

1 ***Transcriptomic profiling and microsatellite identification in cobia***
2 ***(Rachycentron canadum), using high throughput RNA-sequencing***

3

4 ***David Aciole Barbosa*** ¹**#**(0000-0003-3875-2307), ***Bruno C. Araújo*** ²**#**(
5 ***0000-0002-6432-4405***), ***Giovana Souza Branco*** ³(***0000-0002-4481-2436***);
6 ***Alexandre S. Simeone***¹(***0000-0003-3725-0671***), ***Alexandre W. S. Hilsdorf***
7 ***1***(***0000-0001-9565-8072***), ***Daniela L. Jabes*** ¹(***0000-0001-7297-0784***), ***Luiz***
8 ***R. Nunes*** ⁴(***0000-0001-9619-269X***), ***Renata G. Moreira*** ³(***0000-0001-8139-***
9 ***1776***), ***Fabiano B. Menegidio*** ¹*****(***0000-0002-4705-8352***)

10

11 ¹*Center of Biotechnology, University of Mogi das Cruzes, Mogi das Cruzes, 08780-*
12 *911, SP, Brazil*

13 ²*Cawthron Institute, Nelson, 7010, New Zealand*

14 ³*Department of Physiology, Bioscience Institute, University of São Paulo, São Paulo,*
15 *05508-090, SP, Brazil*

16 ⁴*Center for Natural and Human Sciences, Federal University of ABC, Santo André,*
17 *09210-580, SP, Brazil*

18

19 ** To whom correspondence should be addressed: Fabiano B. Menegidio, PhD, Center of*
20 *Biotechnology, University de Mogi das Cruzes (UMC), Av. Dr. Cândido X. de Almeida e*
21 *Souza, 200 - Centro Cívico, Mogi das Cruzes - SP, 08780-911, Brazil;*
22 *E-mail: fabianomenegidio@umc.br, fabiano.menegidio@bioinformatica.com.br*

23

24 *# These authors contributed equally to this article*

25

26

27 **ABSTRACT:** *Cobia (Rachycentron canadum)* is a marine teleost species with great
28 productive potential worldwide. However, the genomic information currently available
29 for this species in public databases is limited. Such lack of information hinders gene
30 expression assessments that might bring forward novel insights into the physiology,
31 ecology, evolution, and genetics of this potential aquaculture species. In this study,
32 we report the first *de novo* transcriptome assembly of *R. canadum* liver, improving
33 the availability of novel gene sequences for this species. Illumina sequencing of liver
34 transcripts generated 1,761,965,794 raw reads, which were filtered into
35 1,652,319,304 high-quality reads. *De novo* assembly resulted in 101,789 unigenes
36 and 163,096 isoforms, with an average length of 950.61 and 1,617.34 nt,
37 respectively. Moreover, we found that 126,013 of these transcripts bear potentially
38 coding sequences, and 125,993 of these elements (77.3%) correspond to
39 functionally annotated genes found in six different databases. We also identified 701
40 putative ncRNA and 35,414 putative lncRNA. Interestingly, homologues for 410 of
41 these putative lncRNAs have already been observed in previous analyzes with *Danio*
42 *rerio*, *Lates calcarifer*, *Seriola lalandi dorsalis*, *Seriola dumerili* or *Echeneis*
43 *naucrates*. Finally, we identified 7,894 microsatellites related to cobia's putative
44 lncRNAs. Thus, the information derived from the transcriptome assembly described
45 herein will likely assist future nutrigenomics and breeding programs involving this
46 important fish farming species.

47

48 **Keywords:** *Rachycentron canadum*; lncRNA; Transcriptome; Cobia; Microsatellites;
49 Aquaculture

50

51 1. INTRODUCTION

52 Currently, aquaculture is the most prominent food production industry, with
53 significant growth worldwide. Global aquaculture production increased by 5.3%
54 percent per year between 2001-2018, with a historical record of 114.5 million tons of
55 farmed species, including almost 17.7 million tons due to finfish production (FAO
56 2020). Marine aquaculture plays an essential role in the effort of providing the
57 increasing world's demand for animal-based protein. *Cobia (Rachycentron*
58 *canadum)* is a carnivorous marine fish of worldwide distribution, and it is the sole
59 representative of the Rachycentridae family, among farmed fish species. Currently,

60 *R. canadum* is regarded as the most promising marine fish species in Brazil, mostly
61 due to its fast growth rate (reaching about 4 to 6 kg per year), excellent meat quality
62 (with regards to color, texture and flavor), and high market value (Arnold et al. 2002;
63 Benetti et al. 2008; Nunes 2014). However, cobia production is still hindered by the
64 lack of nutrition information, which constraints this species' productivity in industrial
65 aquaculture operations (Fraser and Davies 2009).

