

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

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00310R1	
Full Title:	Finding Nemo: Hybrid assembly with Oxford Nanopore and Illumina reads greatly improves the Clownfish (<i>Amphiprion ocellaris</i>) genome assembly	
Article Type:	Data Note	
Funding Information:	Monash University Malaysia (Tropical and Biology Multidisciplinary Platform)	Not applicable
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 the range of estimated k-mer-predicted genome sizes of between 791 to 967 Mbp. 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 94% fewer 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 assembly contiguity and completeness. The near-complete assembly of the <i>A. ocellaris</i> genome 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:	
Response to Reviewers:	<p>Reviewer #1</p> <p>Comment 1: In general, given the prominence (with 'dramatically') of nanopore data in the title, I would like to encourage the authors to elaborate on this aspect of the study in the Conclusion section. For example, why did you use this particular strategy (MaSuRCA assembler), and what are its strengths and weaknesses? How does long-read coverage affect the assembly process (this study uses only three nanopore flowcells - would this be a recommended efficient strategy to 'fix' any Illumina-based assembly)? How far are we from non-hybrid nanopore-based assemblies?</p> <p>Response: We've added discussion on this matter in the Conclusion. We chose MaSuRCA due to its demonstrated accuracy in addition to its ability to better utilize Nanopore reads in its initial assembly step. Despite producing a highly contiguous genome assembly of the clown fish, MaSuRCA hybrid assembly is an extremely computationally expensive and time-consuming. We are hopeful that the cost of Nanopore sequencing will reduce further especially with the recent release of PromethION which may revolutionize routine long read sequencing. However, as of now, high coverage long read sequencing is still financial challenging for smaller research group despite its indubitable value in assembling challenging (repetitive and/or heterozygous) genome. "Hybrid assembly of Illumina and Nanopore reads is one of the new features of the MaSuRCA assembler version 3.2.2 that works by constructing long and accurate mega-reads from the combination of long and short read data. Although this is a relatively computationally-intensive strategy with long run-times, we observed substantial improvement in the genome statistics when compared to Illumina-only assembly. As Nanopore long technology becomes more mature, it is likely that future de novo genome assembly will shift towards high coverage long read-only assembly followed by genome polishing using Illumina reads."</p> <p>-----</p> <p>Comment 2: Finally, a very similar genome project manuscript was recently posted on BioRxiv: Anna Marcionetti et al., First draft genome assembly of an iconic clownfish species (<i>Amphiprion frenatus</i>), doi 10.1101/205443, 18 October 2017 The manuscripts do not cite each other, but arrive at similar genome assembly qualities using similar strategies.</p> <p>Response: Although the manuscript was posted on BioRxiv, the genome assembly itself is currently unavailable to the public and, according to the manuscript, will only be available in DRYAD repository (under 'Data Accessibility'). So far, we were not able to locate the data on DRYAD.</p> <p>-----</p> <p>Comment 3: 1. Line 128: an upper limit of 1 Mbp reads probably did not exclude anything. What was the actual longest read length?</p> <p>Response: We had initially obtained reads up to 40Mbp base-called from an older version of Albacore. But after re-analysis with version 2.0.1 for this assembly, we had gotten far more sensible Nanopore read lengths and therefore did not apply any maximum length cutoff. We have removed the details on upper limit in the revised manuscript. The actual longest mappable read length for our Nanopore read dataset was 101,379bp, aligned to scaffold6249 (276,565bp) and contains genes with IDs: AMPOCE_00020294 and AMPOCE_00020295.</p> <p>-----</p> <p>Comment 4: 2. Line 131.../Supplemental Figure 1. Not all Illumina data were apparently used for the k-mer profile. Does this perhaps explain the considerable difference in estimated</p>

genome size and assembled genome size? If not, is there another explanation? Also, the legend to the figure ('genome profiling') could be more informative (e.g. genome size estimate...)

Response:

Based on this reviewer comment, we re-ran GenomeScope using the k-mer profile from all Illumina reads, which estimated genome sizes of 806 to 812 Mbp with different k-mer sizes, not too different from the initial number. A separate independent analysis was performed with BMap, which estimated a haploid genome size of 967Mb. Given this result, the assembled genome size is well within the range of genome size estimated based on different methods. We have added the results from the BMap analysis in the manuscript at lines 139 to 140. Supplemental Figure 1 has also been improved.

Comment 5:

3. A (supplementary) table with sequencing statistics (yield for each type of data, incl. RNA-seq) would be appropriate.

Response:

Sequencing statistics has been summarized for each sample ID in the new Supplemental Table 1.

