# GigaScience

# High-quality genome assembles from key Hawaiian coral species
## --Manuscript Draft--

| Manuscript Number: | GIGA-D-22-00143 | |
|---|---|---|
| Full Title: | High-quality genome assembles from key Hawaiian coral species | |
| Article Type: | Data Note | |
| Funding Information: | Paul G. Allen Family Foundation | Dr Eva Majerová |
| | USDA National Institute of Food and Agriculture (1017848) | Dr Hollie M. Putnam |
| | National Science Foundation (NSF-OCE 1756623) | Dr Hollie M. Putnam |
| | National Science Foundation (NSF-OCE 1756616) | Dr Debashish Bhattacharya |
| | Catalyst Science Fund (2020-008) | Dr Debashish Bhattacharya |
| | National Institute of Food and Agriculture and United States Department of Agriculture (NJ01180) | Dr Debashish Bhattacharya |
| | National Aeronautics and Space Administration (80NSSC19K0462) | Dr Debashish Bhattacharya |
| | Ministry of Oceans and Fisheries (20180430) | Dr Hwan Su Yoon |
| | National Research Foundation of Korea (2020R1C1C1010193) | Dr JunMo Lee |
| | Korea Ministry of Environment (2021003420004) | Dr JunMo Lee |

| Abstract: | Background<br>Coral reefs house about 25% of marine biodiversity and are critical for the livelihood of many communities by providing food, tourism revenue, and protection from wave surge. These magnificent ecosystems are under existential threat from anthropogenic climate change. Whereas extensive ecological and physiological studies have addressed coral response to environmental stress, high-quality reference genome data are lacking for many of these species. The latter issue hinders efforts to understand the genomic and genetic basis of stress resistance and to design informed coral conservation strategies.<br>Results<br>We report genome assemblies from four key Hawaiian coral species, Montipora capitata, Pocillopora acuta, Pocillopora meandrina, and Porites compressa. These species, or members of these genera, are distributed worldwide and therefore of broad scientific and ecological importance. For M. capitata, an initial assembly was generated from short-read Illumina and long-read PacBio data, which was then scaffolded into 14 putative chromosomes using Omni-C sequencing. For Poc. acuta, Poc. meandrina, and Por. compressa, high-quality assemblies were generated using short-read Illumina and long-read PacBio data. The Poc. acuta assembly is from a triploid individual, making it the first reference genome of a non-diploid coral animal.<br>Conclusions<br>These assemblies are significant improvements over available data and provide invaluable resources for supporting multi-omics studies into coral biology, not just in Hawai'i, but also in other regions, where related species exist. The Poc. acuta assembly gives us, for the first time, a platform for studying polyploidy in corals, and its role in genome evolution and stress adaptation in these organisms. | |

| Corresponding Author: | Timothy Gordon Stephens, Ph.D.<br>Rutgers University: Rutgers The State University of New Jersey<br>New Brunswick, New Jersey UNITED STATES | |

| Corresponding Author Secondary Information: | |
| --- | --- |
| Corresponding Author's Institution: | Rutgers University: Rutgers The State University of New Jersey |
| Corresponding Author's Secondary Institution: | |
| First Author: | Timothy Gordon Stephens, Ph.D. |
| First Author Secondary Information: | |
| Order of Authors: | Timothy Gordon Stephens, Ph.D. |
| | JunMo Lee |
| | YuJin Jeong |
| | Hwan Su Yoon |
| | Hollie M. Putnam |
| | Eva Majerová |
| | Debashish Bhattacharya |
| Order of Authors Secondary Information: | |

| Additional Information: | |
| --- | --- |
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>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.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>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. | Yes |

| | |
|---|---|
| Have you included the information requested as detailed in our [Minimum Standards Reporting Checklist](#)? | |
| **Availability of data and materials**<br><br>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.<br><br>Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)? | Yes |

# High-quality genome assembles from key Hawaiian coral species

Timothy G. Stephens (ts942@sebs.rutgers.edu)[1,*], JunMo Lee (leejunmo331@gmail.com)[2], YuJin Jeong (lpple0826@knu.ac.kr)[2], Hwan Su Yoon (hsyoon2011@skku.edu)[3], Hollie M. Putnam (hputnam@uri.edu)[4], Eva Majerová (majerova@hawaii.edu)[5], and Debashish Bhattacharya (dbhattac@rutgers.edu)[1]

[1]Department of Biochemistry and Microbiology, Rutgers University, New Brunswick, NJ 08901, USA.

[2]Department of Oceanography, Kyungpook National University, Daegu, Buk-gu 41566, Korea.

[3]Department of Biological Sciences, Sungkyunkwan University, Suwon 16419, Korea.

[4]Department of Biological Sciences, University of Rhode Island; Kingston, RI 02881, USA.

[5]Hawaiʻi Institute of Marine Biology, PO Box 1346 Kāneʻohe HI 96744, USA.

*Corresponding author (ts942@sebs.rutgers.edu)

## Abstract

**Background**

Coral reefs house about 25% of marine biodiversity and are critical for the livelihood of many communities by providing food, tourism revenue, and protection from wave surge. These magnificent ecosystems are under existential threat from anthropogenic climate change. Whereas extensive ecological and physiological studies have addressed coral response to environmental stress, high-quality reference genome data are lacking for many of these species. The latter issue hinders efforts to understand the genomic and genetic basis of stress resistance and to design informed coral conservation strategies.