66 Next-generation sequencing (NGS) studies have become an essential molecular tool
67 in aquaculture, assisting in the production of several commercial fish species, such
68 as *Sparus aurata* (Calduch-Giner et al. 2013), *Dicentrarchus labrax* (Magnanou et al.
69 2014), and *Salmo salar* (Glencross et al. 2015; Andrew et al. 2021). *De novo*
70 transcriptome assembly can be used in several different contexts, like
71 genomics/gene expression analyses and may be applied in many key areas of study,
72 such as conservation genetics, selective breeding, reproductive biology, and nutrition
73 (Leaver et al. 2008; Calduch-Giner et al. 2013; Fox et al. 2014). For example,
74 identifying genes associated with proteins and lipid metabolism in the liver can assist
75 in the development of specific diets to improve the productive chain of commercial
76 aquaculture species, such as cobia. Unfortunately, genomic information in cobia is
77 still scarce, limiting the development of such studies. Thus, to help overcome such
78 limitations, the present manuscript describes an assembled/annotated reference
79 transcriptome of hepatic cells in cobia juveniles. The information contained in this
80 dataset contributes to improve our knowledge regarding the biological and
81 physiological aspects of this fish species and establishes a solid foundation for future
82 studies involving population genomics, breeding programs, and nutrigenomics
83 involving this important marine farming fish.

84

85 **2. MATERIALS AND METHODS**

86 **2.1. Sample collection and RNA preparation**

87 Ninety cobia juveniles (128.85 ± 18.43 g) were obtained from a commercial hatchery
88 (Redemar Alevinos, SP, Brazil) and randomly allocated in three 2,000 L tanks. The
89 animals were kept under a mean temperature of 23 ± 1.5 °C and a photoperiod of
90 12L:12D throughout the trial, at the Marine Biology Center of the University of São
91 Paulo (CEBIMar). Animals were equally hand-fed twice a day, until apparent satiety,
92 with a commercial marine fish diet (Guabipirá, Guabi Nutrição e Saúde Animal S.A.,
93 SP, Brazil). After six weeks (42 days), fish were anesthetized with benzocaine (0.4 g

94 * mL⁻¹) and then euthanized by spinal cord section. Hepatic tissue samples from all
95 experimental animals were collected, immediately frozen in liquid nitrogen, and
96 subsequently stored at -80 °C for further analyses. Total RNA from hepatic tissue
97 samples was extracted using Rneasy Lipid Tissue kit (Qiagen), following the
98 manufacturer's instructions. RNA samples had their concentration determined with a
99 NanodropTM Spectrophotometer (ND-1000). RNA integrity was assessed using a
100 2100 Bioanalyzer System (Agilent Technologies, USA). This study's experimental
101 procedures were conducted according to the guidelines and approval of the
102 Institutional Animal Care and Use Ethics Committee (#008/2017).

103

104 **2.2. Library construction and sequencing**

105 RNA extracted from liver samples from all animals (90 fish, 30 per tank), were
106 equally diluted to 1,000 ng/μl concentration and pooled for library construction, using
107 the TruSeq RNA Sample Preparation kit, according to the manufacturer's
108 specifications (Illumina Inc., USA). Library quality was validated using a Bioanalyzer
109 2100 (Agilent Technologies, USA), and only samples with an RNA Integrity Number
110 (RIN) equal or above 7.5 were used. Finally, paired-end sequencing (2×75 bp) of
111 these cDNA libraries was conducted in an Illumina Nextseq® platform, according to
112 the manufacturer's recommendations. Minimum information about any (x) sequence
113 (MIxS) data for this study is available in Supplementary Table ST1-S1.

114

115 **2.3. Bioinformatics Analysis of Raw Data**

116 FASTQ format raw sequencing data was processed in a Public Galaxy Server
117 available at <https://usegalaxy.eu>. Initially, the quality of raw sequences was
118 assessed using FastQC (Andrews 2010) and MultiQC (Ewels et al. 2016). Fastp
119 (Chen et al. 2018) was then used to remove low-quality reads (Q<30), adapters, and
120 other contaminant sequences. Trinity software (Haas et al. 2013) was then used for
121 *de novo* transcriptome assembly of filtered reads, and assembly metrics were
122 obtained using the TrinityStats script. Transcriptome completeness was finally
123 assessed using the Benchmarking Universal Single-Copy Orthologs (BUSCO)
124 software (Seppey et al. 2019) based on OrthologDB version 9 (Zdobnov et al. 2017).

125

126 **2.4. Transcriptome Annotation**

127 Elements of the assembled transcriptome were functionally annotated by similarity
128 searches (blastn, e-value $\leq e^{-5}$) performed against RefSeq RNA and “*Lates/Seriola*”
129 - a custom database containing RefSeq transcripts from the cobia-related species
130 *Lates calcarifer*, *Seriola dumerili*, and *Seriola lalandi dorsalis*, which are in the
131 Kegg’s Organisms Complete Genomes (see
132 https://www.genome.jp/kegg/catalog/org_list.html). Additional annotations were
133 obtained with the aid of the Eukaryotic Non-Model Transcriptome Annotation
134 Pipeline (EnTAP) (Hart et al. 2020). Contigs were queried (blastx; using e-value $\leq e^{-5}$
135 and $\geq 50\%$ coverage) for similarity against the National Center for Biotechnology
136 Information non-redundant protein database (NCBI nr), NCBI proteins reference
137 database (RefSeq), the curated Swiss-Prot database from UniProt Knowledgebase
138 (UniProtKB) (UniProt Consortium 2019), and the EggNOG proteins database
139 (Huerta-Cepas et al. 2016). The EggNOG hits also helped to assign biological
140 function to individual elements, identifying their respective Gene Ontology (GO) (The
141 Gene Ontology Consortium 2019) and KEGG (Kyoto Encyclopedia of Genes and
142 Genomes) (Kanehisa and Goto 2000; Kanehisa 2019; Kanehisa et al. 2021) terms.
143 The EnTAP functional annotation process was carried out using a Dugong container
144 environment (Menegidio et al. 2018). Transcripts not annotated by EnTAP were
145 evaluated using the cmscan program (default parameters) by Infernal (Nawrocki and
146 Eddy 2013), for classification in the different families of non-coding RNAs, defined in
147 the Rfam database (Kalvari et al. 2018).

