

Finding Nemo 2.0: Hybrid assembly with Oxford Nanopore and Illumina reads dramatically improves the Clownfish (*Amphiprion ocellaris*) genome --Manuscript Draft--

Manuscript Number:	GIGA-D-17-00310
Full Title:	Finding Nemo 2.0: Hybrid assembly with Oxford Nanopore and Illumina reads dramatically improves the Clownfish (<i>Amphiprion ocellaris</i>) genome
Article Type:	Data Note
Funding Information:	
Abstract:	<p>Background: Some of the most widely recognised coral-reef fishes are clownfish or anemonefishes, members of the family Pomacentridae (subfamily: Amphiprioninae). They are popular aquarium species due to their bright colours, adaptability to captivity and fascinating behavior. Their breeding biology (sequential hermaphrodites) and symbiotic mutualism with sea anemones have attracted much scientific interest. Moreover there are some curious geographic-based phenotypes which warrant investigation. Leveraging on the advancement in Nanopore long read technology, we report the first hybrid assembly of the clown anemonefish (<i>Amphiprion ocellaris</i>) genome utilizing Illumina and Nanopore reads, further demonstrating the substantial impact of modest long read sequencing data sets on improving genome assembly statistics.</p> <p>Findings: We generated 43 Gb of short Illumina reads and 9 Gb of long Nanopore reads representing an approximate genome coverage of 54× and 11×, respectively, based on an estimated k-mer predicted genome size of 794 Mb. The final assembled genome size is contained in 6,404 scaffolds with an accumulated length of 880 Mb (96.3% BUSCO-calculated genome completeness). Compared to the Illumina-only assembly, the hybrid approach generated 93% less scaffolds with 18-fold increase in N50 length (401 kb) and increased the genome completeness by an additional 16%. A total of 27,240 high quality protein-coding genes were predicted from the clown anemonefish, 26,211 (96%) of which were annotated functionally with information from either sequence homology or protein signature searches.</p> <p>Conclusions: We present the first genome of any anemonefish and demonstrate the value of low coverage (~11×) long Nanopore reads sequencing in improving both genome contiguity and completeness. The near-complete genome of <i>A. ocellaris</i> will be an invaluable molecular resource for supporting a range of genetic, genomic and phylogenetic studies specifically for clownfish and more generally for other related fish species of the family Pomacentridae.</p>
Corresponding Author:	Han Ming Gan AUSTRALIA
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	
Corresponding Author's Secondary Institution:	
First Author:	Mun Hua Tan
First Author Secondary Information:	
Order of Authors:	Mun Hua Tan Christopher Austin Michael Hammer Yin Peng Lee

	Laurence Croft
	Han Ming Gan
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum</p>	Yes

[Standards Reporting Checklist?](#)

1 **Finding Nemo 2.0: Hybrid assembly with Oxford Nanopore and Illumina reads**
2 **dramatically improves the Clownfish (*Amphiprion ocellaris*) genome**

3
4
5 4 Mun Hua Tan^{1,2,3#}, Christopher M. Austin^{1,2,3#}, Michael P. Hammer⁴, Yin Peng Lee^{2,3},
6 Laurence J. Croft^{1,5}, Han Ming Gan^{1,2,3*}

7
8
9
10
11 7 ¹ Centre for Integrative Ecology, School of Life and Environmental Sciences, Deakin
12 University, Geelong, Victoria 3220, Australia

13
14 9 ² Genomics Facility, Tropical Medicine and Biology Platform, Monash University
15 Malaysia, Jalan Lagoon Selatan, Bandar Sunway 47500, Petaling Jaya, Selangor,
16
17
18 11 Malaysia

19
20 12 ³ School of Science, Monash University Malaysia, Jalan Lagoon Selatan, Bandar
21 Sunway 47500, Petaling Jaya, Selangor, Malaysia

22
23 14 ⁴ Museum and Art Gallery of the Northern Territory, Darwin 0801, Australia

24
25 15 ⁵ Malaysian Genomics Resource Centre Berhad, Boulevard Signature Office, Kuala
26 Lumpur, Malaysia

27
28
29
30
31 18 # Equal contribution

32
33
34
35
36 21 *** Corresponding author:**

37
38 22 Name: Han Ming Gan, PhD

39
40 23 Address: Building Ka, Level 4, Room 4.338, Centre for Integrative Ecology,
41 School of Life and Environmental Sciences, Deakin University, Waurn
42 Ponds, Victoria 3216, Australia

43
44 25
45 26 Phone: (+61) 490786277

46
47 27 Email: han.gan@deakin.edu.au, ORCID: 0000-0001-7987-738X

35

36 **Abstract**

37 **Background:** Some of the most widely recognised coral-reef fishes are clownfish or
38 anemonefishes, members of the family Pomacentridae (subfamily: Amphiprioninae).
39 They are popular aquarium species due to their bright colours, adaptability to
40 captivity and fascinating behavior. Their breeding biology (sequential
41 hermaphrodites) and symbiotic mutualism with sea anemones have attracted much
42 scientific interest. Moreover there are some curious geographic-based phenotypes
43 which warrant investigation. Leveraging on the advancement in Nanopore long read
44 technology, we report the first hybrid assembly of the clown anemonefish
45 (*Amphiprion ocellaris*) genome utilizing Illumina and Nanopore reads, further
46 demonstrating the substantial impact of modest long read sequencing data sets on
47 improving genome assembly statistics.