**Results**

We report genome assemblies from four key Hawaiian coral species, *Montipora capitata*, *Pocillopora acuta*, *Pocillopora meandrina*, and *Porites compressa*. These species, or members of these genera, are distributed worldwide and therefore of broad scientific and ecological importance. For *M. capitata*, an initial assembly was generated from short-read Illumina and long-read PacBio data, which was then scaffolded into 14 putative chromosomes using Omni-C sequencing. For *Poc. acuta*, *Poc. meandrina*, and *Por. compressa*, high-quality assemblies were generated using short-read Illumina and long-read PacBio data. The *Poc. acuta* assembly is from a triploid individual, making it the first reference genome of a non-diploid coral animal.

**Conclusions**

These assemblies are significant improvements over available data and provide invaluable resources for supporting multi-omics studies into coral biology, not just in Hawaiʻi, but also in other regions, where related species exist. The *Poc. acuta* assembly gives us, for the first time, a platform for studying polyploidy in corals, and its role in genome evolution and stress adaptation in these organisms.

2

## Background

32

33 *Montipora capitata*, *Pocillopora acuta*, *Pocillopora meandrina*, and *Porites compressa* are

34 species of scleractinian corals that are widespread in the Hawaiian Islands, with *M. capitata* and

35 *Por. compressa* being dominant reef builders. These species are members of cosmopolitan

36 genera, with closely related taxa inhabiting reefs across the Great Barrier Reef and the Coral

37 Triangle [1-3], as well as other regions, such as *Pocillopora* in Panama [4]. In recent years, due

38 to their critical importance to Hawaiian reef ecosystems, the ease of accessibility of these species

39 to researchers working in the United States, and the growing threat that anthropogenic climate

40 change poses to global reef ecosystem, these four species have become the subject of many stress

41 (including thermal [5-7] and acidification [8, 9]), microbiome [10, 11], and population genomic

42 [12-15] studies (among many others). Given the significant interest in these species as models

43 for coral biology, there is a pressing need to generate high-quality reference data to provide a

44 solid foundation for future research.

45

46 A genome assembly for *M. capitata* was published in 2019 by our group [16] using Pacific

47 Biosciences (PacBio) RSII data. This assembly was significantly larger (886 Mbp) than any of

48 the other coral genomes available at that time (ca. 300-500 Mbp), and is larger than any

49 *Montipora* species genome [17, 18] that has since been published. The published assembly

50 contains a high number (>18% [19]) of duplicated BUSCO genes, suggesting the presence of

51 haplotigs (i.e., sequences derived from different homologous chromosomes) that were not

52 removed during the assembly process. There are currently published genomes for three

53 *Pocillopora* [4, 20, 21] species, none of which are from Hawaiʻi. One of these is a *Poc. acuta*

54 isolate collected from Lombok, Indonesia [22] that was generated using Illumina short-read data.

55 This genome assembly is highly fragmented, consisting of 168,465 scaffolds, and whereas it

56 does have a scaffold N50 of 147 Kbp, the contig N50 is only 9,649 bp. The completeness of the

57 genes predicted in this genome is not high, with only 56% of the core eukaryotic genes [20]

58 identified in the reported "*ab initio*" predicted gene set. A second set of predicted genes inferred

59 using RNA-seq evidence (termed the "experimental" set) contains 93% of core eukaryotic genes,

60 however, this set does not have predicted open reading frames (i.e., it includes both coding and

61 non-coding genes), making it difficult to make a direct comparison with other published

3

62  genomes. There are currently three *Porites* species with published genomes [23-25], while they

63  are all of high completeness and reasonable contiguity, none are from Hawai'i.

64

65  As the cost of genome sequencing, in particular, long-read methods continues to decrease,

66  opportunities arise to generate genome data from understudied species or species that have

67  genomes of lower quality that would benefit from the improvement gained from the newer

68  technologies. Furthermore, technologies such as Dovetail Omni-C, which provides long range

69  linkage information, enables the generation of genome assemblies that are at (or near)

70  chromosomal-level resolution. In this study, we generated an improved reference genome

71  assembly for our previously published Hawaiian *M. capitata* using long-read PacBio, short-read

72  Illumina, and newly generated Omni-C data, that is of chromosome-level resolution. The 14

73  largest scaffolds resulting from this assembly likely represent the 14 chromosomes predicted in

74  *Montipora* species [26]. We also generated, using PacBio HiFi data (i.e., circular consensus

75  corrected PacBio reads), high-quality genome assemblies for two *Pocillopora* and one *Porites*

76  species. The sequenced *Poc. acuta* isolate is a triploid, making it the first non-diploid coral

77  genome to be published.

78

## Data description

### Sample collection and processing

81  The four coral species targeted in this study were collected from Kāneʻohe Bay, Hawaiʻi. For *M.*

82  *capitata*, the initial PacBio and Illumina-based assembly was generated using sperm DNA (see

83  [16]). Input DNA for the Dovetail Genomics approach (https://dovetailgenomics.com), using the

84  Omni-C assay and workflow, was a bleached nubbin (a ~5 x 5cm fragment) from a colony that

85  was greatly reduced in algal symbionts (GPS coords: 21.474465, -157.834468; SRA BioSample:

86  SAMN21845729). This fragment was collected under Hawaiʻi Department of Aquatic Resources

87  Special Activity Permit 2019-60, snap frozen in liquid nitrogen, and stored at -80°C before it was

88  shipped on dry ice to Dovetail Genomics (https://dovetailgenomics.com) for processing using

89  their Omni-C assay and workflow.