148

149 **2.5. Coding Potential Calculator and lncRNA Discovery**

150 Transcripts not annotated by EnTAP and Infernal were evaluated for their respective
151 coding potential (CP) with the aid of three tools: Coding Potential Calculator (CPC2)
152 (Kang et al. 2017), Coding-Potential Assessment Tool (CPAT) (Wang et al. 2013)
153 and RNASamba (Camargo et al. 2020). The transcripts identified as having non-
154 coding potential by all of these tools were separated for functional annotation
155 analysis. In this work, we considered putative lncRNAs transcripts with ≥ 200 nt that
156 were identified as non-coding by all of CP tools and not annotated by EnTAP /
157 Infernal. To discover conserved interspecies lncRNAs, we aligned putative lncRNA
158 sequences against the Zebrafish lncRNA Database (ZFLNC; Hu et al. 2018) using
159 blastn (Boratyn et al. 2013) with a cut-off value $\leq e^{-5}$ and $\geq 50\%$ identity (Fan et al.
160 2018). Similar, blastn searches were employed against ncRNA sequences available

161 at Ensembl from *Danio rerio*, *Lates calcarifer*, *Echeneis naucrates*, *Seriola lalandi*

162 *dorsalis* and *S. dumerili*.

163

164

165 **2.6. Detection of SSRs in lncRNAs**

166 The MlcroSATellite (MISA) software (Beier et al. 2017) was used to identify
167 microsatellites in the putative lncRNAs sequences. The Simple Sequence Repeats
168 (SSR) loci detection was done by searching for two- to six-nucleotide motifs, with a
169 minimum of 1/10, 2/6, 3/5, 4/5, 5/4 and 6/4 (motifs/repeats), as suggested by Gui et
170 al., (2013).

171

172 **3. RESULTS AND DISCUSSION**

173 **3.1 Transcriptome assembly and completeness**

174 Sequencing of the cDNA libraries derived from *R. canadum* liver material resulted in
175 1,761,965,794 raw reads. After high-quality-read selection and trimming, we were
176 left with a total of 1,652,319,304 reads (93.77% of raw reads), which were used for
177 *de novo* transcriptome assembly, using Trinity software (Haas et al. 2013). General
178 features of the *R. canadum* liver transcriptome are summarized in Table 1,
179 consisting of 101,789 unigenes and 163,096 isoforms (likely derived from cryptic
180 transcription start sites, alternative splicing or differential polyadenylation events).
181 The median (N50)/average length of these elements was 7,843/1,617.34 nt for
182 unigenes and 2,312/950.61 nt for isoforms. A total of 95,075 transcripts (58.29%)
183 were ≥ 500 nt. Identification of 83.8% of the complete universal genes (3,839 out of
184 the total 4,584 genes from Actinopterygii odb9 lineage) supported the high quality
185 and completeness of this transcriptome assembly (Fig. 1a). Among the 3,839
186 conserved BUSCO genes, 38.1% were single copy, while 45.7% were duplicated
187 (Supplementary Table ST1-S2).

188

189 **3.2 Functional annotation**

190 Transcriptome annotation against a series of databases (NCBI nr, NCBI RefSeq
191 RNA, Swiss-Prot, GO, KEGG, *Lates/Seriola*) resulted in functional assignment for
192 125,993 transcripts (77.3%). Most sequence homologies were found against NCBI
193 RefSeq RNA (122,741 transcripts, or 97.4%), followed by *Lates* and *Seriola* species
194 databases (118,570 transcripts, or 94.1%), NCBI nr (29,155 transcripts, or 23.14%),
195 NCBI RefSeq (25,728, or 20.42%) and Swiss-Prot (15,507 transcripts, or 12.30%)
196 (Table1; Supplementary Table ST1-S3a-S5). Sequence homologies identified by
197 EnTAP were distributed through many bony fish species, of which *Seriola dumerili*

198 was the most frequent (nr = 39.83%, RefSeq = 40.59%), followed by *Seriola lalandi*
199 *dorsalis* (nr = 23.27%, RefSeq = 25.86%), *Echeneis naucrates* (nr 17.9%, RefSeq =
200 18.37%) and *Larimichthys crocea* (nr = 4.99%, RefSeq = 3.13%) (Fig. 1b). Most
201 transcripts (75,060) were functionally annotated with GO terms by eggNOG (Huerta-
202 Cepas et al. 2016), which assigned 48,550 transcripts (65%) to biological processes,
203 34,492 to cellular components (46%) and 47,766 to molecular functions (64%). The
204 ten most representative functional groups within each category are shown in Fig. 1c.
205 A total of 23,936 isoforms were annotated into at least one KEGG pathway term
206 (Table1; Supplementary Table ST1-S6).

