1

A Hybrid *de novo* Assembly of the Sea Pansy (*Renilla muelleri*) Genome

3

Justin Jiang[1], Andrea M. Quattrini[1*], Warren R. Francis[2], Joseph F. Ryan[3], Estefanía Rodríguez[4],

Catherine S. McFadden[1]

6

[1]Department of Biology, Harvey Mudd College, 1250 N. Dartmouth Ave, Claremont, CA 91711,

USA

[2] University of Southern Denmark, Dept. of Biology, Campusvej 55, Odense M 5230, Denmark

[3] Whitney Laboratory for Marine Bioscience, University of Florida, 9505 Ocean Shore Blvd.

St. Augustine, FL 32080, USA

[4] Division of Invertebrate Zoology, American Museum of Natural History, Central Park West at

79th Street, New York, NY 10024, USA

14

Justin Jiang: jjiang990@gmail.com

Andrea Quattrini: aquattrini@g.hmc.edu

Warren R. Francis: wfrancis@biology.sdu.dk

Joseph F. Ryan: joseph.ryan@whitney.ufl.edu

Estefanía Rodríguez: erodriguez@amnh.org

Catherine S. McFadden: mcfadden@g.hmc.edu

21

*Corresponding Author

23

24    **Abstract**

25    **Background:** Over 3,000 species of octocorals (Cnidaria, Anthozoa) inhabit an expansive range

26    of environments, from shallow tropical seas to the deep-ocean floor. They are important

27    foundation species that create coral "forests" which provide unique niches and three-dimensional

28    living space for other organisms. The octocoral genus *Renilla* inhabits sandy, continental shelves

29    in the subtropical and tropical Atlantic and eastern Pacific Oceans. *Renilla* is especially

30    interesting because it produces secondary metabolites for defense, exhibits bioluminescence, and

31    produces a luciferase that is widely used in dual-reporter assays in molecular biology. Although

32    several cnidarian genomes are currently available, the majority are from hexacorals. Here, we

33    present a *de novo* assembly of the *R. muelleri* genome, making this the first complete draft

34    genome from an octocoral.

35    **Findings:** We generated a hybrid *de novo* assembly using the Maryland Super-Read Celera

36    Assembler v.3.2.6 (MaSuRCA). The final assembly included 4,825 scaffolds and a haploid

37    genome size of 172 Mb. A BUSCO assessment found 88% of metazoan orthologs present in the

38    genome. An Augustus *ab initio* gene prediction found 23,660 genes, of which 66% (15,635) had

39    detectable similarity to annotated genes from the starlet sea anemone, *Nematostella vectensis,* or

40    to the Uniprot database. Although the *R. muelleri* genome is smaller (172 Mb) than other

41    publicly available, hexacoral genomes (256-448 Mb), the *R. muelleri* genome is similar to the

42    hexacoral genomes in terms of the number of complete metazoan BUSCOs and predicted gene

43    models.

44    **Conclusions:** The *R. muelleri* hybrid genome provides a novel resource for researchers to

45    investigate the evolution of genes and gene families within Octocorallia and more widely across

46  Anthozoa. It will be a key resource for future comparative genomics with other corals and for

47  understanding the genomic basis of coral diversity.

48

49  Keywords: octocoral, hybrid assembly, gene prediction, Augustus, PacBio, MaSuRCA

50

51  **Data Description**

52  *Organism Description*

53      Octocorallia is a subclass of Anthozoa (Phylum: Cnidaria) that is comprised of three

54  orders: Alcyonacea, Helioporacea, and Pennatulacea [1]. The Pennatulacea, commonly known as

55  sea pens, are a monophyletic group [1, 2] and are the most morphologically distinct group of

56  octocorals [1, 3]. Sea pens differ from other octocorals by exhibiting the most integrated colonial

57  behavior, with colonies arising from an axial polyp that develops into a peduncle—used to

58  anchor the animal into soft-sediments or onto hard surfaces—and a rachis that supports

59  secondary polyps [1, 3-4]. There are 14 valid families of Pennatulacea distinguished by the

60  arrangement of the secondary polyps around the rachis [1, 4]. The monogeneric family

61  Renillidae Lamarck, 1816 consists of seven species [5], unique because of their foliate colony

62  growth form [1, 4].

63      *Renilla* is found naturally on sandy, shallow sea floors along the Atlantic and Pacific

64  coasts of North and South America [3, 4, 6]. The brilliant bioluminescence and endogenous

65  fluorescence of these animals have led to them becoming important organisms in microscopy

66  and molecular biology. Isolated originally from *R. reniformis*, the enzyme luciferase (Renilla-

67  luciferin 2-monooxygenase) is used in dual luciferase reporter assays, which are commonly used

68  to study gene regulation and expression, signaling pathways, and the structure of regulatory

69    genes [7-8]. The green fluorescent protein from *Renilla* has medical applications as well as

70    general molecular biology and imagery uses [9]. In addition, the compounds produced by *Renilla*

71    for chemical defense [10] may be important sources for discovery of marine natural products

72    [11]. Thus, a genome of the octocoral *Renilla* is highly valuable to the scientific community,

73    providing a novel resource that has a range of important uses— from molecular biology to

74    comparative genomics.

75          Due to the known difficulties of resolving lengthy repeat regions with Illumina-only data

76    [12-13], we used a hybrid assembly approach [13-14], combining long-read Pacific Biosciences

77    (PacBio) data with short-read Illumina data. Studies have shown that a hybrid approach results in

78    a more complete assembly with less genome fragmentation [15-17]. Our hybrid approach used

79    low coverage PacBio reads (15x coverage) along with high coverage Illumina HiSeq reads (105x

80    coverage) to assemble a draft genome of *R. muelleri* Schultze in Kölliker, 1872, a sea pen

81    common to shallow waters of the Gulf of Mexico [6].

