1 # A comprehensive non-redundant reference transcriptome for

2 # the Atlantic silverside *Menidia menidia*

3

4

5 Nina Overgaard Therkildsen[a,1,*] and Hannes Baumann[b]

6

7

8 [a]Department of Biology, Hopkins Marine Station, Stanford University, 120 Ocean View

9 Blvd, CA-93950 Pacific Grove, USA

10 [b]Department of Marine Sciences, University of Connecticut, 1080 Shennecossett Road,

11 CT-06340 Groton, USA

12

13 [1]Current address: Department of Natural Resources, Cornell University, 208 Fernow Hall,

14 NY-14853 Ithaca, USA

15

16 *Correspondence to: Nina Overgaard Therkildsen (nt246@cornell.edu)

17 Department of Natural Resources, Cornell University, 208 Fernow Hall, NY-14853 Ithaca,

18 USA

19

20

21

22 **Declarations of interest**: none

23

**ABSTRACT**

The Atlantic silverside (*Menidia menidia*) has been the focus of extensive research efforts in ecology, evolutionary biology, and physiology over the past three decades, but lack of genomic resources has so far hindered examination of the molecular basis underlying the remarkable patterns of phenotypic variation described in this species. We here present the first reference transcriptome for *M. menidia*. We sought to capture a single representative sequence from as many genes as possible by first using a combination of Trinity and the CLC Genomics Workbench to *de novo* assemble contigs based on RNA-seq data from multiple individuals, tissue types, and life stages. To reduce redundancy, we passed the combined raw assemblies through a stringent filtering pipeline based both on sequence similarity to related species and computational predictions of transcript quality, condensing an initial set of >480,000 contigs to a final set of 20,998 representative contigs, amounting to a total length of 53.3 Mb. In this final assembly, 91% of the contigs were functionally annotated with putative gene function and gene ontology (GO) terms and/or InterProScan identifiers. The assembly contains complete or nearly complete copies of >95% of 248 highly conserved core genes present in low copy number across higher eukaryotes, and partial copies of another 3.8%, suggesting that our assembly provides relatively comprehensive coverage of the *M. menidia* transcriptome. The assembly provided here will be an important resource for future research.

**INTRODUCTION**

The Atlantic silverside *Menidia menidia* (Atherinidae) is an abundant forage fish that inhabits nearshore environments along the east coast of North America, from northern Florida, USA to the Gulf of St. Lawrence, Canada (Hice *et al.* 2012). Its broad distribution along one of the steepest latitudinal temperature gradients in the world, combined with its ecological importance, its semelparous annual life cycle, and the relative ease with which it can be reared in the laboratory has made the Atlantic silverside a valuable model species for ecological and evolutionary research over the past three decades. Extensive laboratory and field studies have, for example, shown that the Atlantic silverside shows a remarkable degree of either co-gradient or counter-gradient variation in a suite of traits across latitudes, including growth rates, fecundity, metabolic rates, vertebral counts, swimming performance, and predator avoidance (reviewed in Conover *et al.* 2005). Common garden experiments have established that these trait differences have a clear genetic basis and often vary between locations less than 100 km apart—spatial scales across which silverside populations mix extensively (Hice *et al.* 2012). This detailed demonstration of pronounced genetic trait differences maintained despite strong gene flow has played an important role in shifting earlier perceptions about local adaptation being rare or absent in highly connected marine environments (Conover *et al.* 2005).

The Atlantic silverside has also been important for studies of rapid adaptation, providing crucial experimental evidence for fishing causing rapid evolution in the exploited populations (Conover & Munch 2002). It also provided the first discovery of temperature-dependent sex determination in fishes (reviewed in Conover *et al.* 2005). More recently, it has been an important model for quantifying novel effects of climate stressors such as

72    ocean warming, acidification and reduced oxygen levels (e.g. Murray *et al.* 2017;

73    Baumann *et al.* 2018) and for examining the geographic distribution of environmental

74    contaminants (Baumann *et al.* 2016). Its close relative, the inland silverside (*M. beryllina*)

75    is also frequently used in ecotoxicology studies (e.g. Jeffries et al. 2015). The Atlantic

76    silverside is therefore a central species for diverse research programs, yet genomic

77    resources have not been available for exploring the molecular basis underlying the many

78    fascinating evolutionary, ecological, and physiological patterns it exhibits.