48

49 **Findings:** We generated 43 Gb of short Illumina reads and 9 Gb of long Nanopore
50 reads representing an approximate genome coverage of 54× and 11×, respectively,
51 based on an estimated k-mer predicted genome size of 794 Mb. The final assembled
52 genome size is contained in 6,404 scaffolds with an accumulated length of 880 Mb
53 (96.3% BUSCO-calculated genome completeness). Compared to the Illumina-only
54 assembly, the hybrid approach generated 93% less scaffolds with 18-fold increase in
55 N₅₀ length (401 kb) and increased the genome completeness by an additional 16%. A
56 total of 27,240 high quality protein-coding genes were predicted from the clown
57 anemonefish, 26,211 (96%) of which were annotated functionally with information
58 from either sequence homology or protein signature searches.

59

60 **Conclusions:** We present the first genome of any anemonefish and demonstrate the
61 value of low coverage (~11×) long Nanopore reads sequencing in improving both
62 genome contiguity and completeness. The near-complete genome of *A. ocellaris* will
63 be an invaluable molecular resource for supporting a range of genetic, genomic and
64 phylogenetic studies specifically for clownfish and more generally for other related
65 fish species of the family Pomacentridae.

66

67 **Keywords:** clownfish, long reads, genome, transcriptome, hybrid assembly

68

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

69 **Data description**

1
2 70 The clown anemonefish, *Amphiprion ocellaris* (NCBI Taxon ID: 80972, Fish Base
3
4 71 ID:6509), is a well-known tropical marine fish species among the non-scientific
5
6 72 community especially following the Pixar film “Finding Nemo” and its sequel
7
8 73 “Finding Dory” [1]. The visual appeal of *A. ocellaris* due to its bright coloration and
9
10 74 behaviour and ease of husbandry has maintained a strong global demand for this
11
12 75 species in the marine aquarium trade driving a fine balance between positive
13
14 76 environmental awareness versus sustainable ornamental use [1, 2]. Further, given high
15
16 77 survival rates and ability to complete their lifecycle in captivity, captive-breeding
17
18 78 programs to partially sustain their global trade have been successful [3]. For the
19
20 79 scientific community, *A. ocellaris* or anemonefishes in general, are actively studied
21
22 80 due to their intriguing reproductive strategy i.e. sequential hermaphroditism [4-7] and
23
24 81 mutualistic relationships with sea anemones [8-12]. Phenotypic body-colour variation
25
26 82 based on host-anemone use and geography also pose additional questions regarding
27
28 83 adaptive genetic variation [13]

29
30 84 In recent years, concurrent with the advent of long read sequencing
31
32 85 technologies [14], several studies have explored combining short but accurate
33
34 86 Illumina reads with long but less accurate Nanopore/PacBio reads to obtain genome
35
36 87 assemblies that are usually more contiguous with higher completeness than
37
38 88 assemblies based on Illumina-only reads [15-19]. To further contribute to the
39
40 89 evaluation of long read technology in fish genomics [15], we sequenced the whole
41
42 90 genome of *A. ocellaris* using Oxford Nanopore and Illumina technologies and
43
44 91 demonstrate that hybrid assembly of long and short reads greatly improved the quality
45
46 92 of genome assembly.

47 **Whole genome sequencing**

48
49 95 Tissues for genome assembly and as reference material were sourced from the
50
51 96 collection of the Museum and Art Gallery of the Northern Territory (NTM). The
52
53 97 samples used for DNA extraction and subsequent whole genome sequencing were
54
55 98 from freshly-vouchered captive bred *A. ocellaris* specimens representing a unique
56
57 99 black and white colour phenotype found only in the Darwin Harbour region, Australia
100 (NTM A3764, A4496, A4497).