Reviewer #2

Comment 1:

Title:

The reference to Nemo 2.0 and the phrase "dramatically improves" led me to believe that this was the second version of an already existing genome assembly. However, I could not find any other Amphiprion ocellaris assemblies by googling, besides a bioRxiv preprint of Amphiprion frenatus (<https://www.biorxiv.org/content/early/2017/10/18/205443>). I am not sure if this warrants changing the title, but please be aware of it. Also, I dislike using "genome" to refer to the genome assembly. I don't think you actually improve the genome present in the species.

Response:

The "2.0" in the title was included due to this manuscript (10.1016/j.gene.2006.03.028) with a similar title "Finding Nemo" and was not due to the recent A. frenatus genome in bioRxiv (see also Response to Reviewer's 1 Comment 2). However, given that our study is not a follow-up of the phylogenetic study reported by Santini and Polacco (2006), we agree that this can be confusing to readers and have removed "2.0" from the title. We have also added "assembly" into the title as per reviewer's suggestion and slightly modified the title.

Comment 2:

Abstract:

Line 54: "93 % less scaffolds". This should be "fewer" if I'm not mistaken.

Response:

Replaced "less" with "fewer".

Comment 3:

Lines 60-65: I prefer to see "genome assembly" instead of just "genome". I find it more accurately descriptive.

Response:

Edited title as per suggestion.

Comment 4:

Lines 120-125: The MaSuRCA quick start guide (ftp://ftp.genome.umd.edu/pub/MaSuRCA/MaSuRCA_QuickStartGuide.pdf) explicitly says that Illumina reads should not be pre-processed before providing them to MaSuRCA. It is not clear whether or not the "clean" reads were used in the assemblies. Were the "clean" reads used? Or were they only used for genome size estimation?

	<p>Response: Since the raw reads obtained from the Illumina MiSeq already had adapter sequences trimmed off by default, reads used as input to MaSuRCA were only adapter-trimmed. However, no quality-trim/cleaning/error correction was performed on these reads since the program performs its own error correction steps. "Clean" reads which excluded bacteria/virus contaminants and mitochondrial origins were used mainly for genome size estimation to not underestimate the genome size when max kmer coverage is applied in GenomeScope. "Kraken-unclassified reads i.e. non-microbial/viral origin were aligned to the complete mitogenome of NTM A3764 (See "Mitogenome Assembly") to exclude sequences of organellar origin. This results in a total of 42.35 Gb "clean" short reads." -----</p> <p>Comment 5: Line 149: "10 iterations of Pilon". Did you actually see any improvements after this many iterations? How did you assess the improvements?</p> <p>Response: We have added details on the assemblies obtained after each pilon iteration as Supplemental Table 3. While the contiguity of the assembly (indicated by N50) is not improved much after each iteration, the number of gaps and 'N's are reduced and these numbers are observed to plateau at later iterations (i9 or i10). In addition, pilon reports the changes made in the assembly. Based on this output, the number of changes made decreases with each iteration and eventually almost plateaus as well (Supplemental Table 3). -----</p> <p>Comment 6: Line 155: Here you claim a "94 % decrease in the number of scaffolds", while you claim 93 % in the abstract. Which is correct? I guess both if you use different criteria for which scaffolds are included or not (>500 bp).</p> <p>Response: 94% decrease is correct - fixed this in abstract</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
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.	
Have you included all the information requested in your manuscript?	
Resources	Yes
A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely	

<p>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>	
<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 Standards Reporting Checklist?</p>	<p>Yes</p>

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

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
7 5 Laurence J. Croft^{1,5}, Han Ming Gan^{1,2,3*}

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
16 10 Malaysia, Jalan Lagoon Selatan, Bandar Sunway 47500, Petaling Jaya, Selangor,
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

48
49
50
51
52
53 30 Additional ORCID details:

54
55 31 Mun Hua Tan: 0000-0003-3396-8213; Laurence J. Croft: 0000-0001-8471-2408.

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 the range of estimated k-mer-predicted genome sizes of between 791 to 967
52 Mbp. The final assembled genome size is contained in 6,404 scaffolds with an
53 accumulated length of 880 Mb (96.3% BUSCO-calculated genome completeness).
54 Compared to the Illumina-only assembly, the hybrid approach generated 94% fewer
55 scaffolds with 18-fold increase in N₅₀ length (401 kb) and increased the genome
56 completeness by an additional 16%. A total of 27,240 high quality protein-coding
57 genes were predicted from the clown anemonefish, 26,211 (96%) of which were
58 annotated functionally with information from either sequence homology or protein
59 signature searches.

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

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

69

70

71 **Data description**

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

86 In recent years, concurrent with the advent of long read sequencing
87 technologies [14], several studies have explored combining short but accurate
88 Illumina reads with long but less accurate Nanopore/PacBio reads to obtain genome
89 assemblies that are usually more contiguous with higher completeness than
90 assemblies based on Illumina-only reads [15-19]. To further contribute to the
91 evaluation of long read technology in fish genomics [15], we sequenced the whole
92 genome of *A. ocellaris* using Oxford Nanopore and Illumina technologies and
93 demonstrate that hybrid assembly of long and short reads greatly improved the quality
94 of genome assembly.