82

83    **Methods and Results**

84    *Data Collection*

85          A live specimen of *R. muelleri* was obtained from Gulf Specimen Marine Lab (Panacea,

86    FL, USA), which collects specimens off the panhandle of Florida in the Gulf of Mexico. Upon

87    receiving the specimen, it was flash frozen in liquid nitrogen. Genomic DNA was then extracted

88    using a modified CTAB protocol [18]. A total of 5.6 μg of DNA was sent to Novogene

89    (Sacramento, CA, USA) for library preparation and sequencing. 350 bp insert DNA libraries that

90    were PCR free were prepared and then multiplexed with other organisms on two lanes of an

91    Illumina HiSeq 2500 (150 bp PE reads). In addition, Illumina MiSeq and PacBio sequencing

92    were performed at the Weill Cornell Medicine Epigenomics Core Facility in New York. For the

93    Illumina MiSeq run, the *Renilla* library was prepared with TruSeq LT and then multiplexed with

94    eight other corals and sequenced (300 bp PE reads, MiSeq v3 Reagent kit). For PacBio

95    sequencing, a DNA library was prepared from 5 ug of DNA using the SMRTbell template prep

96    kit v 1.0. Sequencing was carried out on 10 SMRT cells on a RSII instrument using P6-C4

97    chemistry. PacBio SMRT Analysis 2.3 subread filtering module was used to produce the subread

98    files for assembly.

99        As part of another study, we sequenced total RNA from a congeneric species, *R.*

100   *reniformis*. The specimen was collected alive on the beach in North Flagler County, Florida,

101   USA. RNA was extracted from the whole adult colony and sequenced on a NextSeq500 (150 bp

102   PE reads) instrument. Library preparation and sequencing were performed at the University of

103   Florida's Interdisciplinary Center for Biotechnology.

104

105   *DNA Read Processing*

106       A total of 246,744,426 PE reads were obtained from the HiSeq and 6,725,072 PE reads

107   were obtained from the MiSeq. In total, we generated 39,029,185,500 bases of Illumina data.

108   Adapters were trimmed from all raw Illumina reads using Trimmomatic v.0.35

109   (*ILLUMINACLIP:2:30:10 LEADING:5 TRAILING:5 SLIDINGWINDOW:4:20 MINLEN:3*;

110   Trimmomatic, RRID:SCR_011848) [19], resulting in 38.98 Gb of reads. These reads were then

111   filtered with Kraken v.1.0 (Kraken, RRID:SCR_005484) [20] using the MiniKraken 8GB

112   database [21] to screen for possible microbial, viral and archaeal contamination. A total of 960

113   Mb were removed from the read files, resulting in 36.23 Gb of 150 bp reads and 1.79 Gb of 300

114   bp reads.

115    A total of 1,227,306 PacBio subreads were obtained and screened against the NCBI

116    environmental nucleotide database (env_nt.00 to env_nt.23) [22] using BLASTn v.2.2.31 (-

117    *evalue 1e-10, -out_fmt 5*, RRID:SCR_001598) [23] to identify reads with environmental

118    contaminants. The subreads that did not contain contaminants were extracted using

119    MEtaGenome ANalyzer v.6.4.16 (MEGAN, RRID:SCR_011942) [24-25], resulting in 5.22 Gb

120    in 1,195,521 reads.

121

122    *RNA-Seq Read Processing*

123    We generated 119,604,588 PE reads of RNA-Seq data. We used Trimmomatic (version

124    0.36 (-phred33, ILLUMINACLIP:/usr/local/Trimmomatic-0.32/adapters/TruSeq3-

125    PE.fa:2:30:12:1:true, MINLEN:36; Trimmomatic, RRID:SCR_011848) [19] to remove Illumina

126    adapters. Trinity v 2.4.0 (--seqType fq --max_memory 250G --CPU 6 --left trim.R1.fq --right

127    trim.R2.fq --full_cleanup; Trinity RRID_SCR:013048) [26] was used to assemble the

128    transcriptome.

129

130    *Hybrid Genome Assembly*

131    Two hybrid *de novo* assemblies were performed, one with the Maryland Super-Read

132    Celera Assembler v.3.2.6 (MaSuRCA, RRID:SCR_010691) [27] and the other with SPAdes

133    v.3.11.0 (SPAdes, RRID:SCR_000131; *k-mer lengths 21,33,55,77*) [28]. The Benchmarking

134    Universal Single-Copy Orthologs v.3.0.2 (BUSCO, RRID:SCR_015008) [29] program with

135    default settings was used to screen the *Renilla* genome assemblies for 978 orthologs from the

136    Metazoan dataset as a method to evaluate the completeness of each assembly. BUSCO used

137    BLAST v.2.2.31 [23] and HMMER v.3.1.b2 (HMMER, RRID:SCR_005305) [30] in its pipeline.

138    The stats.sh program from BBMAP v.36.14 (bbmap) [31] was used to generate general assembly

139    statistics for genomes produced by both programs (Table 1).

140        The MaSuRCA assembly resulted in a 147-fold decrease in the number of scaffolds

141    generated, and a 70-fold increase in the N50 contig size (70.522 KB) as compared to the SPades

142    assembly (1.007 KB); it also had more complete BUSCOs present (Table 1). Other statistics also

143    indicate that the MaSuRCA assembly is much less fragmented than the SPAdes assembly (Table

144    1). Therefore, we used the MaSuRCA assembly in further analyses.