79

80    This paper outlines how we generated and annotated the first comprehensive, non-

81    redundant *de novo* reference transcriptome for the Atlantic silverside. We needed this

82    reference for 'in silico exome capture' (Therkildsen & Palumbi 2017, Therkildsen et al.

83    2019) that would let us survey genome-wide patterns of variation underlying local

84    adaptation and rapid fisheries-induced evolution in this species in a cost-effective way. *De*

85    *novo* transcriptome assemblies often contain considerable redundancy with different allelic

86    variants, transcript splice variants, or overlapping fragments of the same transcript being

87    represented by separate contigs. This redundancy provides important information for

88    some types of analysis, but because we wanted to use the transcriptome as a reference

89    for population genomic analysis, our goal was to identify just a single representative

90    complete transcript for each gene, so that genomic sequencing reads could map to unique

91    positions. By pooling RNA samples from multiple individuals, tissue types, and life stages,

92    assembling contigs with two different *de novo* assembly algorithms, and passing the

93    resulting raw assemblies through a stringent filtering pipeline based both on sequence

94    similarity to related species and computational predictions of transcript quality, we

95    successfully minimized redundancy while capturing and retaining maximal diversity of

96 transcripts in our final assembly. The *de novo* Atlantic silverside transcriptome presented

97 here with associated functional annotation will be an important resource for future

98 research.

99

100 **DATA DESCRIPTION**

101 **SAMPLES, LIBRARY PREPARATION, AND SEQUENCING**

102 To capture a broad diversity of transcripts expressed at different life stages and in different

103 tissue types, we based our RNA sequencing on five larval and three adult *M. menidia*

104 (Table 1). The adults were collected directly from the wild at Poquot Beach (NY, N

105 40.9475, W 73.1025) in June 2013. At the same time, we also collected three pairs of

106 parent fish that were strip-spawned to produce three groups of full-sib larvae for rearing in

107 the laboratory following the procedure described by Murray et al. (2014). Twelve days post

108 hatching, the larvae were sacrificed and all samples were stored in RNAlater. All animal

109 handling was in accordance with NIH guidelines and approved under Institutional Animal

110 Care and Use Committee (IACUC) protocol 2010-1842-F at Stony Brook University.

111

112 We extracted total RNA from all samples with the Qiagen RNeasy Plus Universal Tissue

113 Mini Kit (Qiagen GmbH, Hilden, Germany). For each of five larvae, we used the entire

114 animal in a single extraction. For each adult, we did separate extractions for different

115 tissue types (including brain, heart, liver, gonad, muscle, gill, skin, spinal cord, fin, eye)

116 and pooled even quantities of these extracts for each individual. We then prepared a

117 single individually indexed cDNA library for each fish (pooled extracts from all tissue types

118 in a single library) with Illumina's TruSeq RNA sample prep kit v2 (Illumina Inc., San

119 Diego, CA, USA). All eight libraries were sequenced in a single lane of 100 bp paired-end

reads on an Illumina HiSeq 2000 at the University of Utah's Bioinformatics Core Facility,

yielding a total of 170 million raw sequence read pairs (between 16 and 26 million read

pairs per individual, amounting to a total of 34 Gb sequence). Our workflow for processing

the raw reads are shown in Fig. 1

**DATA QUALITY FILTERING**

After removing exact duplicate read pairs (~13% of the total) with the program Fastuniq

v1.1 (Xu *et al.* 2012), we used Trimmomatic v0.32 (Bolger *et al.* 2014) to trim off adapter

sequence and the first base of each read (because of a highly inflated C-content at this

position). We also used Trimmomatic's sliding window approach to trim off the rest of the

read if the average sequence quality over any four bases fell below 20, and discarded

reads shorter than 50 bp after this filtering (~7% of reads). We conservatively discarded a

further 2.6% of reads because they mapped to potential contamination databases (human,

bacterial, viral, rRNA, and Artemia (feed for the larvae)) with bowtie2 v2.2.3 (Langmead &

Salzberg 2012) in 'sensitive' preset mode. Finally, we used the program FLASH v. 1.2.9

(Magoč & Salzberg 2011) with default settings to merge overlapping read ends into single

consensus sequences (merging 68% of the remaining pairs), resulting in a filtered data set

of 92 million merged reads (length 50-240 bp) and 43 million read pairs amounting in total

to 21.1 Gb sequence.