95

96 **Whole genome sequencing**

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

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

103 Genomic DNA was extracted from multiple fin clip and muscle samples using
104 E.Z.N.A.® Tissue DNA Kit (Omega Bio-tek, Norcross, GA). For Illumina library
105 prep, approximately 1µg of gDNA from isolate A3764 was sheared to 300 bp using a
106 Covaris Focused-ultrasonicator (Covaris, Woburn, MA) and subsequently processed
107 using TruSeq DNA sample prep kit (Illumina, San Diego, CA) according to the
108 manufacturer's instructions. Paired-end sequencing was performed on a single lane of
109 HiSeq 2000 (Illumina, San Diego, CA) located at the Malaysian Genomics Resource
110 Centre Berhad. Two additional libraries were constructed from specimen NTM
111 A3764 and both libraries were sequenced on the MiSeq (2×300 bp setting) located at
112 the Monash University Malaysia Genomics Facility.

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

121 **Sequence read processing**

122 Raw Illumina short reads were adapter-trimmed with Trimmomatic v.0.36
123 (*ILLUMINACLIP:2:30:10, MINLEN:100*) (Trimmomatic , RRID:SCR_011848)[20]
124 followed by a screening for vectors and contaminants, using Kraken v.0.10.5 (Kraken,
125 RRID:SCR_005484)[21] based on the MiniKraken DB. Kraken-unclassified reads i.e.
126 non-microbial/viral origin were aligned to the complete mitogenome of NTM A3764
127 (See “Mitogenome Assembly”) to exclude sequences of organellar origin. This results
128 in a total of 42.35 Gb “clean” short reads. Nanopore reads were base-called from their
129 raw FAST5 files using the Oxford Nanopore propriety base-caller, Albacore version
130 2.0.1. Applying a minimum length cut-off of 500 bp, this study produced a total of
131 8.95 Gbp in 895,672 Nanopore reads (N₅₀: 12.7 kb). A table with sequencing statistics
132 is available as Supplemental Table 1.

134 **Genome size estimation**

135 K-mer counting with the “clean” Illumina reads was performed with Jellyfish v.2.2.6
136 (Jellyfish, RRID:SCR_005491)[22], generating k-mer frequency distributions of 17-,

1 137 21- and 25-mers. These histograms were processed by GenomeScope [23] that
2 138 estimated a genome size of 791 to 794 Mbp with approximately 80% of unique
3 139 content and a heterozygosity level of 0.6% (Supplemental Figure 1). Given that we
4 140 had previously excluded adapters as well as sequences from contaminant or organellar
5 141 sources, the max kmer coverage filter was not applied (*max kmer coverage: -1*). A
6 142 separate estimation performed by BMap [24] estimated a haploid genome size of
7 143 967 Mbp. The genome sizes estimated from both approaches are within the range of
8 144 sizes listed for other *Amphiprion* species (792 Mb - 1.2 Gb) as reported on the Animal
9 145 Genome Size Database (<http://www.genomesize.com> accessed on 11th November
10 146 2017)

148 **Hybrid genome assembly**

149 Short reads used for assemblies described in this study were only trimmed for
150 adapters, but not for quality. Both short-read-only and hybrid *de novo* assemblies
151 were performed with MaSuRCA v.3.2.2 (MaSuRCA, RRID:SCR_010691)[25].
152 During hybrid assembly, errors were encountered in the fragment correction step of
153 the Celera Assembler (CA)(Celera assembler, RRID:SCR_010750). To overcome
154 this, given that the CA assembler is no longer maintained, we disabled the *frgcorr*
155 step based on one of the developer's recommendations and the hybrid assembly was
156 subsequently improved with 10 iterations of Pilon v.1.22 (Pilon ,
157 RRID:SCR_014731)[26], using short reads to correct bases, fix mis-assemblies and
158 fill assembly gaps. To assess the completeness of the genome, BUSCO v.3.0.2
159 (BUSCO , RRID:SCR_015008)[27] was used to locate the presence or absence of the
160 Actinopterygii-specific set of 4,584 single copy orthologs (OrthoDB v9).

161 The short-read-only and hybrid assemblies yielded total assembly sizes of 851
162 Mb and 880Mb, respectively. Statistics for assemblies for each Pilon iteration are
163 available in Supplemental Table 2. Inclusion of Nanopore long reads for hybrid
164 assembly representing approximately 11× genome coverage led to a 94% decrease in
165 the number of scaffolds (> 500 bp) from 106,526 to 6,404 scaffolds and an 18-fold
166 increase in the scaffold N₅₀ length from 21,802 bp to 401,715 bp (Table 1). In
167 addition, the genome completeness was also substantially improved in the hybrid
168 assembly, with BUSCO detecting complete sequences of 96.3% (4,417/4,584) of
169 single-copy orthologs in the Actinopterygii-specific dataset.