145        To improve the quality of the draft MaSuRCA assembly, six iterations of Pilon v.1.21

146    (Pilon, RRID:SCR_014731) [32] were used to fix assembly errors and fill assembly gaps.

147    Bowtie2 v.2.3.2 (Bowtie2, RRID:SCR_016368) [33] was used to align Illumina HiSeq and

148    Illumina MiSeq genomic reads to the draft assembly, and the resulting alignments were input to

149    Pilon, which was run on default settings. A total of 52,668 SNPs were corrected, along with

150    14,702 small insertions and 11,841 small deletions (Supplementary Table S1).

151        To remove haploid contigs that were not merged during assembly, we ran BLASTn

152    against the contigs themselves *(-max_target_seqs 10, -evalue 1e-40)* to find contigs that were

153    highly similar. The custom script *haplotypeblastn.py* version 1.0 [34] filtered the BLASTn

154    results by flagging matches that were greater than 75% identical and longer than 500 bp in

155    length. The contigs that were identified as unmerged were subsequently removed using the

156    *select_contigs.pl* script [35]. A total of 59 scaffolds, which amounted to 67 contigs and 384 kb,

157    were removed from the assembly.

158        The bbmap program stats.sh was used to generate assembly statistics on the haplotype-

159    removed assembly [i.e., "final assembly", (Table 1)]. BUSCO analysis using the metazoan

160    orthologs was again used to estimate the completeness of the final assembly, with the flag *-long*

7

161    to produce higher quality training data for the downstream annotation. 857 (87.63%) orthologs

162    were present in the final assembly (Table 1). This final *R. muelleri* assembly was masked, using

163    RepeatMasker v.open-4.0.6 (-*species eukaryota -gccalc -div 50;* RepeatMasker,

164    RRID:SCR_012954) [36], for downstream gene annotation. The final annotation consists of

165    172,512,580 bp in 4,925 scaffolds.

166

*Genome Annotation*

168        Stampy v.1.0.31 (Stampy, RRID:SCR_005504) [37] was used to align 18.06 Gb of RNA-

169    Seq data from *R. reniformis* to the masked genome to generate intron hints. The resulting bam

170    file was processed by filtering out raw alignments using *filterBam* [38] per the recommended

171    Augustus procedures [39]. A total of 1,837,637 intron hints were generated.

172        Augustus v.3.3 (--*UTR=off --allow_hinted_splicesites=atac --alternatives-from-*

173    *evidence=true;* Augustus, RRID:SCR_008417) [40] was used to predict a gene model for *R.*

174    *muelleri*. Augustus training was performed with the hint data from *R. reniformis*, as it has been

175    shown to improve *ab initio* predictions [40-41]. The BUSCO-generated training data was also

176    included to help predict a gene model. A modified extrinsic weight file was used in Augustus to

177    penalize predicted introns that were unsupported by hint evidence and reward predicted introns

178    that were supported by hint evidence by 1e2.

179        Augustus predicted 23,660 genes that had an average exon length of 249 bp and an

180    average intron length of 524 bp as calculated by *gfstats.py* [42] (Table 2). BUSCO with the

181    metazoan lineage (-*m prot*) orthologs was used to assess the quality of the prediction, finding

182    84.87% (830/978) orthologs (Table 2).

183

184    *Functional Annotations*

185        BLASTp v.2.2.31+ (*-evalue 1e-10 -seg yes -soft_masking true -lcase_masking*, BLASTp,

186    RRID:SCR_001010) [23] was used to map the predicted gene models of *R. muelleri* to filtered

187    protein models of another anthozoan, the sea anemone, *Nematostella vectensis* (Joint Genome

188    Institute, JGI, v 1.0) [43]. A total of 63% (14,931) of the predicted genes (23,660) mapped to *N.*

189    *vectensis* proteins (27,273). A custom python script, *filterGenes.py* [44] was used to filter the

190    matches by selecting the highest bit score; in cases where bit scores were identical, the match

191    with the highest percent length of all matches was used as a tiebreaker. Of the 14,931 genes that

192    mapped to *N. vectensis* proteins, 12,279 genes were annotated with GO function, KOG function

193    and/or InterPro domains; 8,101 genes were assigned GO terms; 11,067 genes were assigned

194    KOG functions; and 10,126 genes were assigned InterPro domains (Supplemental File 1). The

195    8,729 genes that did not hit *N. vectensis* proteins were remapped with BLASTp using a lower e-

196    value (*1e-5*) and filtered with the aforementioned python script with the same settings; an

197    additional 2,002 of the genes mapped to *N. vectensis*. Of these, 1,512 genes were annotated with

198    GO functions, KOG functions and/or InterPro domains (Supplemental File 1). The remaining

199    6,727 genes that did not match *N. vectensis* annotations were mapped to the UniProt database

200    (UniProt, RRID:SCR_002380) [45-46] with BLASTp (*-evalue 1e-5*), and 1,844 of these were

201    assigned a UniProt function. In total, 79.36% (18,777/23,660) of the predicted gene models were

202    mapped to either *N. vectensis* predicted proteins or the UniProt database, and 66.08%

203    (15,635/23,660) of the predicted *Renilla* genes have either functional annotations or InterPro

204    domain information associated with them.