58
59 101 Genomic DNA was extracted from multiple fin clip and muscle samples using
60
61 102 E.Z.N.A.® Tissue DNA Kit (Omega Bio-tek, Norcross, GA). For Illumina library
62
63
64
65

103 prep, approximately 1 µg of gDNA from isolate A3764 was sheared to 300 bp using a
104 Covaris Focused-ultrasonicator (Covaris, Woburn, MA) and subsequently processed
105 using TruSeq DNA sample prep kit (Illumina, San Diego, CA) according to the
106 manufacturer's instructions. Paired-end sequencing was performed on a single lane of
107 HiSeq 2000 (Illumina, San Diego, CA) located at the Malaysian Genomics Resource
108 Centre Berhad. Two additional libraries were constructed from specimen NTM
109 A3764 and both libraries were sequenced on the MiSeq (2×300 bp setting) located at
110 the Monash University Malaysia Genomics Facility.

111 To generate Oxford Nanopore long reads, approximately 5 µg of gDNA was
112 extracted from isolates NTM A4496 and A4497, size-selected (8 – 30 kb) with a
113 BluePippin (Sage Science, Beverly, MA) and processed using the Ligation
114 sequencing 1D kit (Oxford Nanopore, UK) according to the manufacturer's
115 instructions. Three libraries were prepared and sequenced on three different R9.4
116 flowcells using the MinION portable DNA sequencer (Oxford Nanopore, UK) for 48
117 hours.

118

119 **Sequence read processing**

120 Raw Illumina short reads were adapter-trimmed with Trimmomatic v.0.36
121 (*ILLUMINACLIP:2:30:10, MINLEN:100*) [20] followed by a screening for vectors
122 and contaminants, using Kraken v.0.10.5 [21] based on the MiniKraken DB. Kraken-
123 unclassified reads i.e. non-microbial/viral origin were aligned to the complete
124 mitogenome of NTM A3764 (See "Mitogenome Assembly") to exclude sequences of
125 organellar origin. This results in a total of 42.35 Gb "clean" short reads. Nanopore
126 reads were base-called from their raw FAST5 files using the Oxford Nanopore
127 propriety base-caller, Albacore version 2.0.1. Applying a minimum and maximum
128 length cut-off of 500 bp and 1Mbp, respectively, this study produced a total of 8.95
129 Gbp in 895,672 Nanopore reads (N₅₀: 12.7 kb)

130

131 **Genome size estimation**

132 K-mer counting with the filtered Illumina reads was performed with Jellyfish v.2.2.6
133 [22], generating k-mer frequency distributions of 17-, 21- and 25-mers. These
134 histograms were processed by GenomeScope [23] that estimated a genome size of
135 791 to 794 Mb with approximately 80% of unique content and a heterozygosity level
136 of 0.6% (Supplemental Figure 1). Given that we had previously excluded adapters as

137 well as sequences from contaminant or organellar sources, the max kmer coverage
138 filter was not applied (*max kmer coverage: -1*). The genome size estimated from this
139 approach is within the range of sizes listed for other *Amphiprion* species (792 Mb -
140 1.2 Gb) as reported on the Animal Genome Size Database
141 (<http://www.genomesize.com> accessed on 11th November 2017)

143 **Hybrid genome assembly**

144 Both short-read-only and hybrid *de novo* assemblies were performed with MaSuRCA
145 v.3.2.2 [24]. During hybrid assembly, errors were encountered in the fragment
146 correction step of the Celera Assembler (CA). To overcome this, given that the CA
147 assembler is no longer maintained, we disabled the *frgcorr* step based on one of the
148 developer's recommendations and the hybrid assembly was subsequently improved
149 with 10 iterations of Pilon v.1.22 [25], using short reads to correct bases, fix mis-
150 assemblies and fill assembly gaps. To assess the completeness of the genome,
151 BUSCO v.3.0.2 [26] was used to locate the presence or absence of the Actinopterygii-
152 specific set of 4,584 single copy orthologs (OrthoDB v9).

153 The short-read-only and hybrid assemblies yielded total assembly sizes of 851
154 Mb and 880Mb, respectively. Inclusion of Nanopore long reads for hybrid assembly
155 representing approximately 11× genome coverage led to a 94% decrease in the
156 number of scaffolds (> 500 bp) from 106,526 to 6,404 scaffolds and an 18-fold
157 increase in the scaffold N₅₀ length from 21,802 bp to 401,715 bp (Table 1). In
158 addition, the genome completeness was also substantially improved in the hybrid
159 assembly, with BUSCO detecting complete sequences of 96.3% (4,417/4,584) of
160 single-copy orthologs in the Actinopterygii-specific dataset.