170

171 **Transcriptome sequencing and assembly**

172 Total RNA extraction from RNAsShield-preserved whole body and muscle tissues of
173 isolate A4496 used Quick-RNA MicroPrep (Zymo Research Corpt, Irvine, CA)
174 according to the manufacturer's protocols. After assessing total RNA intactness on
175 the TapeStation2100 (Agilent), mRNA was enriched using NEBNext Poly(A) mRNA
176 magnetic isolation kit (NEB, Ipswich, MA) and processed with NEBNext Ultra RNA
177 library prep kit for Illumina (NEB, Ipswich, MA). Libraries from both whole body and
178 muscle tissues were sequenced on a fraction of MiSeq V3 flowcell (1×150 bp).
179 Single-end reads from both libraries in addition to two publicly available *A. ocellaris*
180 transcriptome sequencing data (SRR5253145 and SRR5253146, Bioproject ID:
181 PRJNA374650) were individually assembled using Scallop v0.10.2 [28] based on
182 HiSat2 [29] alignment of RNA-sequencing reads to the newly generated *A. ocellaris*
183 genome. The transcriptome assemblies were subsequently merged using the tr2aacds
184 pipeline from the EvidentialGene [30] package and similarly assessed for
185 completeness using BUSCO version 3 [27]. The final non-redundant transcriptome
186 assembly, which was subsequently used to annotate the *A. ocellaris* genome, contains
187 25,264 contigs/isotigs (putative transcripts) with an accumulated length of 68.4 Mb
188 and BUSCO-calculated completeness of 92.8% (Table 1).

189

190 **Genome annotation**

191 Protein-coding genes were predicted with the MAKER v.2.31.9 genome annotation
192 pipeline (MAKER, RRID:SCR_005309)[31]. A total of three passes were run with
193 MAKER2; the first pass was based on hints from the assembled transcripts as RNA-
194 seq evidence (*est2genome*) and protein sequences from 11 fish species downloaded
195 from Ensembl (Ensembl, RRID:SCR_002344)[32] (*protein2genome*), whereas the
196 second and third passes included gene models trained from the first (and then, second)
197 passes with *ab initio* gene predictors SNAP (SNAP, RRID:SCR_002127) [33] and
198 Augustus (Augustus: Gene Prediction, RRID:SCR_008417)[34]. In the final set of
199 genes predicted, sequences with Annotation Edit Distance (AED) values less than 0.5
200 were retained. A small AED value suggests a lesser degree of difference between the
201 predicted protein and the evidences used in the prediction (i.e. fish proteins,
202 transcripts). This resulted in a final set of 27,240 protein-coding genes with an
203 average AED of 0.14 (Table 1). A BUSCO analysis on the completeness of the

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

204 predicted protein dataset detected the presence of 4,259 (92.9%) single-copy
205 orthologs from the Actinopterygii-specific dataset.

206 Further, to infer putative function of these predicted proteins, NCBI's *blastp*
207 v.2.6.0 (*-evalue 1e-10, -seg yes, -soft_masking true, -lcase_masking*) (BLASTP ,
208 RRID:SCR_001010)[35] was used to find homology to existing vertebrate sequences
209 in the non-redundant (NR) database. Applying a hit fraction filter to include only hits
210 with ≥ 70 % target length fraction, the remaining un-annotated sequences were
211 subsequently aligned to all sequences in the NR database. With this method, 20,107
212 proteins (74%) were annotated with a putative function based on homology.

213 Additionally, InterProScan v.5.26.65 (InterProScan , RRID:SCR_005829)[36] was
214 used to examine protein domains, signatures and motifs present in the predicted
215 protein sequences. This analysis detected domains, signatures or motifs for 26,211
216 proteins (96%). Overall, 96% of the predicted clown fish protein-coding genes were
217 functionally annotated with information from at least one of the two approaches.

218 219 **Mitogenome recovery via genome skimming**

220 Genome skimming [37, 38] was performed on three additional *A. ocellaris* individuals
221 from known localities (Supplemental Table 3). Mitogenome assembly was performed
222 with MITObim version 1.9 (MITObim , RRID:SCR_015056)[39] using the complete
223 mitogenome of *A. ocellaris* (GenBank: NC009065.1) as the bait for read mapping.
224 The assembled mitogenomes were subsequently annotated with MitoAnnotator [40].
225 Consistent with original broodstock collection from northern Australia, the captive
226 bred black and white *A. ocellaris* NTM A3764 exhibits strikingly high whole
227 mitogenome nucleotide identity (99.98%) to sample NTM A3708 as a wild collection
228 from Darwin Harbour, Australia. In addition, the overall high pair-wise nucleotide
229 identity (> 98%) of NTM A3764 to newly generated and publicly available *A.*
230 *ocellaris* whole mitogenomes further supports its morphological identification as *A.*
231 *ocellaris* (Supplemental Table 3).