205        We also used BLASTp (*-evalue 1e-10 -seg yes -soft_masking true -lcase_masking)* to

206    map the predicted genes against a newer *N. vectensis* dataset that was generated using RNA-Seq

9

207 (hereafter called the Vienna dataset) [47-48]. A total of 63% (15,001) of the predicted genes

208 (23,660) mapped to the Vienna dataset (25,729) (Supplemental File 2). As above, the predicted

209 genes that did not map were remapped with a lower e-value (*1e-5)*, resulting in 2,071 additional

210 predicted genes mapping to the *N. vectensis* Vienna dataset. In total, 72.15% (17,072) of

211 predicted genes mapped to the Vienna dataset. This dataset did not have associated functional

212 annotations. Combining all gene model annotation methods, 79.82% (18,886) of genes from the

213 Augustus gene model were mapped to the JGI *N. vectensis* annotations, the *N. vectensis* Vienna

214 dataset, or the UniProt database (Supplemental Files 1-3).

215

216 *Genome Assembly Comparisons*

217 We compared the *R. muelleri* genome assembly to previously published anthozoan (e.g.,

218 corals, anemones) genomes using a variety of assessment statistics (Supplemental Table S2).

219 BUSCO was used (*-m geno)* to assess the completeness of six hexacoral genomes and a draft *R.*

220 *reniformis* genome and compare these results to the hybrid *R. muelleri* assembly (Fig. 1). We

221 found the BUSCO-completeness of our *R. muelleri* assembly (857) to be more similar to the

222 well-curated assembly of the model organism *N. vectensis* (893) [49-50] than to the other

223 anthozoans. BUSCOs from the other five hexacoral genomes were less complete, with complete

224 BUSCOs ranging from 728 (*Acropora digitifera*) to 839 (*Discosoma* sp.) [50-57]. Only 800

225 complete BUSCOs were recovered from the other hybrid assembly, the hexacoral *Montastraea*

226 *cavernosa* [57]. The only other publicly-available octocoral genome, *R. reniformis*, had

227 considerably fewer complete BUSCOs (356, Fig.1) [58].

228 The number of predicted genes was highly similar across all anthozoan genomes

229 (Supplementary Table S2). The range of predicted genes was 21,372 to 30,360 across the six

230  hexacorals. The number of predicted genes (23,360) for *R. muelleri* was similar to the 23,668

231  genes predicted for *A. digitifera.*

232  Interestingly, the genome size of *R. muelleri* is considerably smaller (172 Mb) than other

233  hexacoral genomes (256-448 Mb). Of the hexacorals, the anemone *Exaiptasia pallida* has the

234  smallest genome size of 256 Mb, while the others have genome sizes >300 Mb. As indicated by

235  [56], *E. pallida* has smaller and less frequent introns. Similar to *E. pallida,* exon sizes were

236  larger in *R. muelleri* (249 bp) compared to the hexacorals (208 to 230 bp). These results suggest

237  that there may be comparatively fewer non-coding regions in *R. muelleri* because the number of

238  predicted gene models in *R. muelleri* is similar to hexacorals, yet the exon sizes are larger and

239  the genome size is smaller in *R. muelleri.* In addition, repetitive elements in the *R. muelleri*

240  genome may be less frequent, however, this remains to be further examined.

241  We also compared the mitochondrial genome to the previously published mitogenome of

242  *R. muelleri* [59]. We used BLASTn to search for the mitogenome among the contigs (included as

243  the last contig in the assembly) and recovered the entire 18,641 bp circularized, mitogenome.

244  Compared to the published mitogenome, there were just two, single bp differences and one bp

245  indel.

246

**Conclusions**

248  We present the first octocoral genome assembly and showcase the feasibility of the

249  MaSuRCA hybrid assembler for marine invertebrate genomics. The *R. muelleri* genome is one of

250  the smallest anthozoan genomes discovered to date, yet it is comparable to other coral and

251  anemone genomes in terms of predicted gene models. The identification of 88% of complete

252  metazoan BUSCOs in the *R. muelleri* genome highlights that a quality genome assembly can be

253    obtained from relatively low coverage sequencing of short and long read data. Although more

254    data are needed to further increase size and reduce number of scaffolds, and further functional

255    annotation is needed, the genome of the sea pansy, *R. muelleri,* provides a novel resource for the

256    scientific community to further investigations of gene family evolution, comparative genomics,

257    and the genomic basis of coral diversity.

258

259    **Availability of supporting data**

260    The final hybrid assembly and predicted proteins generated by this study are in the *Giga*DB

261    repository [60] and on the reefgenomics website [61]. Raw Illumina and PacBio reads are

262    available in NCBI's Sequence Read Archive (PRJNA491947). RNA-Seq reads have been

263    uploaded to the European Nucleotide Archive (PRJEB28688).

264

265    **Abbreviations**

266         bp: base pair, BUSCO: Benchmarking Universal Single-Copy Orthologs, Gb: gigabp,

267    Mb: megabp, MY: million years, PE: paired end, Pacbio: Pacific Biosciences

268

269    **Additional Files**

270    **Supplementary Table S1.** Summary of Pilon changes per iteration

271    **Supplemental Table S2.** *Renilla muelleri* genome assembly and annotation comparisons to

272    other anthozoan genomes.

273    **Supplemental File 1.** Gene model annotations of *Renilla muelleri* using the *Nematostella*

274    *vectensis* Joint Genome Institute filtered protein model.

275    **Supplemental File 2.** Gene annotations of *Renilla muelleri* using the *Nematostella vectensis*

276 Vienna dataset.

277 **Supplemental File 3.** Reference file that includes annotations for the predicted gene models.