90

91  For *Poc. meandrina*, one nubbin (a ~5 x 5cm fragment) was collected from an adult colony from

92  Reef 13 (GPS coords: 21.450803, -157.794692) on 2020-09-05 (SRA BioSample:

93     SAMN21845732, SAMN21845733, and SAMN21845734) under DAR-2021-33, Amendment

94     No. 1 to HIMB. High molecular weight DNA was extracted using the QIAGEN Genomic-tip

95     100/G (Cat #: 10223), the QIAGEN Genomic DNA Buffer Set (Cat #: 19060), QIAGEN RNase

96     A (100mg/mL concentration: Cat #: 19101), QIAGEN Proteinase K (Cat #: 19131), and DNA lo-

97     bind tubes (Eppendorf Cat #: 022431021). In brief, a clipping of the coral fragment was placed in

98     a cleaned and sterilized mortar and pestle and ground to powder on liquid nitrogen. High

99     molecular weight DNA was then extracted according to the manufacturer's instructions for

100    preparation of tissue samples in the QIAGEN Genomic DNA Handbook (version 06/2015).

101    For *Poc. acuta*, one nubbin was collected from an adult colony from a reef next to the Hawai'i

102    Institute of Marine Biology (GPS coords: 21.436056, -157.786861) on 2018-09-05 (SRA

103    BioSample: SAMN22898959) under Special Activity Permit 2019-60. High molecular weight

104    DNA was extracted using the QIAGEN Genomic-tip 100/G approach outlined for *Poc.*

105    *meandrina* above. High molecular weight DNA from *Poc. meandrina* and *Poc. acuta* was sent to

106    DNA Link Sequencing Lab (https://www.dnalinkseqlab.com) for sequencing on their PacBio

107    Sequel 2 and Illumina NovaSeq 6000 platforms.

108

109    For *Por. compressa*, DNA was extracted from sperm released at 11 pm on 09 June 2017 from a

110    single colony in Kāne'ohe Bay, O'ahu. Total genomic DNA was extracted using the CTAB

111    protocol and the DNeasy Blood and Tissue Kit (Qiagen, Germany) with subsequent clean-up

112    steps. Genomic data were generated using the PacBio RSII platform. To increase the sequence

113    quality of the assembly, a polishing step was done using the Arrow consensus caller. To this end,

114    we generated a total of 20 Gbp of high-throughput sequencing data (Illumina HiSeq2000; 100 bp

115    paired-end library) as follows. The whole-genome sequencing library of *Por. compressa* was

116    prepared using the Truseq Nano DNA Prep Kit (550bp) protocol following the manufacturer's

117    instructions. Randomly sheared genomic DNA was ligated with index adapters and purified. The

118    ligated products were size-selected for 300-400 bp and amplified using the adapter-specific

119    primers. Library quality was checked using a 2100 BioAnalyzer (Agilent Technologies, Santa

120    Clara, CA, USA).

121

122    **RNA Extractions**

123   RNA was extracted by clipping a small piece of coral using clippers sterilized in 10% bleach,

124   deionized water, isopropanol, and RNAse free water, and then placed in 2 mL Fisherbrand™

125   Pre-Filled Bead Mill microcentrifuge tube containing 0.5mm glass beads (Fisher Scientific

126   Catalog. No 15-340-152) with 1000 µL of Zymo DNA/RNA shield. A two-step extraction

127   protocol was used to extract RNA and DNA, with the first step as a "soft" homogenization to

128   reduce shearing of RNA or DNA. Tubes were vortexed at high speed for 1 and 2 minutes for

129   *Poc. acuta* and *M. capitata* fragments, respectively. The supernatant was removed and

130   designated as the "soft extraction". Second, an additional 500 µL of Zymo DNA/RNA shield was

131   added to the bead tubes and placed in a Qiagen TissueLyser for 1 minute at 20 Hz. The

132   supernatant was removed and designated as the "hard extraction". Subsequently, 300 µL of

133   sample from both soft and hard homogenate was extracted with the Zymo Quick-DNA/RNA

134   Miniprep Plus Kit (Zymo Cat D7003) Protocol with the following modifications. RNA quantity

135   (ng_µL) was measured with a ThermoFisher Qubit Fluorometer, DNA quality was assessed

136   using gel electrophoresis, and RNA quality was measured with an Agilent TapeStation System.

137

138   **Haploid genome assembly of Hawaiian coral species**

139   The long-read genome sequencing data (PacBio) of the Hawaiian coral species were initially

140   assembled using CANU (v2.2; default options) [27]. The PacBio reads for *M. capitata* (78.3

141   Gbp; SRR17163565) and *Por. compressa* (63.3 Gbp; SRR12695159 – SRR12695166) were

142   generated using the PacBio RSII platform ('-pacbio' option for CANU assembler). The PacBio

143   reads for *Poc. meandrina* (311.8 Gbp; SRR16077713), and *Poc. acuta* (239.1 Gbp;

144   SRR16077715) were generated using PacBio HiFi platform ('-pacbio-hifi' option for CANU

145   assembler). An error correction step (nucleotide correction of assembly) using the initial

146   assemblies of *M. capitata* (1.2 Gbp; Supplementary Table S1), *Por. compressa* (1.0 Gbp), *Poc.*

147   *meandrina* (0.7 Gbp), and *Poc. acuta* (1.1 Gbp) was done using bowtie2 (v2.4.2; default options)

148   [31] and the Pilon program (v1.23; default options) [28] with the Illumina short-read sequencing

149   data (HiSeq2500: 27.4 Gbp of *M. capitata* [SRR8497577]; HiSeq 2000: 20.9 Gbp of *Por.*

150   *compressa* [SRR12695158]; NovaSeq 6000: 27.2 Gbp of *Poc. meandrina* [SRR16077712], and

151   23.0 Gbp of *Poc. acuata* [SRR16077714]). Before using the Illumina data, quality trimming and

152   adapter clipping of the raw reads were done using Trimmomatic (v0.39; default options) [29]. To

153   remove potential contaminant sequences, assembly results were analyzed using BLASTn (*e-*

154 value cutoff = 1.e$^{-10}$ cutoff) analysis with the nr database (downloaded: Feb. 2019). To estimate

155 genome size and ploidy of the Hawaiian coral species, *k*-mer analysis was done using Jellyfish

156 (21-mer) [30] with the Illumina short-read data.