232 233 **Identification of the *cyp19a1a* gene associated with sexual differentiation**

234 The validated *cyp19a1a* enzyme of *Danio rerio* (Uniprot: O42145) was used as the
235 query (E-value = 1e-10) for *blastp* search against the predicted *A. ocellaris* proteins.
236 The top *blast* hit, AMPOCE_00012675-RA (71.5% protein identity to O42145), was
237 searched (*tblastn*) against the NCBI TSA database (Taxon: *Amphirion*) and showed

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

238 strikingly high protein identity (99%) to a translated RNA transcript from *Amphiprion*
239 *bicinctus* (c183337_g1_i2 : GDCV01327693) [5]. The *cyp19a1a* gene codes for a
240 steroidogenic enzyme that converts androgens into estrogens [41] and was recently
241 shown to be instrumental during sex change in *Amphiprion bicinctus* as evidenced by
242 significant correlation and differential expression of this gene between male and
243 mature females [5]. We also observed similar profile based on mapping of RNA reads
244 from the publicly available male and female transcriptomes of *A. ocellaris* to the
245 *cyp19a1a* gene region as visualized using Integrative Genomics Viewer [42] (Figure
246 2). The *A. ocellaris cyp19a1a* gene is located on a 419 kb scaffold and is spanned by
247 multiple Minimap2-aligned Nanopore reads [43]. It is noteworthy that in the Illumina-
248 only assembly, this gene is fragmented and located on 3 relatively short scaffolds
249 (Figure 2).

250

251 **Conclusion**

252 We present the first clownfish genome co-assembled with high coverage Illumina
253 short reads and low coverage (~11×) Nanopore long reads. Hybrid assembly of
254 Illumina and Nanopore reads is one of the new features of the MaSuRCA assembler
255 version 3.2.2 that works by constructing long and accurate mega-reads from the
256 combination of long and short read data. Although this is a relatively
257 computationally-intensive strategy with long run-times, we observed substantial
258 improvement in the genome statistics when compared to Illumina-only assembly. As
259 Nanopore long technology becomes more mature, it is likely that future *de novo*
260 genome assembly will shift towards high coverage long read-only assembly followed
261 by multiple iterations of genome polishing using Illumina reads.

262

263 **Availability of supporting data**

264 Data supporting the results of this article is available in the GigaDB repository [44].
265 Raw Illumina and Nanopore reads generated in this study are available in the
266 Sequence Read Archive (SRP123679) whereas Whole Genome Shotgun project has
267 been deposited at DDBJ/EMBL/GenBank under the accession NXFZ00000000, both
268 under BioProject PRJNA407816.

269

270 **Acknowledgements**

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

271 This study was funded by Monash University Malaysia Tropical and Biology

272 Multidisciplinary Platform.

273

274 **Competing interests**

275 The authors declare that they have no competing interests

276

277

278 **Figure Legends**

1
2 279

3
4 280 Figure 1. The clown anemonefish (*Amphiprion ocellaris*): Photo by Michael P.

5
6 281 Hammer

7
8 282

9 283 Figure 2. Mapping of MinION long reads, Illumina-assembled scaffolds and RNA-

10
11 284 sequencing reads of male and female *A. ocellaris* to genomic region containing the

12
13 285 *cyp19a1a* gene. Transcripts per million (TPM) values were calculated using Kallisto

14
15 286 version 0.43.1 [45].

16
17 287

18 288 Supplemental Figure 1. Genome profiling of *A. ocellaris* based on Illumina short