278 This dataset includes GO terms, KOG IDs, and InterPro domains as annotated in the

279 *Nematostella vectensis* filtered protein models (Joint Genome Institute).

280

281 **Competing interests**

282 The authors declare no competing interests.

283

284 **Funding**

285 This study was funded by NSF-DEB Award 1457817 to C.S. McFadden and NSF-DEB Award

286 1457581 to E. Rodriguez. Additional funding came from startup funds from the University of

287 Florida DSP Research Strategic Initiatives #00114464 and University of Florida Office of the

288 Provost Programs to J.F. Ryan.

289

290 **Authors' Contributions**

291 **Justin Jiang**: Conceptualization, Investigation, Formal Analysis, Software Programming,

292 Methodology, Validation, Data Curation, Writing - Original Draft Preparation, Writing - Review

293 & Editing, Visualization

294 **Andrea M. Quattrini**: Conceptualization, Supervision, Investigation, Formal Analysis,

295 Methodology, Validation, Data Curation, Writing - Original Draft Preparation, Writing - Review

296 & Editing, Visualization

297 **Warren R. Francis**: Software Programming, Methodology, Validation, Writing - Review &

298 Editing

299    **Joseph F. Ryan**: Methodology, Validation, Data Curation, Writing - Review & Editing

300    **Estefania Rodriguez**: Conceptualization, Writing - Review & Editing

301    **Catherine S. McFadden**: Conceptualization, Supervision, Writing - Original Draft Preparation,

302    Writing - Review & Editing

303

309

## References

311    1. Daly M, Brugler MR, Cartwright P et. al. The phylum Cnidaria: A review of

312    phylogenetic patterns and diversity 300 years after Linnaeus. Zootaxa. 2007;1668:127-

313    182.

314    2. McFadden CS, France SC, Sánchez JA et. al. A molecular phylogenetic analysis of the

315    Octocorallia (Cnidaria: Anthozoa) based on mitochondrial protein-coding sequences.

316    Molecular Phylogenetics and Evolution. 2006;41(3):513:527.

317    3. Williams GC. The global diversity of sea pens (Cnidaria: Octocorallia: Pennatulacea).

318    PLoS One. 2011;6:e22747

319    4. Williams GC. Living genera of sea pens (Coelenterata: Octocorallia: Pennatulacea):

320    illustrated key and synopsis. Zoological Journal of the Linnean Society. 1995;113:93-

321    140.

322     5.  World Register of Marine Species: World List of Octocorallia Renillidae.

323         http://marinespecies.org/aphia.php?p=taxdetails&id=266953, Accessed 19 Aug 2018.

324     6.  Cairns SD, Bayer FM. Octocorallia (Cnidaria) of the Gulf of Mexico. In: Felder DL,

325         Camp DK, editors. Gulf of Mexico–Origins, Waters, and Biota. Volume 1.

326         Biodiversity. College Station, Texas: Academic; 2009:321-331.

327     7.  Sherf BA, Navarro SL, Hannah RR, Wood KV. Dual-luciferase reporter assay: an

328         advanced co-reporter technology integrating firefly and Renilla luciferase

329         assays. Promega Notes. 1996;56:2.

330     8.  Saito K, Chang YF, Horikawa K et al. Luminescent Proteins for High-Speed Single-Cell

331         and Whole-Body Imaging. Nature Communications. 2012; doi:10.1038/ncomms2248.

332     9.  Stepanenko OV, Verkhusha VV, Kuznetsova IM, Uversky VN, Turoverov KK. Current

333         Protein & Peptide Science. 2008; doi:10.2174/138920308785132668

334     10. Clavico EE, De Souza AT, Da Gama BA, Pereira RC. Antipredator defense and

335         phenotypic plasticity of sclerites from *Renilla muelleri*, a tropical sea pansy. The

336         Biological Bulletin, 2007;213(2):135-140.

337     11. Ledoux JB, Antunes A. Beyond the beaten path: improving natural products

338         bioprospecting using an eco-evolutionary framework–the case of the octocorals. Critical

339         Reviews in Biotechnology. 2018;38(2):184-198.

340     12. Pop M, Salzberg SL. Bioinformatics challenges of new sequencing technology. Trends in

341         Genetics. 2008;24(3):142-149.

342     13. Koren S, Schatz MC, Walenz BP, Martin J, Howard JT, Ganapathy G, Phillippy AM.

343         Hybrid error correction and de novo assembly of single-molecule sequencing

344         reads. Nature biotechnology, 2012;30(7):693.

345    14. English AC, Richards S, Han Y et al. Mind the gap: Upgrading genomes with Pacific

346         Biosciences RS long-read sequencing technology. PLoS ONE. 2012;

347         doi:10.1371/journal.pone.0047768.

348    15. Bashir A, Klamer AA, Robins WP et al. A hybrid approach for the automated finishing of

349         bacterial genomes. Nature Biotechnology. 2012; doi:10.1038/nbt.2288.

350    16. Giordano F, Aigrain L, Quail MA et al. De novo yeast genome assemblies from MinION,

351         PacBio and MiSeq platforms. Scientific Reports. 2017; doi:10.1038/s41598-017-03996-z.

352    17. Tan MH, Austin CM, Hammer MP et al. Finding Nemo: hybrid assembly with Oxford

353         Nanopore and Illumina reads greatly improves the clownfish (*Amphiprion ocellaris*)

354         genome assembly. GigaScience. 2018; doi:10.1093/gigascience/gix137.

355    18. McFadden CS, Alderslade P, Ofwegen LP van, Johnsen H, Rusmevichientong A.