207 The final functional annotation allowed us to identify several marker genes for future
208 nutrigenomics initiatives, including elements involved in the following GO Biological
209 Processes: (i) GO:0007586 - Digestion of Nutrients (ex.: Carboxypeptidases, Trypsin
210 and Trypsin-like homologues, Chymotrypsin-like Elastase Family members, etc.); (ii)
211 GO:0042445 - Proteins Involved with Hormone Function/Metabolism (ex. Calcitonin,
212 Estrogen Receptors, Hecpidins, Insulin/Glucagon homologues and receptors, etc.);
213 (iii) GO:0006629/GO:0005975 Lipid/Carbohydrate Metabolic Processes (ex.
214 Phospholipases A/B/C/D, Bile Salt-stimulated Lipases, Chitinases, Galactosidases,
215 Alpha/Beta Glucosidases, etc.); (iv) GO:0044765 - Nutrient Transport (ABC
216 Transporters, Amino Acid Permeases, Apolipoprotein/Hemoglobin homologues,
217 etc.); (v) GO:0048644 - Muscle Structure and Morphogenesis (ex.
218 Actin/Myosin/Formin/Growth Factor homologues, etc.) and (vi) GO:0048565 -
219 Digestive Tract Development (ex. Kruppel-like Factors, Digestive Organ Expansion
220 Factor homologues, Pancreatic and Duodenal Homeobox 1-containing proteins,
221 Hepatocyte Growth Factors, etc), among others (see Supplementary Table ST1-S3b,
222 for details).

223

224 **3.3 Known Non-Coding RNAs**

225 For the 37,103 transcripts not annotated in the previous steps, the Infernal tool suite
226 (Nawrocki and 2013) was used to filter the presence of known non-coding RNAs
227 available from the RFAM database (Kalvari et al. 2018). The cmscan script allowed
228 annotation of 33,677 such sequences. Among these, we obtained 699 significant hits
229 (based on an e-value threshold of 0.001), allowing their classification as putative
230 non-coding RNAs (Table1; Supplementary Table ST1-S7). Among the ncRNAs
231 identified in our *R. canadum* transcriptome are small nucleolar RNAs (snoRNAs),

232 such as SNORD5, SNORD21, SNORD27, SNORD31, SNORD36, SNORD48,
233 SNORD52, SNORD63, SNORD78, SNORD88 and SNORD103 - which have also
234 been described in *Danio rerio*, *Gasterosteus aculeatus*, *Tetraodon nigroviridis* and
235 *Oryzias latipes*, according to the snoOPY database (<http://snoopy.med.miyazaki-u.ac.jp>).
236 Cajal small body RNAs (scaRNAs), which are involved in the modification of
237 snoRNAs, were also identified in our data. Among the scaRNAs, we found a
238 homologue for SCARNA7, a C/D RNA box, which localizes to Cajal bodies in HeLa
239 cells and is conserved in all vertebrates (Marz et al. 2011) MicroRNA (miRNA),
240 accounted for 170 non-coding transcripts, distributed across 111 miRNA types.
241 MicroRNAs have been found to be related to several relevant biological aspects in
242 fish, including regulation of growth, muscle, vascular and cardiac tissues, among
243 many others (Rasal et al. 2016). Moreover, Herkenhoff et al. (2018) suggested that
244 miRNAs may also serve as biomarkers for selection of adaptive traits for
245 aquaculture.

246

247 **3.4 Novel Non-Coding RNAs and Long Non-Coding RNAs**

248 Next, transcripts not annotated by EnTAP and Infernal were subjected to a coding
249 potential (CP) analysis, using four different CP calculator tools: Coding Potential
250 Calculator (CPC2) (Kang et al. 2017), Coding-Potential Assessment Tool (CPAT)
251 (Wang et al. 2013) and RNASamba (Camargo et al. 2020), which subdivided these
252 novel elements into coding and non-coding transcripts (Supplementary Table ST1-
253 S8). Among the remaining 36,404 unannotated elements, 0.054% were identified as
254 protein coding (representing 20 putative new genes), 2.65% were of undetermined
255 nature (968 transcripts) and 97.29% were classified as non-protein coding elements
256 (35,416 transcripts). The distribution of transcripts classified as non-protein coding
257 elements (non-coding RNAs) can be observed in Fig. 1d (the 35,416 elements
258 specified at this point are represented by the intersection of all CP tools). Among
259 these, 35,414 were larger than 200 bp, and were finally classified as putative long-
260 noncoding RNAs (lncRNAs) (Supplementary Table ST1-S9). Most cobia lncRNAs
261 (26,255, 74.14%) are larger than 400 bp; 5,582 of them are \geq 600 bp (15.76%),
262 1,871 (5.28%) are \geq 800 bp and 778 (2.20%) are \geq 1,000 bp, while 928 (2.62%) have
263 sizes \geq 2,000 bp.