19
20 289 reads

21
22 290

23
24 291

25
26 292

27
28 293

29
30 294

31
32 295

33
34 296

35
36 297

37
38 298

39
40 299

41
42 300

43
44 301

45
46 302

47
48 303

49
50 304

51
52 305

53
54 306

55
56 307

57
58 308

59
60 309

61
62 310

63
64 311

65

312 **References**

313

- 314 1. Miltz TA and Foale S. The “Nemo Effect”: Perception and reality of Finding
315 Nemo's impact on marine aquarium fisheries. *Fish and Fisheries*. 2017;18
316 3:596-606.
- 317 2. Madduppa HH, von Juterzenka K, Syakir M and Kochzius M. Socio-economy
318 of marine ornamental fishery and its impact on the population structure of the
319 clown anemonefish *Amphiprion ocellaris* and its host anemones in Spermonde
320 Archipelago, Indonesia. *Ocean & Coastal Management*. 2014;100 Supplement
321 C:41-50. doi:<https://doi.org/10.1016/j.ocecoaman.2014.07.013>.
- 322 3. Hall H and Warmolts D. The Role of Public Aquariums in the Conservation
323 and Sustainability of the Marine Ornamentals Trade. *Marine Ornamental*
324 *Species*. Blackwell Publishing Company; 2008. p. 305-24.
- 325 4. Madhu R, Madhu K and Rethesh T. Life history pathways in false clown
326 *Amphiprion ocellaris* Cuvier, 1830: A journey from egg to adult under captive
327 condition. *Journal of the Marine Biological Association of India*. 2012;54
328 1:77-90.
- 329 5. Casas L, Saborido-Rey F, Ryu T, Michell C, Ravasi T and Irigoien X. Sex
330 Change in Clownfish: Molecular Insights from Transcriptome Analysis.
331 *Scientific Reports*. 2016;6:35461. doi:10.1038/srep35461.
- 332 6. Buston P. Social hierarchies: size and growth modification in clownfish.
333 *Nature*. 2003;424 6945:145-6.
- 334 7. Kobayashi Y, Horiguchi R, Miura S and Nakamura M. Sex- and tissue-
335 specific expression of P450 aromatase (*cyp19a1a*) in the yellowtail clownfish,
336 *Amphiprion clarkii*. *Comp Biochem Physiol A Mol Integr Physiol*. 2010;155
337 2:237-44.
- 338 8. Davenport D and Norris KS. Observations on the symbiosis of the sea
339 anemone *Stoichactis* and the pomacentrid fish, *Amphiprion percula*. *The*
340 *Biological Bulletin*. 1958;115 3:397-410.
- 341 9. Arvedlund M and Nielsen LE. Do the anemonefish *Amphiprion ocellaris*
342 (*Pisces: Pomacentridae*) imprint themselves to their host sea anemone
343 *Heteractis magnifica* (*Anthozoa: Actinidae*)? *Ethology*. 1996;102 2:197-211.
- 344 10. Mariscal RN. An experimental analysis of the protection of *Amphiprion*
345 *xanthurus* Cuvier & Valenciennes and some other anemone fishes from sea
346 anemones. *Journal of Experimental Marine Biology and Ecology*. 1970;4
347 2:134-49.
- 348 11. Hattori A. Coexistence of two anemonefishes, *Amphiprion clarkii* and *A.*
349 *perideraion*, which utilize the same host sea anemone. *Environmental Biology*
350 *of Fishes*. 1995;42 4:345-53.
- 351 12. Schmiege PF, D'Aloia CC and Buston PM. Anemonefish personalities
352 influence the strength of mutualistic interactions with host sea anemones.
353 *Marine Biology*. 2017;164 1:24.
- 354 13. Allen GR. *Damselfishes of the world*. Melle, Germany: Mergus Publishers;
355 1991.
- 356 14. Heather JM and Chain B. The sequence of sequencers: The history of
357 sequencing DNA. *Genomics*. 2016;107 1:1-8.
358 doi:<https://doi.org/10.1016/j.ygeno.2015.11.003>.
- 359 15. Austin CM, Tan MH, Harrisson KA, Lee YP, Croft LJ, Sunnucks P, et al. De
360 novo genome assembly and annotation of Australia's largest freshwater fish,

361 the Murray cod (*Maccullochella peelii*), from Illumina and Nanopore
362 sequencing read. *GigaScience*. 2017;6 8:1-6. doi:10.1093/gigascience/gix063.

16. Gan HM, Lee YP and Austin CM. Nanopore Long-Read Guided Complete
363 Genome Assembly of *Hydrogenophaga intermedia*, and Genomic Insights into
364 4-Aminobenzenesulfonate, p-Aminobenzoic Acid and Hydrogen Metabolism
365 in the Genus *Hydrogenophaga*. *Frontiers in Microbiology*. 2017;8 1880
366 doi:10.3389/fmicb.2017.01880.

17. Zimin AV, Puiu D, Hall R, Kingan S, Clavijo BJ and Salzberg SL. The first
367 near-complete assembly of the hexaploid bread wheat genome, *Triticum*
368 *aestivum*. *GigaScience*. 2017:gix097-gix. doi:10.1093/gigascience/gix097.

18. Zimin AV, Stevens KA, Crepeau MW, Puiu D, Wegrzyn JL, Yorke JA, et al.
369 An improved assembly of the loblolly pine mega-genome using long-read
370 single-molecule sequencing. *GigaScience*. 2017;6 1:1-4.
371 doi:10.1093/gigascience/giw016.

19. Zimin AV, Puiu D, Luo M-C, Zhu T, Koren S, Marcais G, et al. Hybrid
372 assembly of the large and highly repetitive genome of *Aegilops tauschii*, a
373 progenitor of bread wheat, with the mega-reads algorithm. *Genome Research*.
374 2017; doi:10.1101/gr.213405.116.

20. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for
375 Illumina sequence data. *Bioinformatics*. 2014;30 15:2114-20.

21. Wood DE and Salzberg SL. Kraken: ultrafast metagenomic sequence
376 classification using exact alignments. *Genome biology*. 2014;15 3:R46.

22. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel
377 counting of occurrences of k-mers. *Bioinformatics*. 2011;27 6:764-70.
378 doi:10.1093/bioinformatics/btr011.

23. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski
379 J, et al. GenomeScope: fast reference-free genome profiling from short reads.
380 *Bioinformatics*. 2017;33 14:2202-4. doi:10.1093/bioinformatics/btx153.

24. Bushnell B. BBMap short read aligner. University of California, Berkeley,
381 California URL <http://sourceforge.net/projects/bbmap>. 2016.

25. Zimin AV, Marçais G, Puiu D, Roberts M, Salzberg SL and Yorke JA. The
382 MaSuRCA genome assembler. *Bioinformatics*. 2013;29 21:2669-77.

26. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al.
383 Pilon: an integrated tool for comprehensive microbial variant detection and
384 genome assembly improvement. *PloS one*. 2014;9 11:e112963.

27. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM.
385 BUSCO: assessing genome assembly and annotation completeness with
386 single-copy orthologs. *Bioinformatics*. 2015;31 19:3210-2.

28. Shao M and Kingsford C. Accurate assembly of transcripts through phase-
387 preserving graph decomposition. *Nature Biotechnology*. 2017;
388 doi:10.1038/nbt.4020
389 <https://www.nature.com/articles/nbt.4020#supplementary-information>.

29. Kim D, Langmead B and Salzberg SL. HISAT: a fast spliced aligner with low
390 memory requirements. *Nature Methods*. 2015;12:357.
391 doi:10.1038/nmeth.3317
392 <https://www.nature.com/articles/nmeth.3317#supplementary-information>.

30. Gilber D. Gene-omes built from mRNA seq not genome DNA.
393 *F1000Research*. 2016;5 1695:1. doi:10.7490/f1000research.1112594.1.