356         Phylogenetic relationships within the tropical soft coral genera *Sarcophyton* and

357         *Lobophytum* (Anthozoa, Octocorallia). Invertebrate Biology 2006;125:288-305.

358    19. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence

359         data. Bioinformatics 2014;30(15):2114–20.

360    20. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using

361         exact alignments. Genome Biology. 2014;15(3):R46.

362    21. Wood DE. Minikraken 8 GB database, Johns Hopkins University,

363         https://ccb.jhu.edu/software/kraken/dl/minikraken_20171019_8Gb.tgz (August 7 2018,

364         date last accessed)

365    22. National Center for Biotechnology Information: Trivial HTTP: env_nt.00 to env_nt.23.

366         ftp://ftp.ncbi.nlm.nih.gov/blast/db/

367    23. Boratyn GM, Camacho C, Cooper PS et al. BLAST: a more efficient report with usability

368        improvements. Nucleic Acids Research 2013;41(W1):W29–33.

369    24. Huson DH, Mitra S, Ruscheweyh HJ et al. Integrative analysis of environmental

370        sequences using MEGAN4, Genome Research, 2011;21:1552-1560.

371    25. Huson DH, Auch AF, Qi J et al. MEGAN analysis of metagenomic data, Genome

372        Research, 2007;17(3):377-86.

373    26. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, MacManes

374        MD. De novo transcript sequence reconstruction from RNA-seq using the Trinity

375        platform for reference generation and analysis. Nature protocol. 2013;8(8):1494.

376    27.  Zimin AV, Marçais G, Puiu D et al. The MaSuRCA genome assembler. Bioinformatics

377        2013;29(21):2669–77.

378    28. Bankevich A, Nurk S, Antipov D, et al. SPAdes: A New Genome Assembly Algorithm

379        and Its Applications to Single-Cell Sequencing; Journal of Computational Biology. 2012;

380        doi:10.1089/cmb.2012.0021

381    29. Simão FA, Waterhouse RM, loannidis P et al. BUSCO: Assessing Genome Assembly

382        and Annotation Completeness with Single-Copy Orthologs. Bioinformatics. 2015;

383        doi:10.1093/bioinformatics/btv351.

384    30. Finn RD, Clements J, Eddy SR et al. HMMER Web Server: Interactive Sequence

385        Similarity Searching. Nucleic Acids Research. 2011; doi:10.1093/nar/gkr367.

386    31. Bushnell B. BBMap Short Read Aligner. Berkeley, CA: University of California; 2016.

387        https://sourceforge.net/projects/bbmap/ (August 7 2018, date last accessed).

388  32. Walker BJ, Abeel T, Shea T et al. Pilon: an integrated tool for comprehensive microbial

389      variant detection and genome assembly improvement. PLoS One. 2014;

390      doi:10.1371/journal.pone.0112963

391  33. Langmead B, Salzberg SL. Fast Gapped-Read Alignment with Bowtie 2. Nature

392      Methods. 2012; doi:10.1038/nmeth.1923.

393  34. Francis WR *haplotypeblastn.py;*

394      https://bitbucket.org/wrf/sequences/raw/f23b4dd3c965cc1774b9e10eb433242a18c13c65/

395      haplotypeblastn.py (August 7 2018, date last accessed).

396  35. Hahn C *select_contigs.pl;* https://github.com/chrishah/phylog/blob/master/scripts-

397      external/select_contigs.pl (August 7 2018, date last accessed).

398  36. Smit AFA, Hubley R, Green P. RepeatMasker; http://repeatmasker.org

399  37. Lunter G, Goodson M. Stampy: a statistical algorithm for sensitive and fast mapping of

400      Illumina sequence reads. Genome Research. 2011;21(6):936-939.

401  38. Pena-Centeno T; *filterBam,*

402      https://github.com/nextgenusfs/augustus/tree/master/auxprogs/filterBam

403  39. https://computationalbiologysite.wordpress.com/2013/07/25/incorporating-rnaseq-tophat-

404      to-augustus, (August 7 2018, date last accessed).

405  40. Stanke M, Steinkamp R, Waack S et al. AUGUSTUS: a web server for gene finding in

406      eukaryotes. Nucleic Acids Research 2004;32(suppl–2):W309–12.

407  41. Stanke M, Schöffmann O, Morgenstern B, Waack S. Gene prediction in eukaryotes with

408      a generalized hidden Markov model that uses hints from external sources. BMC

409      Bioinformatics. 2006; doi:10.1186/1471-2105-7-62.

42. Francis WR, Wörheide G. Similar ratios of introns to intergenic sequence across animal genomes. Genome Biology and Evolution; 2017;9(6):1582-1598.

43. Joint Genomics Institute: Trivial HTTP, Nemve1. https://genome.jgi.doe.gov/portal/Nemve1/Nemve1.download.ftp.html (7 August 2018, date last accessed)

44. Macdonald B. filterGenes.py. https://github.com/mcfaddenlab/filterGenes.py/blob/master/README.md (August 7, 2018, date last accessed)

45. Uniprot Consortium. UniProt: the Universal Protein Knowledgebase. Nucleic Acids Research. 2018; doi:10.1093/nar/gky092

46. UniProt Consortium, Reviewed Swiss-Prot data, ftp://ftp.uniprot.org/pub/databp/uniprot/current_release/knowledgebase/complete/uniprot _sprot.fasta.gz (August 7, 2018, date last accessed)

47. https://ndownloader.figshare.com/files/1215191, (August 7 2018, date last accessed).