264

265 **3.5 Long Non-Coding RNAs Annotation**

266 We performed an orthologous analysis of our putative lncRNAs using the ZFLNC
267 database and a custom ncRNA sequences database from *Danio rerio*, *Lates*
268 *calcarifer*, *Seriola lalandi dorsalis*, *Seriola dumerili*, *Echeneis naucrates* available at
269 Ensembl (see Methods, for details). Among the 206 ZFLNC hits, the most common
270 annotations were ZFLNCT01535, ZFLNCT11671, ZFLNCT19022 (matches with 5
271 transcripts, each), ZFLNCT02442 (matches with 4 transcripts), ZFLNCT13035 and
272 ZFLNCT16489 (each matching three putative lncRNAs). Interestingly, among the 79
273 hits for *D. rerio* ncRNA sequences present in Ensembl, 57 are annotated as long
274 intervening noncoding RNAs/Long intergenic noncoding RNAs (lincRNAs). The
275 remaining hits correspond to antisense (5), retained_intron (2), misc_RNA (1), sense
276 intronic (1) snoRNA (1) and 12 of them as processed transcript (Supplementary
277 Table ST1-S9). As mentioned above, the majority of transcripts described herein
278 displayed significant similarity to coding DNA of fish species phylogenetically related
279 to cobia, such as *Seriola dumerili*, *Seriola lalandi dorsalis*, *Echeneis naucrates* and
280 *Lates calcarifer*, during our initial annotation efforts. The lncRNA conservation search
281 performed here showed that 157 of 159 hits obtained from *E. naucrates* ncRNAs
282 correspond to sequences already identified as lncRNA in this species (the other 2
283 hits are annotated as small nucleolar RNAs U85 and SNORD10). On the other hand,
284 only few hits were found from *Seriola* spp. and *L. calcarifer*: (i) one snoRNA
285 appeared only for *L. calcarifer* and is the same lncRNA transcript matching *E.*
286 *naucrates* SNORD10; (ii) one misc_RNA, called 7SK RNA, was found for the same
287 transcript for these organisms – which is also the same hit found from *D. rerio*
288 (Supplementary Table ST1-S9).

289

290 **3.6 lncRNA Microsatellites**

291 Sequences from all 35,414 putative lncRNA transcripts were used to discover
292 potential microsatellites in the cobia genome with the aid of MISA. The SSR loci
293 detection was performed by searching for two to six nucleotide motifs, with a
294 minimum of 6,5,5,4 and 4 repeats, respectively. A total of 7,894 microsatellites were
295 detected in the putative lncRNAs (Supplementary Table ST1-S10). Among the
296 microsatellites, mono-nucleotide motifs were the most abundant type detected in
297 lncRNAs (55.41%). Other motifs included di-nucleotide (34.80%), tri-nucleotide
298 (5.9%), tetra-nucleotide (2.33%), penta-nucleotide (1.28%) and hexa-nucleotide
299 (0.28%) motifs (Fig. 1e). The mono-nucleotide repeat T was the most abundant motif

300 detected (50.07%), followed by A (46.27%), C (2.29%) and, finally, G (1.37%)
301 (Supplementary Table ST1-S10).

302

303 **4. CONCLUSION**

304 Our study has built the first liver transcriptome assembly of this important
305 commercial species, providing an important tool for further research with cobia. The
306 availability and deposition of the transcriptome sequence allows to access novel
307 gene sequences, contributing to gene expression assessments, and consequently
308 improving the knowledge regarding cobia physiology and nutrition, since the liver can
309 be considered the main lipogenic tissue in fish. In addition, the provided assembly
310 and genetic markers dataset will be essential as a base for future nutrigenomics
311 projects involving, genetic breeding programs and marker-assisted selection for this
312 species.

313 **DECLARATIONS**

314

315 **Availability of data and material**

316 Sequencing raw data were deposited in the Sequence Read Archive (SRA)
317 repository of the National Center for Biotechnology Information (NCBI), under
318 accession number SRR13009897, SRR13009896, SRR13009895, SRR13009894,
319 SRR13009893, SRR13009892, SRR13009891, SRR13009890, SRR13009889,
320 SRR13009888, SRR13009887 and SRR13009886, associated to the BioProject
321 numbers PRJNA675281 and BioSamples numbers SAMN16708758,
322 SAMN16708759, SAMN16708760, SAMN16708761, SAMN16708762,
323 SAMN16708763, SAMN16708764, SAMN16708765, SAMN16708766,
324 SAMN16708767, SAMN16708768 and SAMN16708769. The Transcriptome
325 Shotgun Assembly (TSA) project has been deposited at DDBJ/EMBL/GenBank
326 under accession number GIWT00000000. The version described in this paper is the
327 first version, GIWT00000000.1. Supplementary Table S1 is available from the
328 Figshare repository (10.6084/m9.figshare.14522781.v2). Additional data derived
329 from this study (including all intermediate data) are also available from the Open
330 Science Framework (OSF) repository (DOI: 10.17605/OSF.IO/BV3WA). Details
331 about the softwares and databases used are available in Supplementary Table ST1-
332 S11.

333

334 **Authors' contributions**

335 B.C.A. sampled the specimens. B.C.A., G.S.B., A.W.S.H. and R.G.M. performed
336 molecular analyses and sequencing. D.A.B., A.S.S., D.L.J., L.R.N. and F.B.M.
337 assembled and evaluated the transcriptome assembly and annotation. All authors
338 wrote the paper. All authors read and approved the final version of the manuscript.

339

340 **Ethics approval**

341 This study's experimental procedures were conducted according to the guidelines
342 and approval of the Mogi das Cruzes University Institutional Animal Care and Use
343 Ethics Committee (#008/2017).

344

345 **Funding**

346 This study was financed in part by the São Paulo Research Foundation (FAPESP:
347 2019/26018-0) and National Council for Scientific and Technological Development
348 (CNPq: 305493/2019-1). D.A.B., B.C.A and A.S.S. are recipients of scholarship
349 grants from Coordination for the Improvement of Higher Education Personnel
350 (CAPES). A.W.S.H. is recipient of CNPq productivity scholarships (304662/2017-8).

351

352 **Conflicts of interest**

353 The authors report no conflicts of interest. The authors alone are responsible for the
354 content and the writing of the paper.