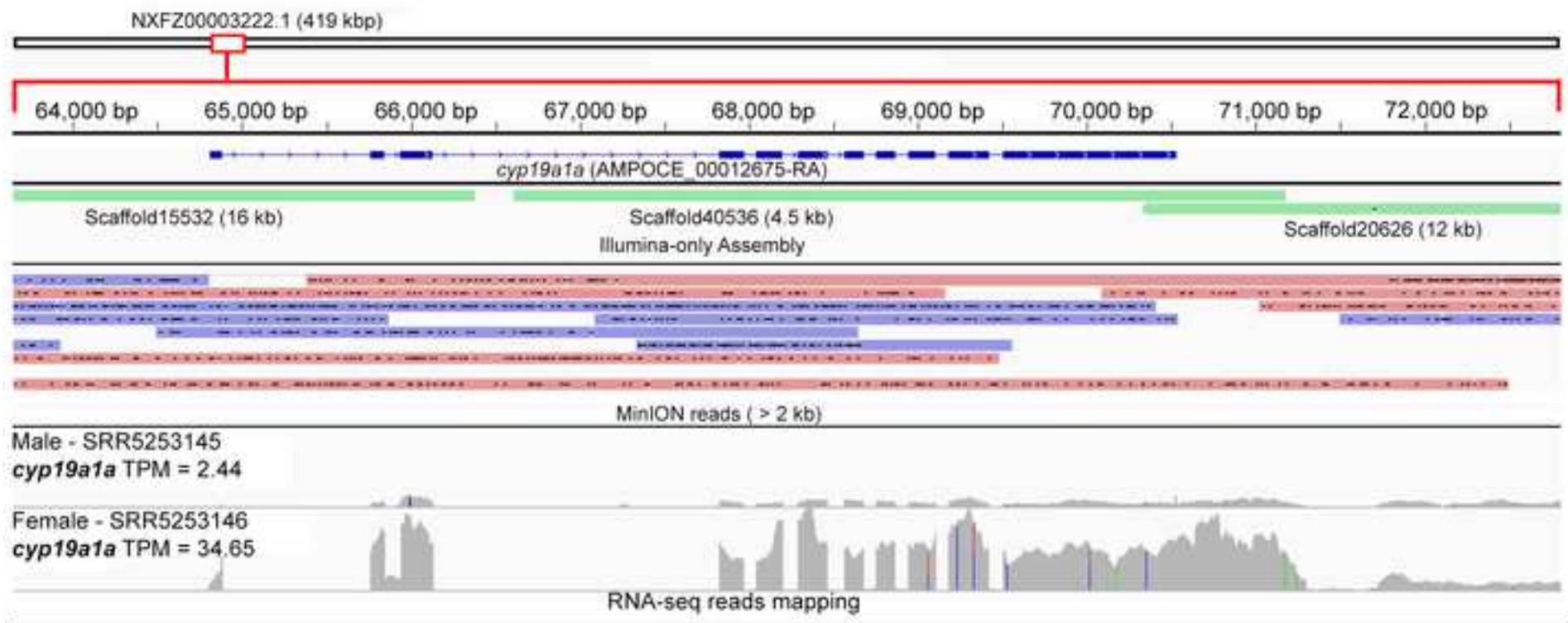
- 409 31. Holt C and Yandell M. MAKER2: an annotation pipeline and genome-
410 database management tool for second-generation genome projects. BMC
411 bioinformatics. 2011;12 1:491.
- 412 32. Hubbard T, Barker D, Birney E, Cameron G, Chen Y, Clark L, et al. The
413 Ensembl genome database project. Nucleic acids research. 2002;30 1:38-41.
- 414 33. Korf I. Gene finding in novel genomes. BMC bioinformatics. 2004;5 1:59.
- 415 34. Stanke M, Steinkamp R, Waack S and Morgenstern B. AUGUSTUS: a web
416 server for gene finding in eukaryotes. Nucleic acids research. 2004;32
417 suppl_2:W309-W12.
- 418 35. Boratyn GM, Camacho C, Cooper PS, Coulouris G, Fong A, Ma N, et al.
419 BLAST: a more efficient report with usability improvements. Nucleic acids
420 research. 2013;41 W1:W29-W33.
- 421 36. Zdobnov EM and Apweiler R. InterProScan—an integration platform for the
422 signature-recognition methods in InterPro. Bioinformatics. 2001;17 9:847-8.
- 423 37. Gan HM, Schultz MB and Austin CM. Integrated shotgun sequencing and
424 bioinformatics pipeline allows ultra-fast mitogenome recovery and confirms
425 substantial gene rearrangements in Australian freshwater crayfishes. BMC
426 evolutionary biology. 2014;14 1:19.
- 427 38. Grandjean F, Tan MH, Gan HM, Lee YP, Kawai T, Distefano RJ, et al. Rapid
428 recovery of nuclear and mitochondrial genes by genome skimming from
429 Northern Hemisphere freshwater crayfish. Zoologica Scripta. 2017.
- 430 39. Hahn C, Bachmann L and Chevreur B. Reconstructing mitochondrial
431 genomes directly from genomic next-generation sequencing reads—a baiting
432 and iterative mapping approach. Nucleic acids research. 2013;41 13:e129-e.
- 433 40. Iwasaki W, Fukunaga T, Isagozawa R, Yamada K, Maeda Y, Satoh TP, et al.
434 MitoFish and MitoAnnotator: A mitochondrial genome database of fish with
435 an accurate and automatic annotation pipeline. Molecular biology and
436 evolution. 2013;30 11:2531-40.
- 437 41. Kallivretaki E, Eggen R, Neuhauss S, Alberti M, Kausch U and Segner H.
438 Aromatase in zebrafish: a potential target for endocrine disrupting chemicals.
439 Mar Environ Res. 2006;62 90:7.
- 440 42. Thorvaldsdóttir H, Robinson JT and Mesirov JP. Integrative Genomics Viewer
441 (IGV): high-performance genomics data visualization and exploration.
442 Briefings in bioinformatics. 2013;14 2:178-92.
- 443 43. Li H. Minimap2: fast pairwise alignment for long nucleotide sequences. arXiv.
444 2017;1708:Artciel 01492.
- 445 44. Tan, M, H; Austin, C, M; Hammer, M, P; Lee, Y, P; Croft, L, J; Gan, H, M
446 (2017): Supporting data for "Finding Nemo: Hybrid assembly with Oxford
447 Nanopore and Illumina reads greatly improves the Clownfish (*Amphiprion*
448 *ocellaris*) genome assembly" GigaScience Database.
449 <http://dx.doi.org/10.5524/100397>
- 450 45. Bray NL, Pimentel H, Melsted P and Pachter L. Near-optimal probabilistic
451 RNA-seq quantification. Nat Biotech. 2016;34 5:525-7. doi:10.1038/nbt.3519
452 [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)
453 [information.](http://www.nature.com/nbt/journal/v34/n5/abs/nbt.3519.html#supplementary-information)
454
455

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)





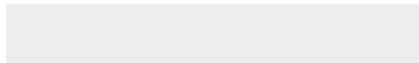


Click here to access/download
Supplementary Material
Supplemental Figure 1_091217.tif



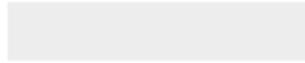


Click here to access/download
Supplementary Material
Supplemental Table 1_091217.docx





Click here to access/download
Supplementary Material
Supplemental Table 2_091217.docx





Click here to access/download
Supplementary Material
Supplemental Table 3_091217.docx