157       To reconstruct haploid genomes using the initial assemblies of the Hawaiian coral

158 species, we used the following protocol. First, we predicted repetitive DNA sequences in the

159 initial assemblies and constructed soft-masked assemblies. Repetitive DNA elements were

160 identified using the RepeatModeler pipeline (v2.0.1;

161 http://www.repeatmasker.org/RepeatModeler/) [31-33] which includes RECON (v1.08) and

162 RepeatScout (v1.0.6) as *de novo* repeat finding programs. We used the default options for l-mer

163 size and removed low-complexity and tandem repeats. To classify repeat content, the libraries

164 were constructed from giri repbase (http://www.girinst.org). The consensus sequences of repeat

165 families were used to analyze corresponding repeat regions with RepeatMasker (v4.1.1; default

166 options with soft-masked; http://www.repeatmasker.org/). The second step in the protocol was to

167 infer assemblies as haploid genomes using the HaploMerger2 (HM2) program (the latest release,

168 20180603) [34] and the soft-masked assemblies. The third step was validation of duplicated

169 eukaryotic core genes in the haploid genome assemblies using the Benchmarking Universal

170 Single-Copy Orthologs (BUSCO) program (v4.1.4; genome-based analysis with

171 eukaryota_odb10 dataset) [35]. The final step was to repeat the HM2 analysis until the number of

172 duplicated eukaryotic core genes decreased to under 1%, or the value could not be decreased any

173 further in the haploid assemblies (Supplementary Table S1). The purged assembly of *M. capitata*

174 was sent to Dovetail Genomics along with an additional coral fragment (see above) that was used

175 for high molecular weight DNA extraction for analysis using their Omni-C assay and HiRise

176 v2.2.0 assembly workflow. A total of 56.5 million read-pairs of Dovetail Genomics Omni-C

177 sequencing data (SRR16077716) were generated and used for scaffolding. This step produced a

178 final genome assembly that was at putative chromosome level resolution for *M. capitata*.

179

180 **Gene prediction and functional annotation**

181 Quality trimming and adapter removal of the RNA sequencing (RNA-seq) data in the Hawaiian

182 coral species (*M. capitata* [77.5 Gbp; SRR14729868 – SRR14729873, SRR14729878,

183 SRR14729881, SRR14729889, SRR14729890, SRR14729893, and SRR14729894], *Por.*

184 *compressa* [76.5 Gbp; SRR14729874 – SRR14729877, SRR14729879 – SRR14729880,

185   SRR14729882 – SRR14729888, SRR14729891, and SRR14729892], *Poc*. *acuta* [656.7 Gbp;

186   SRR14610884 – SRR14610890, SRR14610892 – SRR14610901, SRR14610903 –

187   SRR14610912, SRR14610914 – SRR14610923, SRR14610925 – SRR14610932,

188   SRR14610975, SRR14610977 – SRR14610986, SRR14610988 – SRR14610997, SRR14610999

189   – SRR14611008, SRR14611010 – SRR14611019, SRR14611021 – SRR14611030,

190   SRR14611033 – SRR14611042, SRR14611044 – SRR14611053, and SRR14611055 –

191   SRR14611057], and *Poc*. *meandrina* [10.6 Gbp; SRR16077711]) were done using Trimmomatic

192   (v0.39; default options) [29]. These data were assembled using Trinity v2.11 with the default

193   option of *de novo* transcriptome assembly [36, 37]. The trimmed RNA-seq raw reads, and the

194   assembled transcriptomes were aligned to the haploid genome assemblies using the STAR

195   aligner (v2.6.0c; default options for the raw reads), and the STARlong aligner (v2.6.0c; --

196   runMode alignReads --alignIntronMin 10 --seedPerReadNmax 100000 --seedPerWindowNmax

197   1000 --alignTranscriptsPerReadNmax 100000 --alignTranscriptsPerWindowNmax 10000),

198   respectively [38]. Based on each alignment (i.e., bam file), gene predictions were done using the

199   BRAKER2 pipeline (v2.1.5; http://bioinf.uni-greifswald.de/bioinf/braker) [39], which includes

200   GeneMark-ET [40] and AUGUSTUS [41] with default (automatically optimized) options. From

201   the two types (i.e., RNA-seq, and assembled transcriptome) of gene models from the Hawaiian

202   coral species, the best (longest) gene models were manually selected based on results of

203   BLASTp search (*e*-value cutoff = 1.e$^{-5}$ cutoff). Functional annotation of gene models was done

204   using the NCBI Conserved Domain Search (CD-Search) [42], the eggNOG-mapper [43], and the

205   KEGG Automatic Annotation Server (KAAS) [44].

206

207   **Genomes of corals used for comparative analysis**

208   The genome assemblies and predicted genes from the four *Montipora* (*M. cactus* [17], *M.*

209   *capitata* from the Hawaiian Waiopae tide pools [18], *M. efflorescens* [17], and the previous

210   version of the Hawaiian *M. capitata* isolate [16] that we assembled in this study), three