162 **Transcriptome sequencing and assembly**

163 Total RNA extraction from RNAsShield-preserved whole body and muscle tissues of
164 isolate A4496 used Quick-RNA MicroPrep (Zymo Research Corpt, Irvine, CA)
165 according to the manufacturer's protocols. After assessing total RNA intactness on
166 the TapeStation2100 (Agilent), mRNA was enriched using NEBNext Poly(A) mRNA
167 magnetic isolation kit (NEB, Ipswich, MA) and processed with NEBNext Ultra RNA
168 library prep kit for Illumina (NEB, Ipswich, MA). Libraries from both whole body and
169 muscle tissues were sequenced on a fraction of MiSeq V3 flowcell (1×150 bp).
170 Single-end reads from both libraries in addition to two publicly available *A. ocellaris*

171 transcriptome sequencing data (SRR5253145 and SRR5253146, Bioproject ID:
172 PRJNA374650) were individually assembled using Scallop v0.10.2 [27] based on
173 HiSat2 [28] alignment of RNA-sequencing reads to the newly generated *A. ocellaris*
174 genome. The transcriptome assemblies were subsequently merged using the tr2aacds
175 pipeline from the EvidentialGene [29] package and similarly assessed for
176 completeness using BUSCO version 3 [26]. The final non-redundant transcriptome
177 assembly, which was subsequently used to annotate the *A. ocellaris* genome, contains
178 25,264 contigs/isotigs (putative transcripts) with an accumulated length of 68.4 Mb
179 and BUSCO-calculated completeness of 92.8% (Table 1).

180

181 **Genome annotation**

182 Protein-coding genes were predicted with the MAKER v.2.31.9 genome annotation
183 pipeline [30]. A total of three passes were run with MAKER2; the first pass was
184 based on hints from the assembled transcripts as RNA-seq evidence (*est2genome*) and
185 protein sequences from 11 fish species downloaded from Ensembl [31]
186 (*protein2genome*), whereas the second and third passes included gene models trained
187 from the first (and then, second) passes with *ab initio* gene predictors SNAP [32] and
188 Augustus [33]. In the final set of genes predicted, sequences with Annotation Edit
189 Distance (AED) values less than 0.5 were retained. A small AED value suggests a
190 lesser degree of difference between the predicted protein and the evidences used in
191 the prediction (i.e. fish proteins, transcripts). This resulted in a final set of 27,240
192 protein-coding genes with an average AED of 0.14 (Table 1). A BUSCO analysis on
193 the completeness of the predicted protein dataset detected the presence of 4,259
194 (92.9%) single-copy orthologs from the Actinopterygii-specific dataset.

195 Further, to infer putative function of these predicted proteins, NCBI's *blastp*
196 v.2.6.0 (*-evalue 1e-10, -seg yes, -soft_masking true, -lcase_masking*) [34] was used to
197 find homology to existing vertebrate sequences in the non-redundant (NR) database.
198 Applying a hit fraction filter to include only hits with ≥ 70 % target length fraction, the
199 remaining un-annotated sequences were subsequently aligned to all sequences in the
200 NR database. With this method, 20,107 proteins (74%) were annotated with a putative
201 function based on homology. Additionally, InterProScan v.5.26.65 [35] was used to
202 examine protein domains, signatures and motifs present in the predicted protein
203 sequences. This analysis detected domains, signatures or motifs for 26,211 proteins
204 (96%). Overall, 96% of the predicted clown fish protein-coding genes were

205 functionally annotated with information from at least one of the two approaches.

206

207 **Mitogenome recovery via genome skimming**

208 Genome skimming [36, 37] was performed on three additional *A. ocellaris* individuals
209 from known locality (Supplemental Table 1). Mitogenome assembly was performed
210 with MITObim version 1.9 [38] using the complete mitogenome of *A. ocellaris*
211 (GenBank: NC009065.1) as the bait for read mapping. The assembled mitogenomes
212 were subsequently annotated with MitoAnnotator [39]. Consistent with original
213 broodstock collection from northern Australia, the captive bred black and white *A.*
214 *ocellaris* NTM A3764 exhibits strikingly high whole mitogenome nucleotide identity
215 (99.98%) to sample NTM A3708 as a wild collection from Darwin Harbour,
216 Australia. In addition, the overall high pair-wise nucleotide identity (> 98%) of NTM
217 A3764 to newly generated and publicly available *A. ocellaris* whole mitogenomes
218 further supports its morphological identification as *A. ocellaris* (Supplemental Table
219 1).

220

221 **Identification of the *cyp19a1a* gene associated with sexual differentiation**

222 The validated *cyp19a1a* enzyme of *Danio rerio* (Uniprot: O42145) was used as the
223 query (E-value = 1e-10) for blastp search against the predicted *A. ocellaris* proteins.
224 The top blast hit, AMPOCE_00012675-RA (71.5% protein identity to O42145), was
225 searched (tblastn) against the NCBI TSA database (Taxon: *Amphirion*) and showed
226 strikingly high protein identity (99%) to a translated RNA transcript from *Amphirion*
227 *bicinctus* (c183337_g1_i2 : GDCV01327693) [5]. The *cyp19a1a* gene codes for a
228 steroidogenic enzyme that converts androgens into estrogens [40] and was recently
229 shown to be instrumental during sex change in *Amphirion bicinctus* as evidenced by
230 significant correlation and differential expression of this gene between male and
231 mature females [5]. We also observed similar profile based on mapping of RNA reads
232 from the publicly available male and female transcriptomes of *A. ocellaris* to the
233 *cyp19a1a* gene region as visualized using Integrative Genomics Viewer [41] (Figure
234 2). The *A. ocellaris cyp19a1a* gene is located on a 419 kb scaffold and is spanned by
235 multiple Minimap2-aligned Nanopore reads [42]. It is noteworthy that in the Illumina-
236 only assembly, this gene is fragmented and located on 3 relatively short scaffolds
237 (Figure 2).