48. Moran Y, Fredman D, Praher D et al. Cnidarian MicroRNAs frequently regulate targets by cleavage. Genome Research. 2014; doi:10.1101/gr.162503.113.

49. Joint Genome Institute. *Nematostella vectensis* genome. Version 1. https://genome.jgi.doe.gov/portal/Nemve1/Nemve1.download.html (August 7, 2018, date last accessed).

50. Putnam NH, Srivastava M, Hellsten U, Dirks B, Chapman J, Salamov, A, et al. Sea anemone genome reveals ancestral eumetazoan gene repertoire and genomic organization. Science 2007;317(5834):86-94.

432    51. Shinzato C, Shoguchi E, Kawashima T et al. National Center for Biotechnology

433        Information, *Acropora digitifera* genome Version 1.

434        https://www.ncbi.nlm.nih.gov/nuccore/BACK00000000.1 (November 2015, date last

435        accessed).

436    52. Shinzato C, Shoguchi E, Kawashima T, et al. Using the *Acropora digitifera* genome to

437        understand coral responses to environmental change. Nature. 2011;476:7360-320.

438    53. Liew YJ, Aranda M, Voolstra CR. Reefgenomics.Org - a repository for marine genomics

439        data. Database (Oxford) 2016, 1–4 A*mplexidiscus fenestrafer* and *Discosoma* sp.

440        genomes. http://corallimorpharia.reefgenomics.org (August 7, 2018, date last accessed).

441    54. Wang X, Liew YJ, Li Y, Zoccola D, Tambutte S, Aranda M. Draft genomes of the

442        corallimorpharians *Amplexidiscus fenestrafer* and *Discosoma* sp. Molecular Ecology

443        Resources 2017; 17(6); 187-195.

444    55. Baumgarten E, Simakov O, Esherick LY et al. National Center for Biotechnology

445        Information, (*Ex*)*aiptasia pallida* genome Version 1.1

446        ftp://ftp.ncbi.nlm.nih.gov/sra/wgs_aux/LJ/WW/LJWW01/LJWW01.1.fsa_nt.gz (August

447        7 2018, date last accessed).

448    56. Baumgarten S, Simakov O, Esherick LY et al. The genome of *Aiptasia*, a sea anemone

449        model for coral symbiosis. Proceedings of the National Academy of Sciences

450        2015;112(38):11893-11898.

451    57. Matz Lab. *Montastraea cavernosa* genome. Jul 2018 version.

452        https://matzlab.weebly.com/data--code.html (August 7, 2018, date last accessed).

453   58. Kayal E, Bentlage B, Pankey MS et al. Phylogenomics provides a robust topology of the

454        major cnidarian lineages and insights on the origins of key organismal traits. BMC

455        Evolutionary Biology 2018;18:68.

456   59. Kayal E, Roure B, Phillipe H et al. Cnidarian phylogenetic relationships as revealed by

457        mitogenomics. BMC Evolutionary Biology, 2013;13:5.

458   60. Jiang J, Quattrini AM, Francis WR, et al. A hybrid de novo assembly of the sea pansy

459        (*Renilla muelleri)* genome. GigaScience Database 2018. doi:XXXXX

460   61. Liew YJ, Aranda M, Voolstra CR. Reefgenomics.Org - a repository for marine genomics

461        data. Database (Oxford) 2016, 1–4 *Renilla muelleri* genome http://rmue.reefgenomics.org

462        (August 7, 2018, date last accessed)