***DE NOVO* ASSEMBLY AND REDUNDANCY REDUCTION**

Because different assembly algorithms and parameter settings may recover different

transcripts, we used two different programs to *de novo* assemble the pooled set of filtered

RNA-seq reads from all libraries. First, we generated two assemblies with the CLC

144     Genomic Workbench v6.0.2 (CLC Bio), one using the automatically optimized parameters

145     (word size 25, bubble size 50) and one using a larger word size (k-mer) of 40 to facilitate

146     more contiguous and accurate assembly of highly expressed transcripts. For both

147     assemblies, we mapped reads back to the initial contigs to update the consensus

148     sequence, and we broke up scaffolded sequence with no read support, only maintaining

149     contigs >200 bp. In parallel, we assembled the reads with Trinity v. r20131110 (Grabherr

150     *et al.* 2011) using the default settings (including a fixed k-mer size of 25). The Trinity

151     output explicitly clusters related 'isoforms', and since we were only interested in retaining a

152     single representative transcript for each gene, we mapped all reads back to the assembly

153     and extracted only the isoform with the highest mapped read depth within each

154     subcomponent following the procedure by Yang and Smith (2013). The three *de novo*

155     assemblies contained between 135,931 and 193,079 contigs each, for a total of 483,424

156     contigs in the combined set (Table 2).

157

158     A comparison of the three assemblies (CLC (k-mer 25), CLC (k-mer 40), and Trinity

159     (single isoform per cluster)) with blastn v2.2.29+ (Camacho *et al.* 2009) revealed that only

160     78-87% of the contigs in each assembly had significant hits (e-value $<10^{-3}$) to the other

161     assemblies, indicating that each assembly contained a set of unique transcripts. To

162     maintain maximal transcript diversity, we therefore proceeded with a merged set of all

163     three assemblies. The merged assemblies contain substantial redundancy, so to collapse

164     the contig set into the longest representative for each unique sequence we used cd-hit-est

165     v4.5.4 (Li & Godzik 2006) to remove shorter contigs that showed >95% sequence

166     similarity to other contigs. Due to assembly challenges, some genes could also be

167     presented by multiple different fragments rather than a transcript of complete length, so to

168  join partial assemblies (fragments) of the same transcript, we used CAP3 v12/21/07

169  (Huang 1999) to meta-assemble contigs with >95% similarity over at least 100 bp (an

170  approach shown to improve the quality of transcriptome assemblies e.g. by Melchior et al.

171  (2014)). Since both the *de novo* assembly processes and the meta-assembly may

172  introduce chimeric contigs, we used the method by Yang and Smith (2013) to break up

173  likely chimeras (observed in 0.8% of transcripts) based on separate blastx comparisons to

174  the peptide sets for three reference fish species (see below). The resulting redundancy-

175  reduced contig set contained 177,877 contigs (Table 2).

176

177  **SELECTING PUTATIVE GENE ORTHOLOGS**

178  Because we wanted to reduce our contig set to only include a single representative

179  transcript for each gene, we used a reciprocal best hit blast approach to extract non-

180  redundant putative orthologs to the gene sets in the three most closely related species for

181  which annotated genome assemblies were available at the time: platyfish (*Xiphophorus*

182  *maculatus*), medaka (*Oryzias latipes*), and Nile tilapia (*Oreochromis niloticus*). We

183  compared our contig set against the full peptide set for each reference species

184  (downloaded from Ensembl release 75 (Zerbino *et al.* 2018)) with blastx, and then

185  compared the peptide sequences for each species to our contig set with tblastn. For each

186  reference species, we recorded reciprocal best hits (RBHs) when a contig and a protein

187  had a best match to each other (e-value<$10^{-4}$). We then used a sequential approach to

188  select a combined set of 19,349 contigs in our *Menidia* assembly that were RBHs (and

189  therefore putative orthologs) to a unique peptide sequence in at least one of the reference

190  species (Supplementary Note).