238

239 **Conclusion**

240 We present the first clownfish genome and demonstrate the value of low coverage
241 (~11×) Nanopore long read sequencing in improving both genome contiguity and
242 completeness. The near-complete genome of *A. ocellaris* will be an invaluable
243 molecular resource for supporting a range of genetic, genomic, and phylogenetic
244 studies specifically for clownfish and more generally for other related fish species of
245 the family Pomacentridae.

246

247

248 **Availability of supporting data**

249 Data supporting the results of this article will be available in the GigaDB repository.
250 Raw Illumina and Nanopore reads generated in this study are available in the
251 Sequence Read Archive (SRP123679) whereas Whole Genome Shotgun project has
252 been deposited at DDBJ/EMBL/GenBank under the accession NXFZ00000000, both
253 under BioProject PRJNA407816.

254

255 **Acknowledgements**

256 This study was funded by Monash University Malaysia Tropical and Biology
257 Multidisciplinary Platform.

258

259 **Competing interests**

260 The authors declare that they have no competing interests

261

262

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

263 **Figure Legends**

1
2 264

3
4 265 Figure 1. The clown anemonefish (*Amphiprion ocellaris*): Photo by Michael P.

5
6 266 Hammer

7
8 267

9 268 Figure 2. Mapping of MinION long reads, Illumina-assembled scaffolds and RNA-

10
11 269 sequencing reads of male and female *A. ocellaris* to genomic region containing the

12
13 270 *cyp19a1a* gene. Transcripts per million (TPM) values were calculated using Kallisto

14
15 271 version 0.43.1 [43].

16
17 272

18 273 Supplemental Figure 1. Genome profiling of *A. ocellaris* based on Illumina short

19
20 274 reads

21
22 275

23
24 276

25
26 277

27
28 278

29
30 279

31
32 280

33
34 281

35
36 282

37
38 283

39
40 284

41
42 285

43
44 286

45
46 287

47
48 288

49
50 289

51
52 290

53
54 291

55
56 292

57
58 293

59
60 294

61
62 295

63
64 296

65

297 **References**

298

- 299 1. Miltz TA and Foale S. The “Nemo Effect”: Perception and reality of Finding
300 Nemo's impact on marine aquarium fisheries. *Fish and Fisheries*. 2017;18
301 3:596-606.
- 302 2. Madduppa HH, von Juterzenka K, Syakir M and Kochzius M. Socio-economy
303 of marine ornamental fishery and its impact on the population structure of the
304 clown anemonefish *Amphiprion ocellaris* and its host anemones in Spermonde
305 Archipelago, Indonesia. *Ocean & Coastal Management*. 2014;100 Supplement
306 C:41-50. doi:<https://doi.org/10.1016/j.ocecoaman.2014.07.013>.
- 307 3. Hall H and Warmolts D. The Role of Public Aquariums in the Conservation
308 and Sustainability of the Marine Ornamentals Trade. *Marine Ornamental*
309 *Species*. Blackwell Publishing Company; 2008. p. 305-24.
- 310 4. Madhu R, Madhu K and Rethesh T. Life history pathways in false clown
311 *Amphiprion ocellaris* Cuvier, 1830: A journey from egg to adult under captive
312 condition. *Journal of the Marine Biological Association of India*. 2012;54
313 1:77-90.
- 314 5. Casas L, Saborido-Rey F, Ryu T, Michell C, Ravasi T and Irigoien X. Sex
315 Change in Clownfish: Molecular Insights from Transcriptome Analysis.
316 *Scientific Reports*. 2016;6:35461. doi:10.1038/srep35461.
- 317 6. Buston P. Social hierarchies: size and growth modification in clownfish.
318 *Nature*. 2003;424 6945:145-6.
- 319 7. Kobayashi Y, Horiguchi R, Miura S and Nakamura M. Sex- and tissue-
320 specific expression of P450 aromatase (*cyp19a1a*) in the yellowtail clownfish,
321 *Amphiprion clarkii*. *Comp Biochem Physiol A Mol Integr Physiol*. 2010;155
322 2:237-44.
- 323 8. Davenport D and Norris KS. Observations on the symbiosis of the sea
324 anemone *Stoichactis* and the pomacentrid fish, *Amphiprion percula*. *The*
325 *Biological Bulletin*. 1958;115 3:397-410.
- 326 9. Arvedlund M and Nielsen LE. Do the anemonefish *Amphiprion ocellaris*
327 (*Pisces: Pomacentridae*) imprint themselves to their host sea anemone
328 *Heteractis magnifica* (*Anthozoa: Actinidae*)? *Ethology*. 1996;102 2:197-211.
- 329 10. Mariscal RN. An experimental analysis of the protection of *Amphiprion*
330 *xanthurus* Cuvier & Valenciennes and some other anemone fishes from sea
331 anemones. *Journal of Experimental Marine Biology and Ecology*. 1970;4
332 2:134-49.
- 333 11. Hattori A. Coexistence of two anemonefishes, *Amphiprion clarkii* and *A.*
334 *perideraion*, which utilize the same host sea anemone. *Environmental Biology*
335 *of Fishes*. 1995;42 4:345-53.
- 336 12. Schmiege PF, D'Aloia CC and Buston PM. Anemonefish personalities
337 influence the strength of mutualistic interactions with host sea anemones.
338 *Marine Biology*. 2017;164 1:24.
- 339 13. Allen GR. *Damselfishes of the world*. Melle, Germany: Mergus Publishers;
340 1991.
- 341 14. Heather JM and Chain B. The sequence of sequencers: The history of
342 sequencing DNA. *Genomics*. 2016;107 1:1-8.
343 doi:<https://doi.org/10.1016/j.ygeno.2015.11.003>.
- 344 15. Austin CM, Tan MH, Harrisson KA, Lee YP, Croft LJ, Sunnucks P, et al. De
345 novo genome assembly and annotation of Australia's largest freshwater fish,