355

356 References

357

- 358 1. Andrew SC, Primmer CR, Debes PV, Erkinaro J, Verta JP (2021) The Atlantic
359 salmon whole blood transcriptome and how it relates to major locus
360 maturation genotypes and other tissues. *Mar Genomics*.
361 <https://doi.org/10.1016/j.margen.2020.100809>
- 362 2. Andrews S (2010) FastQC: a quality control tool for high throughput sequence
363 data. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>. Accessed 21
364 April 2021
- 365 3. Arnold CR, Kaiser JB, Holt GJ (2002) Spawning of cobia *Rachycentron*
366 *canadum* in captivity. *J World Aquac Soc*. [https://doi.org/10.1111/j.1749-](https://doi.org/10.1111/j.1749-7345.2002.tb00496.x)
367 [7345.2002.tb00496.x](https://doi.org/10.1111/j.1749-7345.2002.tb00496.x)
- 368 4. Beier S, Thiel T, Münch T, Scholz U, Mascher M (2017) MISA-web: a web
369 server for microsatellite prediction. *Bioinformatics*.
370 <https://doi.org/10.1093/bioinformatics/btx198>
- 371 5. Benetti DD, Orhun, MR, Sardenberg B, O'Hanlon B, Welch A, Hoenig R, Zink
372 I, Rivera JA, Denlinger B, Bacoat D, Palmer K, Cavalin F (2008) Advances in
373 hatchery and grow-out technology of cobia *Rachycentron canadum*
374 (Linnaeus). *Aquac Res*. <https://doi.org/10.1111/j.1365-2109.2008.01922.x>
- 375 6. Boratyn GM et al. (2013) BLAST: a more efficient report with usability
376 improvements. *Nucleic Acids Res*. <https://doi.org/10.1093/nar/gkt282>.
- 377 7. Camargo AP et al. (2020) RNAsamba: neural network-based assessment of
378 the protein-coding potential of RNA sequences. *NAR Genom Bioinform*.
379 <https://doi.org/10.1093/nargab/lqz024>
- 380 8. Calduch-Giner JA, Bermejo-Nogales A, Benedito-Palos L, Estensoro I,
381 Ballester-Lozano G, Sitjà-Bobadilla A, Pérez-Sánchez J (2013) Deep
382 sequencing for *de novo* construction of a marine fish (*Sparus aurata*)
383 transcriptome database with a large coverage of protein-coding transcripts.
384 *BMC Genomics*. <https://doi.org/10.1186/1471-2164-14-178>.
- 385 9. Chen S, Zhou Y, Chen Y, Gu J (2018) fastp: an ultra-fast all-in-one FASTQ
386 preprocessor. *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/bty560>
- 387 10. Ewels P, Magnusson M, Lundin S, Käller M (2016) MultiQC: summarize
388 analysis results for multiple tools and samples in a single report.
389 *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/btw354>.

- 390 11. Fox SE, Christie MR, Marine M, Priest HD, Mockler TC, Blouin MS (2014)
391 Sequencing and characterization of the anadromous steelhead
392 (*Oncorhynchus mykiss*) transcriptome. Mar Genomics.
393 <https://doi.org/10.1016/j.margen.2013.12.001>.
- 394 12. FAO (2020) The State of World Fisheries and Aquaculture 2020.
395 Sustainability in action. <http://www.fao.org/documents/card/en/c/ca9229en>.
396 Accessed 21 April 2021
- 397 13. Fraser TWK, Davies SJ (2009) Nutritional requirements of cobia,
398 *Rachycentron canadum* (Linnaeus): a review. Aquac Res.
399 <https://doi.org/10.1111/j.1365-2109.2009.02215.x>
- 400 14. The Gene Ontology Consortium (2019) The gene ontology resource: 20 years
401 and still GOing strong. Nucleic Acids Res. <https://doi.org/10.1093/nar/gky1055>
- 402 15. Glencross BD, De Santis C, Bicskei B, Taggart JB, Bron JE, Betancor MB,
403 Tocher DR (2015) A comparative analysis of the response of the hepatic
404 transcriptome to dietary docosahexaenoic acid in Atlantic salmon (*Salmo*
405 *salar*) post-smolts. BMC Genomics. [https://doi.org/10.1186/s12864-015-1810-](https://doi.org/10.1186/s12864-015-1810-z)
406 [z](https://doi.org/10.1186/s12864-015-1810-z)
- 407 16. Fan G, Cao Y, Wang Z (2018) Regulation of Long Noncoding RNAs
408 Responsive to Phytoplasma Infection in *Paulownia tomentosa*. Int J
409 Genomics. <https://doi.org/10.1155/2018/3174352>
- 410 17. Gui D et al. (2013) De novo assembly of the Indo-Pacific humpback dolphin
411 leucocyte transcriptome to identify putative genes involved in the aquatic
412 adaptation and immune response. PLoS One.
413 doi:10.1371/journal.pone.0072417
- 414 18. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J,
415 Couger MB, Eccles D, Li B, Lieber M, MacManes MD (2013) *De novo*
416 transcript sequence reconstruction from RNA-seq using the Trinity platform for
417 reference generation and analysis. Nat Protoc.
418 <https://doi.org/10.1038/nprot.2013.084>.
- 419 19. Hart AJ, Ginzburg S, Xu M, Fisher CR, Rahmatpour N, Mitton JB, Paul R,
420 Wegrzyn JL (2020) EnTAP: bringing faster and smarter functional annotation
421 to non-model eukaryotic transcriptomes. Mol Ecol Resour.
422 <https://doi.org/10.1111/1755-0998.13106>.