191

192    Because all our reference species diverged from the silverside >75 million years ago

193    (Near *et al.* 2012; Campanella *et al.* 2015), the RBH contig selection procedure will fail to

194    identify recently diverged genes. To recover additional high quality non-redundant

195    transcripts, we used TransDecoder v. r20131110 (Haas *et al.* 2013) to predict coding

196    regions in our redundancy-reduced contig set on the basis of nucleotide composition, open

197    reading frame (ORF) length and Pfam domain content. Transdecoder predicted candidate

198    coding sequence of at least 100 amino acids in 39,604 contigs and of the 15,222 that

199    contained complete length ORFs, we retained 1,961 which did not have a significant (e-

200    value<$10^{-2}$) blastn hit to the RBH contig set (and therefore are non-redundant). To

201    minimize potential contamination in our final assembly, we compared the joined contig set

202    (RBH contigs and non-redundant contigs with complete ORFs) to the NCBI non-redundant

203    protein database (NR) (downloaded on July 14 2014) with blastx and used the program

204    MEGAN v. 5.7.1 (Huson *et al.* 2011) to identify and remove 311 contigs with best hits to

205    non-chordate taxa or to human sequence, ending up with a final reference transcriptome

206    contig set of 20,998 contigs.

207

208    **FUNCTIONAL ANNOTATION**

209    The final contig set was functionally annotated with the Blast2GO v3.1.2 suite (Conesa *et*

210    *al.* 2005). For each sequence, we imported significant hits (e-value < $10^{-6}$) from blastx

211    searches against the UniProt Swiss-Prot and the NCBI non-redundant (NR) protein

212    databases and used Blast2GO's Blast Description Annotator tool to select the most

213    informative and relevant descriptor before mapping GO (Gene Ontology) terms to the

214    matches and applying the built-in annotation rule with the default parameters and evidence

215    code weights. We also imported GO-terms associated with the reciprocal-best-hit genes in

216 the reference fish species, and merged the combined sets of assigned annotations with

217 GO-terms inferred from InterProScan analysis of each sequence. As a final step, we used

218 the Blast2GO Validate Annotations tool to ensure that no parent-child redundancy was

219 present in the assigned GO-terms, and we applied the Annex tool to augment the

220 annotation based on inference of biological processes from commonly associated

221 molecular functions and cellular components. This way we obtained a total of 490,807 GO-

222 terms annotated to 19,117 of the contigs (91% of all contigs; the median number of GO-

223 terms per contig was 19, Table S1).

224

225 **EVALUATION OF COMPLETENESS AND UNIQUENESS**

226 The final non-redundant transcriptome assembly had significant blastx hits to 84% of gene

227 models in the platyfish genome (81% of these were reciprocal best hits), and for 74% of

228 these genes the top high-scoring segment pair covered >90% of the total length of the

229 reference peptide sequence, indicating complete or nearly complete transcripts. CEGMA

230 v2.5 (Parra *et al.* 2007) also detected complete or nearly complete copies of >95% of 248

231 highly conserved core genes present in low copy number across higher eukaryotes (and

232 partial copies of another 3.8%). Similarly, BUSCO (Simão et al. 2015) analysis flagged

233 only 3.1% of 4,584 highly conserved genes in Actinopterygii (ray-finned fish) species as

234 missing from the assembly (90.5% of these reference genes were detected as complete

235 copies, 6.4% as fragmented), further suggesting that the assembly provides a relatively

236 comprehensive coverage of the *M. menidia* transcriptome. The extensive transcript

237 diversity is likely caused by the inclusion of many different tissue types across two life

238 stages in our RNA-seq libraries.

239    Exposure to a variety of stressors prior to RNA harvesting may have increased transcript

240    diversity further and could be pursued in future work targeting specific response pathways.

241

242    The number of contigs in our final assembly is much closer to the number of coding genes

243    found in related species (21,437 - 23,774 for medaka, platyfish, and tilapia, (Zerbino *et al.*

244    2018)) than any of the larger assemblies. Yet, the strict redundancy reduction did result in

245    a small loss of transcripts diversity as the complete Trinity assembly and the merged set of

246    all raw assemblies actually included up to 100% of the CEGMA genes. However, this loss

247    of diversity was compensated for by much better mapping specificity. For the full Trinity

248    assembly and the total merged assembly, >93% of the cleaned RNA-seq reads mapped

249    back to the *de novo* reference with bowtie2 v2.2.3 (Langmead & Salzberg 2012) in the

250    'sensitive' preset mode, but only 61% (Trinity) or 9% (total merged) of these mapped to a

251    unique position (the remaining reads mapping to multiple contigs, Table 2). In contrast,

252    almost all (98%) of the 74% of RNA-seq reads that mapped to the final assembly mapped

253    only to a single position, suggesting that most genes are only represented by a single

254    contig and that this assembly therefore will be useful reference for mapping genomic

255    reads, as further demonstrated in Therkildsen and Palumbi (2017). In addition to the highly

256    non-redundant assembly that will be useful for population genomics and many other

257    purposes, we are also making each of our intermediate larger assemblies (see Fig. 1 and

258    Table 2) available as supplementary data files (File S2 and S3) for other types of analysis

259    that specifically targets redundancy among similar transcripts, e.g. analysis of splice

260    variation or variation within closely related gene families. With recent technological

261    advances, *de novo* assembly of the entire genome is an increasingly attainable goal for

262    many non-model organisms. Yet, the cost and effort involved in assembling only the

263 transcriptome generally is still much lower, so an important role remains for reference

264 transcriptomes - especially for studies focusing on functional genomic variation.