- 346 the Murray cod (*Maccullochella peelii*), from Illumina and Nanopore
347 sequencing read. *GigaScience*. 2017;6 8:1-6. doi:10.1093/gigascience/gix063.
- 348 16. Gan HM, Lee YP and Austin CM. Nanopore Long-Read Guided Complete
349 Genome Assembly of *Hydrogenophaga intermedia*, and Genomic Insights into
350 4-Aminobenzenesulfonate, p-Aminobenzoic Acid and Hydrogen Metabolism
351 in the Genus *Hydrogenophaga*. *Frontiers in Microbiology*. 2017;8 1880
352 doi:10.3389/fmicb.2017.01880.
- 353 17. Zimin AV, Puiu D, Hall R, Kingan S, Clavijo BJ and Salzberg SL. The first
354 near-complete assembly of the hexaploid bread wheat genome, *Triticum*
355 *aestivum*. *GigaScience*. 2017:gix097-gix. doi:10.1093/gigascience/gix097.
- 356 18. Zimin AV, Stevens KA, Crepeau MW, Puiu D, Wegrzyn JL, Yorke JA, et al.
357 An improved assembly of the loblolly pine mega-genome using long-read
358 single-molecule sequencing. *GigaScience*. 2017;6 1:1-4.
359 doi:10.1093/gigascience/giw016.
- 360 19. Zimin AV, Puiu D, Luo M-C, Zhu T, Koren S, Marcais G, et al. Hybrid
361 assembly of the large and highly repetitive genome of *Aegilops tauschii*, a
362 progenitor of bread wheat, with the mega-reads algorithm. *Genome Research*.
363 2017; doi:10.1101/gr.213405.116.
- 364 20. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for
365 Illumina sequence data. *Bioinformatics*. 2014;30 15:2114-20.
- 366 21. Wood DE and Salzberg SL. Kraken: ultrafast metagenomic sequence
367 classification using exact alignments. *Genome biology*. 2014;15 3:R46.
- 368 22. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel
369 counting of occurrences of k-mers. *Bioinformatics*. 2011;27 6:764-70.
370 doi:10.1093/bioinformatics/btr011.
- 371 23. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski
372 J, et al. GenomeScope: fast reference-free genome profiling from short reads.
373 *Bioinformatics*. 2017;33 14:2202-4. doi:10.1093/bioinformatics/btx153.
- 374 24. Zimin AV, Marçais G, Puiu D, Roberts M, Salzberg SL and Yorke JA. The
375 MaSuRCA genome assembler. *Bioinformatics*. 2013;29 21:2669-77.
- 376 25. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al.
377 Pilon: an integrated tool for comprehensive microbial variant detection and
378 genome assembly improvement. *PloS one*. 2014;9 11:e112963.
- 379 26. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM.
380 BUSCO: assessing genome assembly and annotation completeness with
381 single-copy orthologs. *Bioinformatics*. 2015;31 19:3210-2.
- 382 27. Shao M and Kingsford C. Accurate assembly of transcripts through phase-
383 preserving graph decomposition. *Nature Biotechnology*. 2017;
384 doi:10.1038/nbt.4020
385 <https://www.nature.com/articles/nbt.4020#supplementary-information>.
- 386 28. Kim D, Langmead B and Salzberg SL. HISAT: a fast spliced aligner with low
387 memory requirements. *Nature Methods*. 2015;12:357.
388 doi:10.1038/nmeth.3317
389 <https://www.nature.com/articles/nmeth.3317#supplementary-information>.
- 390 29. Gilber D. Gene-omes built from mRNA seq not genome DNA.
391 *F1000Research*. 2016;5 1695:1. doi:10.7490/f1000research.1112594.1.
- 392 30. Holt C and Yandell M. MAKER2: an annotation pipeline and genome-
393 database management tool for second-generation genome projects. *BMC*
394 *bioinformatics*. 2011;12 1:491.