- 423 20. Herkenhoff ME et al. (2018) Fishing into the MicroRNA transcriptome.
424 Frontiers in genetics. <https://doi.org/10.3389/fgene.2018.00088>
- 425 21. Hu X et al. (2018) ZFLNC: a comprehensive and well-annotated database for
426 zebrafish lncRNA. Database. <https://doi.org/10.1093/database/bay114>
- 427 22. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC,
428 Rattei T, Mende DR, Sunagawa S, Kuhn M, Jensen LJ (2016) eggNOG 4.5: a
429 hierarchical orthology framework with improved functional annotations for
430 eukaryotic, prokaryotic and viral sequences. Nucleic Acids Res.
431 <https://doi.org/10.1093/nar/gkv1248>.
- 432 23. Kalvari I et al. (2018) Rfam 13.0: shifting to a genome-centric resource for
433 non-coding RNA families. Nucleic Acids Res.
434 <https://doi.org/10.1093/nar/gkx1038>
- 435 24. Kanehisa M, Goto S (2000) KEGG: kyoto encyclopedia of genes and
436 genomes. Nucleic acids research. <https://doi.org/10.1093/nar/28.1.27>.
- 437 25. Kanehisa M (2019) Toward understanding the origin and evolution of cellular
438 organisms. Protein Sci. <https://doi.org/10.1002/pro.3715>.
- 439 26. Kanehisa M, Furumichi M, Sato Y, Ishiguro-Watanabe M, Tanabe M (2021)
440 KEGG: integrating viruses and cellular organisms. Nucleic Acids Res.
441 <https://doi.org/10.1093/nar/gkaa970>
- 442 27. Kang YJ et al. (2017) CPC2: a fast and accurate coding potential calculator
443 based on sequence intrinsic features. Nucleic Acids Res.
444 <https://doi.org/10.1093/nar/gkx428>
- 445 28. Leaver MJ, Bautista JM, Björnsson BT, Jönsson E, Krey G, Tocher DR,
446 Torstensen BE (2008) Towards fish lipid nutrigenomics: current state and
447 prospects for fin-fish aquaculture. Rev Fish Sci.
448 <https://doi.org/10.1080/10641260802325278>.
- 449 29. Magnanou E, Klopp C, Noirot C, Besseau L, Falcón J (2014) Generation and
450 characterization of the sea bass *Dicentrarchus labrax* brain and liver
451 transcriptomes. Gene. <https://doi.org/10.1016/j.gene.2014.04.032>.
- 452 30. Marz M et al. (2011) Animal snoRNAs and scaRNAs with exceptional
453 structures. RNA Biol. <https://doi.org/10.4161/rna.8.6.16603>
- 454 31. Menegidio FB, Jabes DL, Costa de Oliveira R, Nunes LR (2018) Dugong: a
455 Docker image, based on Ubuntu Linux, focused on reproducibility and

- 456 replicability for bioinformatics analyses. *Bioinformatics*.
457 <https://doi.org/10.1093/bioinformatics/btx554>.
- 458 32. Nawrocki EP, Eddy SR (2013) Infernal 1.1: 100-fold faster RNA homology
459 searches. *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/btt509>.
- 460 33. Nunes AJP (2014) Ensaios com o beijupirá, *Rachycentron canadum*.
461 Fortaleza: Ministério da Pesca e Aquicultura/CNPQ/UFC.
462 <http://www.repositorio.ufc.br/handle/riufc/8655>. Accessed 21 April 2021
- 463 34. Rasal KD et al. (2016) MicroRNA in aquaculture fishes: a way forward with
464 high-throughput sequencing and a computational approach. *Rev Fish Biol*
465 *Fish*. <https://doi.org/10.1007/s11160-016-9421-6>.
- 466 35. Seppey M, Manni M, Zdobnov EM (2019) BUSCO: assessing genome
467 assembly and annotation completeness. *Methods Mol Biol*.
468 https://doi.org/10.1007/978-1-4939-9173-0_14.
- 469 36. UniProt Consortium (2019) UniProt: a worldwide hub of protein knowledge.
470 *Nucleic acids research*. <https://doi.org/10.1093/nar/gky1049>
- 471 37. Wang L et al. (2013) CPAT: Coding-Potential Assessment Tool using an
472 alignment-free logistic regression model. *Nucleic Acids Res*.
473 <https://doi.org/10.1093/nar/gkt006>
- 474 38. Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simao FA,
475 Ioannidis P, Seppey M, Loetscher A, Kriventseva EV (2017) OrthoDB v9. 1:
476 cataloging evolutionary and functional annotations for animal, fungal, plant,
477 archaeal, bacterial and viral orthologs. *Nucleic Acids Res*.
478 <https://doi.org/10.1093/nar/gkw1119>.
479
480
481

482 **Table 1.** Summary of the *de novo* transcriptome assembly for *Rachycentron*
 483 *canadum*.