265

266

271

272

273 **References**

274 Baumann H, Cross EL, Murray CS (2018) Robust quantification of fish early life $CO_2$

275      sensitivities via serial experimentation. *Biology Letters*, **14**, 20180408.

276 Baumann Z, Mason RP, Conover DO, Balcom P, Chen CY, Buckman KL, Fisher NS,

277      Baumann H (2016) Mercury bioaccumulation increases with latitude in a coastal

278      marine fish (Atlantic silverside, *Menidia menidia*). *Canadian Journal of Fisheries and*

279      *Aquatic Sciences*, **74**, 1009–1015.

280 Bolger AM, Lohse M, Usadel B (2014) Trimmomatic: a flexible trimmer for Illumina

281      sequence data. *Bioinformatics*, **30**, 2114–2120.

282 Camacho C, Coulouris G, Avagyan V, Ma N (2009) BLAST+: architecture and

283      applications. *BMC Bioinformatics*, **10**, 421.

284 Campanella D, Hughes LC, Unmack PJ, Bloom DD, Piller KR, Orti G (2015) Multi-locus

285      fossil-calibrated phylogeny of Atheriniformes (Teleostei, Ovalentaria). *Molecular*

286      *Phylogenetics and Evolution*, **86**, 8–23.

287    Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M (2005) Blast2GO: a

288        universal tool for annotation, visualization and analysis in functional genomics

289        research. *Bioinformatics*, **21**, 3674–3676.

290    Conover DO, Munch SB (2002) Sustaining fisheries yields over evolutionary time scales.

291        *Science*, **297**, 94–96.

292    Conover DO, Arnott SA, Walsh MR, Munch SB (2005) Darwinian fishery science: lessons

293        from the Atlantic silverside (*Menidia menidia*). *Canadian Journal of Fisheries and*

294        *Aquatic Sciences*, **62**, 730–737.

295    Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,

296        Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, Di

297        Palma F, Birren BW, Nusbaum C, Lindblad-Toh K, Friedman N *et al.* (2011) Full-length

298        transcriptome assembly from RNA-Seq data without a reference genome. *Nature*

299        *Biotechnology*, **29**, 644–652.

300    Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB,

301        Eccles D, Li B, Lieber M, Macmanes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks

302        N, Westerman R, William T, Dewey CN, Henschel R *et al.* (2013) De novo transcript

303        sequence reconstruction from RNA-seq using the Trinity platform for reference

304        generation and analysis. *Nature Protocols*, **8**, 1494–1512.

305    Hice LA, Duffy TA, Munch SB, Conover DO (2012) Spatial scale and divergent patterns of

306        variation in adapted traits in the ocean. *Ecology Letters*, **15**, 568–575.

307    Huang X (1999) CAP3: A DNA Sequence Assembly Program. *Genome Research*, **9**, 868–

308        877.

309    Huson DH, Mitra S, Ruscheweyh H-J, Weber N, Schuster SC (2011) Integrative analysis

310        of environmental sequences using MEGAN4. *Genome Research*, **21**, 1552–1560.