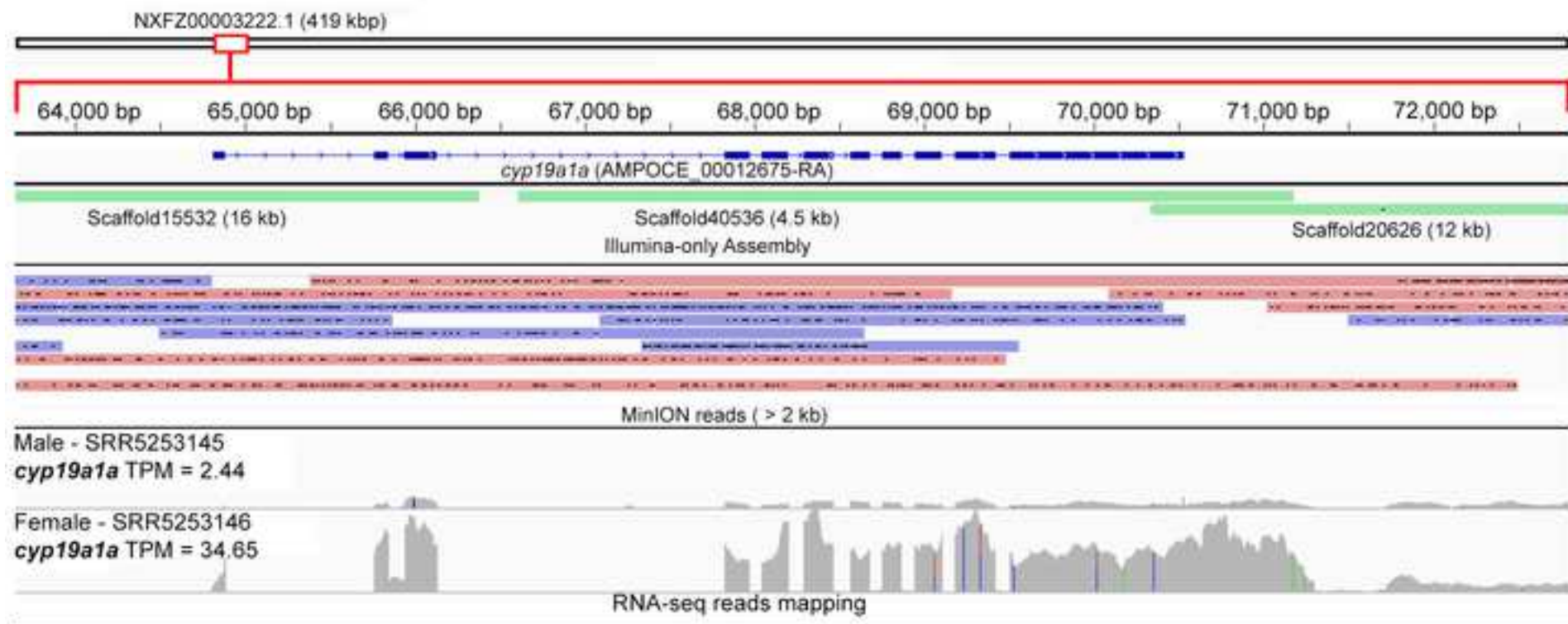
- 395 31. Hubbard T, Barker D, Birney E, Cameron G, Chen Y, Clark L, et al. The
396 Ensembl genome database project. *Nucleic acids research*. 2002;30 1:38-41.
- 397 32. Korf I. Gene finding in novel genomes. *BMC bioinformatics*. 2004;5 1:59.
- 398 33. Stanke M, Steinkamp R, Waack S and Morgenstern B. AUGUSTUS: a web
399 server for gene finding in eukaryotes. *Nucleic acids research*. 2004;32
400 suppl_2:W309-W12.
- 401 34. Boratyn GM, Camacho C, Cooper PS, Coulouris G, Fong A, Ma N, et al.
402 BLAST: a more efficient report with usability improvements. *Nucleic acids
403 research*. 2013;41 W1:W29-W33.
- 404 35. Zdobnov EM and Apweiler R. InterProScan—an integration platform for the
405 signature-recognition methods in InterPro. *Bioinformatics*. 2001;17 9:847-8.
- 406 36. Gan HM, Schultz MB and Austin CM. Integrated shotgun sequencing and
407 bioinformatics pipeline allows ultra-fast mitogenome recovery and confirms
408 substantial gene rearrangements in Australian freshwater crayfishes. *BMC
409 evolutionary biology*. 2014;14 1:19.
- 410 37. Grandjean F, Tan MH, Gan HM, Lee YP, Kawai T, Distefano RJ, et al. Rapid
411 recovery of nuclear and mitochondrial genes by genome skimming from
412 Northern Hemisphere freshwater crayfish. *Zoologica Scripta*. 2017.
- 413 38. Hahn C, Bachmann L and Chevreur B. Reconstructing mitochondrial
414 genomes directly from genomic next-generation sequencing reads—a baiting
415 and iterative mapping approach. *Nucleic acids research*. 2013;41 13:e129-e.
- 416 39. Iwasaki W, Fukunaga T, Isagozawa R, Yamada K, Maeda Y, Satoh TP, et al.
417 MitoFish and MitoAnnotator: A mitochondrial genome database of fish with
418 an accurate and automatic annotation pipeline. *Molecular biology and
419 evolution*. 2013;30 11:2531-40.
- 420 40. Kallivretaki E, Eggen R, Neuhauss S, Alberti M, Kausch U and Segner H.
421 Aromatase in zebrafish: a potential target for endocrine disrupting chemicals.
422 *Mar Environ Res*. 2006;62 90:7.
- 423 41. Thorvaldsdóttir H, Robinson JT and Mesirov JP. Integrative Genomics Viewer
424 (IGV): high-performance genomics data visualization and exploration.
425 *Briefings in bioinformatics*. 2013;14 2:178-92.
- 426 42. Li H. Minimap2: fast pairwise alignment for long nucleotide sequences. *arXiv*.
427 2017;1708:Artciel 01492.
- 428 43. Bray NL, Pimentel H, Melsted P and Pachter L. Near-optimal probabilistic
429 RNA-seq quantification. *Nat Biotech*. 2016;34 5:525-7. doi:10.1038/nbt.3519
430 [http://www.nature.com/nbt/journal/v34/n5/abs/nbt.3519.html#supplementary-](http://www.nature.com/nbt/journal/v34/n5/abs/nbt.3519.html#supplementary-information)
431 [information](http://www.nature.com/nbt/journal/v34/n5/abs/nbt.3519.html#supplementary-information).
432
433