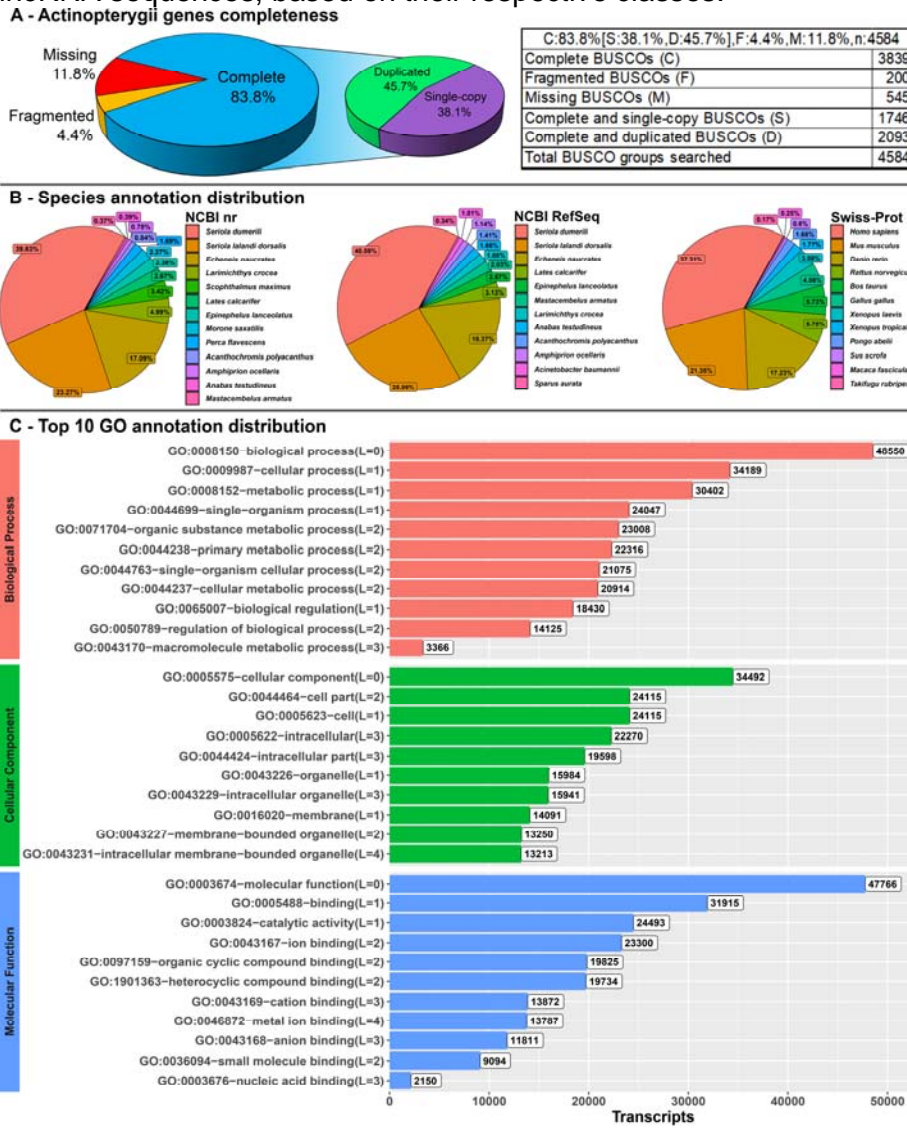
| Illumina sequencing | Raw | High-Quality |
|--------------------------------------|----------------------------|---------------------|
| Number of reads | 1,761,965,794 | 1,652,319,304 |
| Assembly | Unigenes | Isoforms |
| Total sequences | 101,789 | 163,096 |
| N50 | 7,843 | 2,312 |
| Average contig length (bp) | 1,617.34 | 950.61 |
| Median contig length (bp) | 704 | 391 |
| Total assembled bases | 263,781,171 | 96,761,286 |
| CG% | 44.5 | |
| Coding potential | Transcripts | |
| Undetermined transcripts | 968 | |
| Coding transcripts | 126,013 | |
| Non-coding transcripts | 701 | |
| Putative long non-coding transcripts | 35,414 | |
| Functional Annotation | Annotated Sequences | |
| NCBI nr | 29,155 | |
| NCBI RefSeq RNA | 122,741 | |
| Swiss-Prot | 15,507 | |
| Lates/Seriola | 118,570 | |
| GO-Biological Process | 48,550 | |
| GO-Molecular Function | 47,766 | |
| GO-Cellular Component | 34,492 | |
| Kegg | 23,936 | |
| Infernal | 699 | |
| ZFLNC | 206 | |
| ncRNA Ensembl DB | 257 | |

484

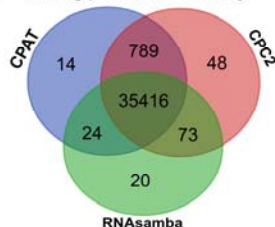
485

486 **Fig. 1** Completeness and homology search of *Rachycentron canadum* liver
 487 transcripts. (A) Percentage of completeness on the core set of genes from *R.*
 488 *canadum* based on Actinopterygii database (orthodb9), using BUSCO. (B) Species
 489 annotation distribution for the best hits from NCBI nr, NCBI RefSeq and Swiss-Prot
 490 databases. (C) Gene ontology distribution for Biological Process, Cellular
 491 Component, Molecular Function categories of assembled transcripts from the *R.*
 492 *canadum* liver transcriptome. (D) Distribution of transcripts classified as non-protein
 493 coding elements. (E) Distribution of microsatellites (SSRs) identified in the putative
 494 lncRNA sequences, based on their respective classes.

495



D - Coding potential summary



E - Distribution of SSRs by type class