211   *Pocillopora* (*Poc. damicornis* [4], *Poc. acuta* [from Indonesia] [22], and *Poc. verrucosa* [21]),

212   and four *Porites* (*Por. astreoides* [25], *Por. australiensis* [24], *Por. lutea* [23], and *Por. rus* [45])

213   species that have been sequenced were retrieved from their respective repositories

214   (Supplementary Table S2) and used for comparative analysis with the assemblies generated in

215   this study. The *M. cactus* and *M. efflorescens* genome assemblies (from

216    https://marinegenomics.oist.jp [17]) were filtered, retaining only scaffolds identified by Yuki, Go

217    [19] as not being haplotigs. The updated gene models from Yuki, Go [19] were used in place of

218    those available with the original assemblies. For species where just the gene modes were

219    provided (in gff format), gffread v0.11.6 (-S -x cdsfile -y pepfile) [46] was used to infer the

220    protein and CDS sequences. Open Reading Frames (ORFs) were predicted in the RNA-Seq

221    based "experimental" genes predicted in the Indonesian *Poc. acuta* isolate [22], using

222    TransDecoder v5.5.0 (https://github.com/TransDecoder/TransDecoder); HMMER v3.1b2 was

223    used to search the candidate ORFs against the Pfam database (release 33.1; i-Evalue < 0.001)

224    and BLASTP (v2.10.1; -max_target_seqs 1 -evalue 1e-5) was used to search candidate ORFs

225    against the SwissProt database (release 2020_05), with the resulting homology information used

226    by TransDecoder to guide ORF prediction. Only the longest transcript per gene had ORFs

227    predicted and single-exon genes without strand information were assumed to be from the

228    forward/positive strand (TransDecoder will change the strand of single exon genes if required,

229    based on the results of ORF prediction).

230

231    **Genome size estimation**

232    The genome size and ploidy of the new (this study) and published *Montipora*, *Pocillopora*, and

233    *Porites* species (except the Indonesian *Poc. acuta* which does not have read data available to

234    download, *Por. rus* which only had reads from the holobiont [i.e., reads from the coral, algal

235    symbiont, and associated bacteria] available, and *Por. astreoides* which only had PacBio long

236    reads available) were estimated using the GenomeScope2 and Smudgeplot tools [47]. For each

237    species, the available short-reads genome sequencing data were retrieved from NCBI SRA

238    (Supplementary Table S2), trimmed using cutadapt v3.5 [48] (-q 20 --minimum-length 25 -a

239    AGATCGGAAGAGCACACGTCTGAACTCCAGTCA -A

240    AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT), and decomposed into *k*-mers using

241    Jellyfish [30] (v2.3.0; k=21). The *k*-mer frequency histogram produced by Jellyfish (using the

242    'jellyfish histo' command) was imported into GenomeScope2 with a theoretical diploid model

243    fitted with the data (Fig. 1C, D, and F and Supplementary Fig. S1); a theoretical triploid model

244    was fitted with the Hawaiian *Poc. acuta* data (Fig. 1E and Supplementary Fig. S1F) because it

245    was found to be a triploid after initial analysis using Smudgeplot and GenomeScope2.

246    Smudgeplot was run using the *k*-mers extracted by Jellyfish (following the workflow from

247    https://github.com/KamilSJaron/smudgeplot/wiki/manual-of-smudgeplot-with-jellyfish), with

248    thresholds for the lower *k*-mer coverage cutoff (just after the minimum between the initial error

249    peak and the first major peak) and upper *k*-mer coverage cutoff (8.5 times the coverage of the

250    first major coverage peak) chosen for each species using the GenomeScope2 profile shown in

251    Supplementary Figure S1. The "smudge plots" shown in Supplementary Figure S1 were

252    generated using the haploid coverage values estimated by GenomeScope2. The cutoffs used

253    when running Smudgeplot for each species are shown in Supplementary Table S2.

254

255    **Conformation of sample ploidy**

256    The program nQuire [49] (retrieved 7/7/2021 from https://github.com/clwgg/nQuire), which uses

257    the frequency distribution of bi-allelic variant sites inferred from aligned reads to model the

258    ploidy of a sample, was used to verify the ploidy of the four genomes sequenced in this study.

259    Briefly, bowtie2 v2.4.4 ('--very-sensitive --no-unal') was used to align the trimmed (by cutadapt;

260    described previously) Illumina short-reads against their respective genome assemblies; aligned

261    reads were coordinate sorted using samtools v1.11 [50]. The aligned and sorted BAM files were

262    converted into "BIN" files using nQuire ('nQuire create -q 20 -c 20 -x'), filtering for reads with a

263    minimum mapping quality of 20 and sites with a minimum coverage of 20. Denoised BIN files

264    were created using the "nQuire denoise" command run on the initial BIN files. The delta Log-

265    Likelihood values for each ploidy model (diploid, triploid, and tetraploid) was calculated by the

266    "nQuire lrdmodel" command for each of the initial and denoised BIN files. The lower the delta

267    Log-Likelihood value of a given model the better fit it is for the frequency distribution of the bi-

268    allelic variant sites extracted from the aligned reads; the ploidy of the sample is there for

269    assumed to be the ploidy model with the lowest delta Log-Likelihood value. The nQuire results

270    are shown in Supplementary Table S4.

271

272    **Assessment of completeness using BUSCO**

273    The "completeness" of the genome assemblies and predicted genes (published in this study and

274    from previous studies; Supplementary Table S3) were assessed using BUSCO v5.0.0 ('--mode

275    genome' and '--mode protein', respectively) with the eukaryota_odb10 (release 2020-09-10) and

276    metazoa_odb10 datasets (release 2021-02-24) [51].