Table 1. Genome and Transcriptome statistics of the clownfish (*Amphiprion ocellaris*) genome

	Illumina (≥500bp)	Illumina + Nanopore (≥500bp)
Genome Assembly		
<u>Contig statistics</u>		
Number of contigs	133,997	7,810
Total contig size (bp)	851,389,851	880,159,068
Contig N ₅₀ size (bp)	15,458	323,678
Longest contig (bp)	204,209	2,051,878
<u>Scaffold statistics</u>		
Number of scaffolds	106,526	6,404
Total scaffold size (bp)	852,602,726	880,704,246
Scaffold N ₅₀ size (bp)	21,802	401,715
Longest scaffold (bp)	227,111	3,111,502
GC / AT / N (%)	39.6 / 60.2 / 0.14	39.4 / 60.5 / 0.06
<u>BUSCO Genome Completeness</u>		
Complete	3,691 (80.5%)	4,417 (96.3%)
Complete and single copy	3,600 (78.5%)	4,269 (93.1%)
Complete and duplicated	91 (2.0%)	148 (3.2%)
Fragmented	534 (11.6%)	63 (1.4%)
Missing	359 (7.9%)	104 (2.3%)
Transcriptome Assembly		
Number of contigs	25,364	
Total length (bp)	68,405,796	
Contig N ₅₀ size (bp)	3,670	
<u>BUSCO completeness</u>		
Complete	4,253 (92.8%)	
Complete and single-copy	4,128 (90.1%)	
Complete and duplicated	125 (2.7%)	
Fragmented	127 (2.8%)	
Missing	204 (4.4%)	
Genome Annotation		
Number of protein-coding genes	27,420	
Number of functionally-annotated proteins	26,211	

Mean protein length	514 aa
Longest protein	29,084 aa (titin protein)
Average number (length) of exon per gene	9 (355 bp)
Average number (length) of intron per gene	8 (1,532 bp)







Click here to access/download
Supplementary Material
SupplementalFigure1.tif





Click here to access/download
Supplementary Material
Supplemental Table 1_151117.docx