311    Jeffries KM, Brander SM, Britton MT, Fangue NA, Connon RE (2015) Chronic exposures
312        to low and high concentrations of ibuprofen elicit different gene response patterns in a
313        euryhaline fish. *Environmental Science and Pollution Research*, **22**, 17397–17413.
314    Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nature*
315        *Methods*, **9**, 357–359.
316    Li W, Godzik A (2006) Cd-hit: a fast program for clustering and comparing large sets of
317        protein or nucleotide sequences. *Bioinformatics*.
318    Magoč T, Salzberg SL (2011) FLASH: fast length adjustment of short reads to improve
319        genome assemblies. *Bioinformatics*, **27**, 2957–2963.
320    Melicher D, Torson AS, Dworkin I, Bowsher JH (2014) A pipeline for the de novo assembly
321        of the *Themira biloba* (Sepsidae: Diptera) transcriptome using a multiple k-mer length
322        approach. *BMC Genomics*, **15**, 1–13.
323    Murray CS, Fuiman LA, Baumann H (2017) Consequences of elevated $CO_2$ exposure
324        across multiple life stages in a coastal forage fish. *ICES Journal of Marine Science*, **74**,
325        1051–1061.
326    Murray CS, Malvezzi A, Gobler CJ, Baumann H (2014) Offspring sensitivity to ocean
327        acidification changes seasonally in a coastal marine fish. *Marine Ecology Progress*
328        *Series*, **504**, 1–11.
329    Near TJ, Eytan RI, Dornburg A, Kuhn KL, Moore JA, Davis MP, Wainwright, P.C.,
330        Friedman M, Smith WL (2012) Resolution of ray-finned fish phylogeny and timing of
331        diversification. *Proceedings of the National Academy of Sciences of the United States*
332        *of America*, **109**, 13698–13703.
333    Parra G, Bradnam K, Korf I (2007) CEGMA: a pipeline to accurately annotate core genes
334        in eukaryotic genomes. *Bioinformatics*, **23**, 1061–1067.

335 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM (2015) BUSCO:

336      assessing genome assembly and annotation completeness with single-copy orthologs,

337      *Bioinformatics*, **31**, 3210–3212.

338 Therkildsen NO, Palumbi SR (2017) Practical low-coverage genomewide sequencing of

339      hundreds of individually barcoded samples for population and evolutionary genomics

340      in nonmodel species. *Molecular Ecology Resources*, **17**, 194–208.

341 Therkildsen NO, Wilder AP, Conover DO, Munch SB, Baumann H, Palumbi SR (2019)

342      Contrasting genomic shifts underlie parallel phenotypic evolution in response to

343      fishing. *Science*, **365**, 487–490.

344 Xu H, Luo X, Qian J, Pang X, Song J, Qian G, Chen J, Chen S (2012) FastUniq: a fast de

345      novo duplicates removal tool for paired short reads. *PLoS ONE*, **7**, e52249.

346 Yang Y, Smith SA (2013) Optimizing de novo assembly of short-read RNA-seq data for

347      phylogenomics. *BMC Genomics*, **14**, 328.

348 Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, Billis K, Cummins C,

349      Gall A, Girón CG, Gil L, Gordon L, Haggerty L, Haskell E, Hourlier T, Izuogu OG,

350      Janacek SH, Juettemann T, To JK, Laird MR *et al.* (2018) Ensembl 2018. *Nucleic*

351      *Acids Research*, **46**, D754–D761.

352

353

354 **Data accessibility**

355 The raw sequence data are deposited in the NCBI Sequence Read Archive (SRA) with

356 accession numbers SRR3990241- SRR3990248 associated with BioProject

357 PRJNA330848. The final assembly of the Atlantic silverside transcriptome (20,998 contigs)

358 is deposited in the NCBI GenBank Transcriptome Shotgun Assembly Sequence Database

359 (TSA) under Accession no. GEVY00000000. The transcriptome annotation table is

360 provided as Supplementary Table S1, and the full merged and redundancy-reduced

361 assemblies are provided in fasta format as Supplementary Files S2 and S3.

362

363

364

365 **List of supplementary files:**

366

367 Table S1. Annotation table for the final transcriptome assembly:

368 MenidiaTranscriptome_AnnotationTable_GO_Interpro.csv

369

370 File S1: Supplementary Note

371

372 File S2: Fasta file with the total combined contig set (483,424 contigs):

373 MenidiaTranscriptome_Complete_Merged_Contig_Set.fa

374

375 File S3: Fasta file with the redundancy-reduced contig set (177,877 contigs):

376 MenidiaTranscriptome_RedundancyReduced_Merged_Contig_Set.fa

377

378

379

380 **Figure legends**

381 **Fig. 1.** Diagram showing the sequence of steps in our bioinformatic workflow for cleaning

382 the RNA-seq read data, *de novo* assembly, and redundancy reduction. Yellow boxes

383 represent RNA-seq read data, blue boxes represent data processing steps, and red boxes

384 represent transcriptome assemblies. Statistics such as the total number of contigs, the

385 total assembled length and the proportion of conserved core genes found in each of the

386 intermediate assemblies are provided in Table 2.