277

278 **Analysis of extra-chromosomal scaffolds**

279 The proteins predicted on the extra-chromosomal scaffolds (i.e., the scaffolds that do not

280 comprise the 14 putative chromosomes) in the *M. capitata* assembly were compared against the

281 proteins from the chromosomal scaffolds using BLASTP v1.10.1 [52]; the resulting hits were

282 filtered using an *e*-value cutoff $< 1\text{x}10^{-5}$. Additional filtering steps were applied to produce two

283 sets of hits: for the first (lenient) set, hits were retained if they had a query coverage of $> 75\%$

284 and an identity $> 75\%$, with the single best (*e*-value-based) top hit kept for each query sequence;

285 for the second (stringent) set, hits were retained if they had a query coverage of $> 95\%$ and an

286 identity $> 95\%$, with the single best (*e*-value-based) top hit kept for each query sequence. The

287 lenient filtered top hits were used to determine if the extra-chromosomal scaffolds tend to encode

288 genes that have similarity to a single, or multiple, chromosomes. For this analysis, only proteins

289 with top hits to the chromosomal scaffolds (i.e., proteins with hits that have an *e*-value $< 1\text{x}10^{-5}$,

290 query coverage $> 75\%$, and an identity $> 75\%$) were considered, and only scaffolds with multiple

291 proteins with top hits were considered.

292

293 # Data Validation and Quality Control

294 *Montipora capitata* **genome assemblies**

295 The *M. capitata* assembly generated in the study (assembly version V3.0; hereinafter the "new"

296 Hawaiian *M. capitata* genome assembly) has fewer assembled bases (781 Mbp vs. 886 Mbp) and

297 scaffolds (1,699 vs. 3,043), and a vastly improved N50 (47.7 Mbp vs. 0.54 Mbp; Supplementary

298 Table S3), compared to the assembly of the same Hawaiian *M. capitata* isolate (hereinafter the

299 "old" Hawaiian *M. capitata* genome assembly) that was previously published by our group [16].

300 The 14 largest scaffolds in the new assembly, ranging in size from ~22 to ~69 Mbp, likely

301 represent the 14 chromosomes predicted in other *Montipora* species (Figs. 1A and B) [26]. These

302 putative chromosomes total 680 Mbp of assembled sequence, which is only slightly larger than

303 the estimated genome size of 644 Mbp (Fig. 1C; estimated by GenomeScope2 [47] using *k*-mers

304 of size 21 bp). The estimated genome size of the other published *Montipora* species is ~700

305 Mbp, whereas the estimated genome size of the new Hawaiian *M. capitata* genome is 644 Mbp

306 (although the assembly is a little larger; see discussion below). This suggests that species in the

307 genus *Montipora* have genomes that are marginally smaller than 700 Mbp in size.

11

308     The *M. capitata* isolate that was sequenced appears to be a diploid, with a good fit

309     between its *k*-mer frequency histogram and the theoretical diploid model implemented in

310     GenomeScope2 (black line in Fig. 1C and Supplementary Fig. S1A), and a clear "smudge"

311     (bright yellow region in Supplementary Fig. S1A) of *k*-mer pairs with a coverage of 2n and a

312     normalized coverage of 1/2; all of which suggests that the sample is diploid. nQuire also

313     predicted that the *M. capitata* sample was a diploid (i.e., the diploid model had the lowest delta

314     Log-Likelihood value; Supplementary Table S4), supporting the results of GenomeScope2 and

315     Smudegeplot.

316     Compared with the old assembly, the new *M. capitata* assembly has a slightly higher

317     BUSCO completeness for both the Metazoa (from 95.2% to 95.7%, respectively) and Eukaryota

318     (from 97.7% to 99.2%, respectively) datasets (Supplementary Table S3) but a significantly

319     reduced number of duplicated BUSCO genes for both the Metazoa (from 21.2% to 1.6%,

320     respectively) and Eukaryota (from 22.0% to 1.2%, respectively) datasets. The high number of

321     duplicated BUSCO genes in the old assembly is likely a result of haplotigs that evaded removal

322     during the assembly process; this problem appears to have been resolved in the new assembly.

323     Compared with the other published *Montipora* genomes, the new *M. capitata* assembly is the

324     most contiguous and complete to date, with a significantly higher N50 (47.7 Mbp compared to

325     the next best of 1.2 Mbp in *M. efflorescens*) and BUSCO completeness (e.g., 99.2% Eukaryota

326     dataset completeness compared to the next best of 92.1% in *M. cactus*). As the same PacBio and

327     Illumina libraries were used to construct the new and old assemblies, the significant

328     improvement observed in the new assembly is attributed to the use of a different hybrid assembly

329     approach, combined with the Dovetail Omni-C library preparation and scaffolding with the

330     HiRise (v2.2.0) software.

331

332     ***Pocillopora* genome assemblies**

333     The *Poc. acuta* genome assembly generated in this study (hereinafter the "Hawaiian *Poc. acuta*")

334     is larger (408 Mbp) than *Poc. acuta* from Indonesia (352 Mbp) [22] (Supplementary Table S3)

335     and then its estimated genome size of 353 Mbp (Fig. 1E). The size of the *Poc. meandrina*

336     genome assembly generated in this study (377 Mbp) is comparable to the published Indonesian

337     *Poc. acuta* (352 Mbp) [22] and *Poc. verrucosa* (381 Mbp) [21] species, but is larger than *Poc.*

338     *damicornis* (234 Mbp) [4] (Supplementary Table S3). Although the latter is likely under-

339 assembled given its smaller size relative to the estimated genome size for that species. Moreover,

340 the estimated genome sizes for these species appears to be around 330-350 Mbp, with the

341 assemblies being 350-380 Mbp in size (excluding the Hawaiian *Poc. acuta* [further discussion

342 below]). This suggests that species in the genus *Pocillopora* have genomes that are ~350 Mbp in

343 size.

