.L.
841	and K.M.H. conducted the bioinformatic and phylogenetic analyses. L.K.
842	and Y.L. prepared the figures and submitted sequences to GenBank. All
843	authors contributed in the writing of the manuscript.
844	



Anadara crebricost	broughtonii (1) (GB) 48.2 k a broughtonii (2) (GB) 47.0 k capharca kagoshimensis (GB) a kagoshimensis 56.2 kb inaequivalvis 45.9 kb Anadara trapezia ata 36.7 kb 33.4 kb	(b) 46.7 kb		
Anadara Vellicata (GB) Potiarca pilula (GB) Tegillarca sp. 1.00/100 Tegillarca granosa (GB) 1.00/100 Tegillarca nodifera Vara transversa 18.8 kb Anadara Vellicata (GB) 17.5 kb Yara transversa	GB) 28.4 kb 50.4 kb 31.6 kb 38.7 kb		"Arcidae"	
semitorta (GB) 19.6 kb mitorta 19.5 kb rirescens 24.5 kb	"Arcinae"		"Arcidae"	
18.7 kb			Noetiidae	
17.9 kb Ana () 19.0 kb		Glyc	ymerididae	
20.5 kb 25.9 kb		C	ucullaeidae	
44.7 kb "Arcinae" 18.0 kb			"Arcidae"	
		Pł	nilobryidae -	—Li
1.00/100	- Mizuhopecten yessoensis - Mizuhopecten yessoensis - Pinc Ostrea dense	 16.2 kb 21.0 kb tada maxima 17. tamellosa 16. 	.0 kb .3 kb	
	1.00/100 C	rassostrea hongkongen ⁰⁰ – Crassostrea gigas	sis () 16.5 kb () 18.2 kb	

















Figure



Figure



CRediT author statement

L.K., Y.L. Q.L. and K.M.H. conceived and designed this study. L.K., Y.Y., L.Q., Y.L., K.M.K. and K.M.H. collected specimens. L.K., Y.L. and K.M.H. conducted the bioinformatic and phylogenetic analyses. L.K. and Y.L. prepared the figures and submitted sequences to GenBank. All authors contributed in the writing of the manuscript.