463

464

465
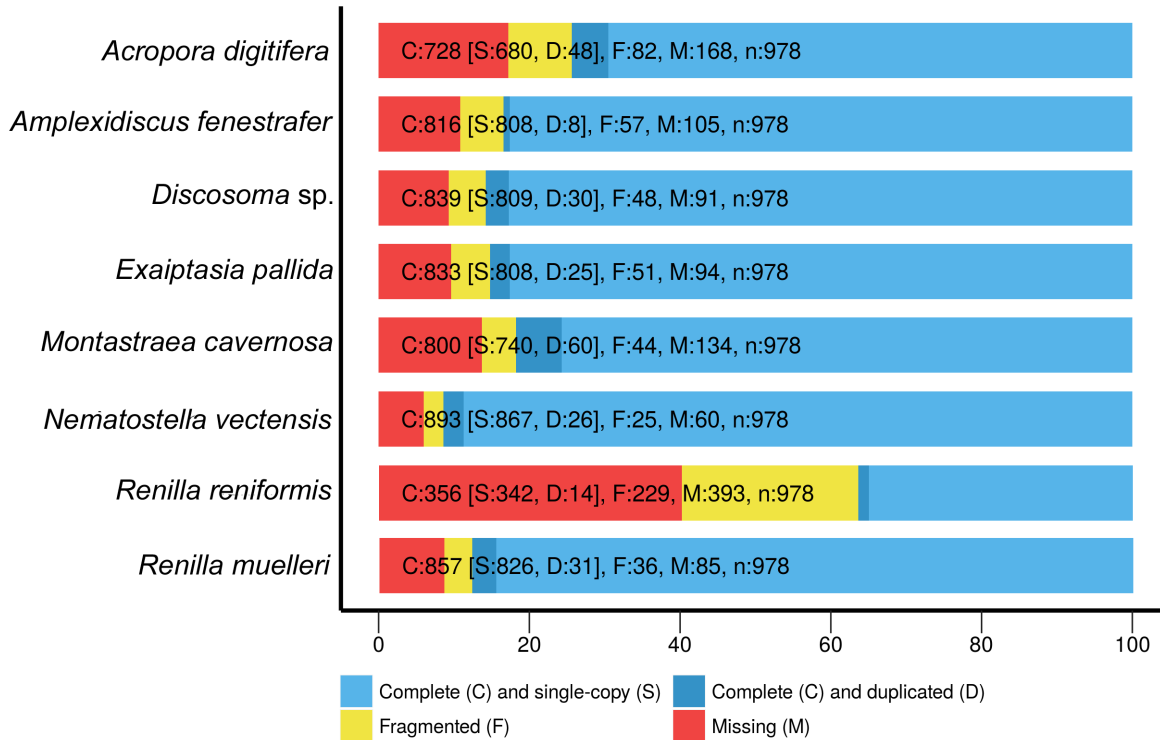
466

467

468

469

470

471

472

473

474

475 **Figure Captions**

476 **Figure 1.** BUSCO-generated chart showing relative completeness of six hexacoral genomes, one

477 octocoral genome, and the *Renilla muelleri* assembly.

478



479

480

481

482

483

484

485

486

487 **Table 1.** General statistics and BUSCO-completeness of both initial hybrid assemblies and the
488 final hybrid assembly.

| | MaSuRCA hybrid | SPAdes hybrid | Final MaSuRCA hybrid |
|---|---|---|---|
| scaffold total | 4,984 | 725,809 | 4,925 |
| contig total | 5,263 | 725,809 | 5,196 |
| scaffold sequence total | 172,512,580 | 231,255,108 | 172,160,214 |
| contig sequence total | 172.472 Mb | 231.255 Mb | 172.091 Mb |
| scaffold L/N50 | 635/70.423 Kb | 33702/1.007 Kb | 633/70.522 Kb |
| contig L/N50 | 687/64.492 Kb | 33702/1.007 Kb | 684/64.781 Kb |
| Max scaffold /contig length | 513.145 Kb | 323.009 Kb | 513.151 Kb |
| Number of scaffolds > 50 Kb | 960 | 14 | 961 |
| % genome in scaffolds > 50 Kb | 61.07% | 0.95% | 61.23% |
| GC% | 36.18% | 36.97% | 36.17% |
| N% | 0.042% | 0.000% | 0.040% |
| BUSCO assessment: | | | |
| Complete | 858 (87.73%) | 508 (51.94%) | 857 (87.63%) |
| Complete and single-copy | 826 (84.46%) | 493 (50.41%) | 826 (84.46%) |
| Complete and Duplicated | 32 (3.27%) | 15 (1.53%) | 31 (3.17%) |
| Fragmented | 36 (3.68%) | 200 (20.45%) | 36 (3.68%) |
| Missing | 84 (8.59%) | 270 (27.61%) | 85 (8.69%) |

489 Unmerged haplotypes were removed in the final assembly, which was also error-corrected with
490 Pilon.
491

492      **Table 2.** Statistics for the gene model predicted by Augustus**.**

|  | Number |
| --- | --- |
| Genes | 23,660 |
| Exons | 140,384 |
| Introns | 117,838 |
| Average Exon Length | 249 |
| Exons Per Gene | 5.93 |
| Average Intron Length | 524 |
| Introns Per Gene | 4.98 |
| BUSCO assessment: |  |
| Complete | 830 (84.87%) |
| Complete and single-copy | 798 (81.60%) |
| Complete and Duplicated | 32 (3.27%) |
| Fragmented | 64 (6.54%) |
| Missing | 84 (8.59%) |

493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513

514  **Supplemental Table S1.** Summary of Pilon changes per iteration
515

|  | First Iteration | Second Iteration | Third Iteration | Fourth Iteration | Fifth Iteration | Sixth Iteration |
|---|---|---|---|---|---|---|
| Single-nucleotide polymorphism changes | 32,292 | 10,039 | 4,688 | 2,790 | 1,697 | 1,152 |
| Ambiguous bp | 567 | 199 | 99 | 50 | 41 | 26 |
| Small Insertions | 9,180 (54,855 bp) | 1,982 (15,381 bp) | 1,231 (14,443 bp) | 858 (11,391 bp) | 810 (12,777 bp) | 641 (10,596 bp) |
| Small Deletions | 6706 (41,808 bp) | 1,925 (16,566 bp) | 1038 (11,922 bp) | 848 (12,916 bp) | 640 (10,603 bp) | 684 (12,319 bp) |

516
517

518

519

520

521

522

523

524

525

526

527

528

529

530 **Supplemental Table S2.** *Renilla muelleri* genome assembly and annotation comparisons to

531 other anthozoan genomes.

| | Genome Size (Mb) | Total # Complete BUSCOs** | Contig N50 (KB) | Exon length (bp) | # Predicted Gene models |
|---|---|---|---|---|---|
| *Acropora digitifera* | 420 | 728 | 10.6 | 230 | 23,668 |
| *Amplexidiscus fenestrafer* | 350 | 816 | 20.0 | 218 | 21,372 |
| *Discosoma* sp. | 428 | 839 | 18.7 | 226 | 23,199 |
| *Exaiptasia pallida* | 256 | 833 | 14.4 | NA | 26,042 |
| *Montastraea cavernosa* | 448 | 800 | 343 | NA | 30,360 |
| *Nematostella vectensis** | 329 | 893 | 19.8 | 208 | 27,273 |
| *Renilla reniformis* | 132 | 356 | 1.8 | NA | 12,689 |
| *Renilla muelleri* | 172 | 857 | 64.8 | 249 | 23,360 |

532
533 * Data taken from [52] and [56]
534 **Complete BUSCOs generated from analysis herein