344      The Hawaiian *Poc. acuta* isolate that was sequenced is a triploid; the presence of three

345 major peaks in the *k*-mer frequency histogram (at ~17x, ~35, and ~51x) which fit the triploid

346 model implemented by GenomeScope2 (black line Fig. 1E and Supplementary Fig. S1F), and the

347 clear "smudge" (bright yellow region in Supplementary Fig. S1F) of *k*-mer pairs with a coverage

348 of ~3n and a normalized coverage of 1/3, all suggests that the sample is triploid. nQuire also

349 predicts that the *Poc. acuta* is a triploid (Supplementary Table S4), supporting the results of

350 GenomeScope2 and Smudegeplot. For *Poc. meandrina*, GenomeScope2 (Fig. 1D), Smudgeplot

351 (Supplementary Fig. S1E), and nQuire (Supplementary Table S4) all predict that the isolate that

352 was sequenced is a diploid.

353      The BUSCO completeness of the Hawaiian *Poc. acuta* genome is improved for both the

354 Metazoa (96.1%), and Eukaryota (98.5%) datasets (Supplementary Table S3) compared to the

355 Indonesian *Poc. acuta* assembly (89.4% and 91.4%, respectively) and the other *Pocillopora*

356 assemblies (~91-95% and 91-98%, respectively). However, the Hawaiian assembly does have a

357 slightly higher proportion of duplicated BUSCO genes (2.5% and 2.0% in the Metazoa and

358 Eukaryota datasets) compared with some (the Indonesian *Poc. acuta* and *Poc. damicornis*

359 genomes which have <1% in both datasets) but not all (the *Poc. verrucosa* genome which has

360 2.9% and 5.5%, respectively) of the published genomes. This is likely a result of the Hawaiian

361 *Poc. acuta* being a triploid; haplotig removal programs (i.e., HaploMerger2 [34]) are generally

362 designed for use with diploid species, so it is unsurprising that they were unable to fully

363 resolving the assembly given the added complexity associated with resolving assemblies of

364 higher ploidy genomes. Regardless, the Hawaiian *Poc. acuta* assembly is more contiguous (i.e.,

365 higher N50 and fewer scaffolds) then the other *Pocillopora* genomes and is the first assembly

366 generated from a non-diploid coral. The *Poc. meandrina* genome has a BUSCO completeness

367 (96.1% for the Metazoa and 98.8% for the Eukaryota datasets) that is just as high as the

368 Hawaiian *Poc. acuta* genome, but with fewer duplicated BUSCO genes (1.2% and 0.4%,

369 respectively), suggesting that this assembly has minimal retained haplotigs.

370

### *Porites compressa* genome assembly

The size of the *Por. compressa* genome assembly generated in this study (593 Mbp) is similar to the published *Por. australiensis* (576 Mbp) [24] and *Por. lutea* (552 Mbp) [23] genomes, and a little smaller than *Por. astreoides* (677 Mbp). The estimated genome sizes for these species appears to be around 525-550 Mbp (excluding *Por. astreoides*, *Por. lutea* and *Por. rus*), with the assemblies coming in at around 550-600 Mbp. The high number of duplicated BUSCO genes in the *Por. astreoides* assembly (11.5% and 14.9% for the Metazoa and Eukaryota datasets, respectively) suggests that its larger assembly size (compared with the other *Porites* species) is likely explained by retained haplotigs. The genome assembly (470 Mbp) and estimated genome size (405 Mbp) of *Por. rus* is smaller than the other *Porites* isolates however, these data were generated from holobiont samples (*i.e.,* samples with both coral, algal symbiont, and associated bacteria DNA present) using a metagenomic binning strategy. The difference in this approach compared with how the other *Porites* genomes were processed likely explain the difference between the sizes. *Por. lutea* has an estimated genome size of 694 Mbp, which is significantly larger than the other *Porites* species and its assembled genome. Whereas this suggests that the *Por. lutea* genome is under-assembled (comprising only ~80% of the estimated genome) its relatively high completeness (95.3% and 98.5% for the Metazoa and Eukaryota datasets, respectively) suggests that the genome size has been overestimated, possibly driven by sequencing error or other factors associated with sample preparation or collection from the field. These results indicate that species in the genus *Porites* have genomes that are just under 600 Mbp in size. For *Por. compressa*, GenomeScope2 (Fig. 1F), Smudgeplot (Supplementary Fig. S1I), and nQuire (Supplementary Table S4) all predict that the isolate sequenced is a diploid.

The BUSCO completeness of the *Por. compressa* assembly is slightly higher (95.5% for the Metazoa and 99.2% for the Eukaryota datasets; Supplementary Table S3) compared to the *Por. astreoides* (93.2% and 98.0%, respectively), *Por. australiensis* (91.6% and 94.9%, respectively), *Por. lutea* (95.3% and 98.5%, respectively), and *Por. rus* (69.6% and 67.1%, respectively) assemblies, but has a much higher N50 (4 Mbp) compared to the published species (0.41, 0.55, 0.66, and 0.14 Mbp, respectively) and fewer scaffolds (608 vs. 3,051, 4,983, 2,975, and 14,982, respectively). The published genome assemblies also have many more gaps (~0-29% of assembled bases are 'N' characters) compared to *Por. compressa* (0%), demonstrating that the

14

401 new assembly is of equally high completeness compared to the published species, but with a

402 much higher contiguity.

403

404 **Predicted protein-coding genes**

405 For *M. capitata*, 54,384 protein-coding genes were predicted in the new assembly compared with

406 the 63,227 predicted in the old version (Supplementary Table S3). The reduction in the number

407 of predicted genes in the new *M. capitata* assembly, compared with the published version, is

408 likely driven by its reduced assembly size, with many of the missing genes likely arising from

409 haplotigs retained in the old assembly, that were removed in the new version. The BUSCO

410 completeness of the predicted genes is improved in the new assembly (95.2% of the Metazoa and

411 96.5% for the Eukaryota BUSCO datasets) compared with the old assembly (94.0% and 93.3%,

412 respectively), and the number of duplicated BUSCO genes is reduced in the new assembly (2.3%

413 and 1.2%, respectively) compared to the published (18.2% and 18.8%, respectively). The

414 predicted gene set from the new Hawaiian *M. capitata* assembly also has > 4.2% and > 3.5%

415 more complete BUSCO genes (from the Metazoa and Eukaryota datasets, respectively)

416 recovered compared to the other published isolates, demonstrating that the gene models

417 predicted in the new assembly are also highly complete. Whereas increase in the number of

418 genes predicted in the new Hawaiian *M. capitata* genome, compared with the published species,

419 could be attributed to differences in the workflows used to predicted the genes in these species

420 [53] however, it is also likely driven by the higher completeness and contiguity of the new

421 genome assembly.

422 There are 33,730 predicted protein-coding genes in the Hawaiian *Poc. acuta* and 31,840

423 in the *Poc. meandrina* genome assemblies (Supplementary Table S3), which is ~4,000–8,000

424 more than predicted in the other *Pocillopora* species. The number of complete BUSCO genes

425 from the Metazoa and Eukaryota BUSCO datasets is > 6% higher in the new Hawaiian *Poc.*

426 *acuta* and *Poc. meandrina* species then in the other *Pocillopora* species; the Hawaiian *Poc. acuta*

427 also has 29.6% and 31.3% (respectively) more complete BUSCO genes recovered than the

428 Indonesian *Poc. acuta*. The number of duplicated BUSCO genes is > 0.7% and > 2.3%

429 (respectively) higher in the Hawaiian *Poc. acuta* gene set compared with the published

430 *Pocillopora* species however, this was expected given the increased size of the genome

431 assembly. The proportion of fragmented BUSCO genes is > 0.9% and > 2% lower (Metazoa and

432    Eukaryota BUSCO datasets, respectively) lower in the Hawaiian *Pocillopora* species compared

433    with the published species. The average transcript length and the number of CDS per transcript

434    of the Hawaiian *Pocillopora* genes (~1,350 bp and ~6.6, respectively) are congruent with the

435    predicted genes of the published *Pocillopora* species (~1,100–1,900 bp and ~5.5-7.5,

436    respectively). This suggests that the higher number of predicted genes in the Hawaiian

437    *Pocillopora* species is not caused by the presence haplotigs in the genome assembly, although

438    this likely contributes to the slights higher number of duplicated BUSCO genes in the Hawaiian

439    *Poc. acuta*, or by the presence of fragmented genes models, since the number of fragmented

440    BUSCO genes and the gene statistics suggest that the majority of genes are full length.

441    Therefore, the higher number of predicted genes in this species can be (at least partially)

442    attributed to the more complete and contiguous genome assemblies of the Hawaiian *Pocillopora*

443    species relative to published species.

444        There are 44,130 predicted protein-coding genes in the Hawaiian *Por. compressa* genome

445    assembly (Supplementary Table S3), which is > 8,000 more genes than predicted in the *Por.*

446    *australiensis* (35,910) and *Por. lutea* (31,126) genomes, 4,677 more than in the *Por. rus* (39,453)

447    genome, and 20,506 less than in the *Por. astreoides* (64,636) genome. The number of complete

448    BUSCO genes from the Metazoa and Eukaryota BUSCO datasets is > 4% higher in *Por.*

449    *compressa* than in the published *Porites* species. The number of duplicated BUSCO genes in

450    *Por. compressa* is similar to *Por. lutea* and *Por. rus* but lower than in *Por. astreoides* and *Por.*

451    *australiensis*, and the number of fragmented BUSCO genes in *Por. compressa* is much lower (>

452    1.9% and > 5.1%, respectively) than in the published species. As with the previous Hawaiian

453    genomes, we attribute the higher number of predicted genes in this species to a more complete

454    and contiguous assembly, relative to the published data.

455

456    **Similarity between *Montipora capitata* chromosomal and extra-chromosomal scaffolds**

457    There are 1,685 scaffolds (totaling ~101 Mbp) in the new *M. capitata* assembly that were not

458    placed into the 14 putative chromosomes by the scaffolding software. Given that the size of the

459    14 chromosomal sequences totals ~680 Mbp, which is close to the estimated genome size of 644

460    Mbp, it is possible that the extra-chromosomal sequences represent retained haplotigs. To

461    explore this issue, we compared the predicted genes in the extra-chromosomal (6,545 protein-

462    coding genes) and chromosomal (47,839) scaffolds to determine how similar the protein content

16

463    is between the two sets of scaffolds and to see if the extra-chromosomal proteins tend to be

464    contained within a single chromosome (suggesting that the extra-chromosomal sequences are

465    likely retained haplotigs). Out of the 6,546 proteins encoded in the extra-chromosomal scaffolds,

466    3,896 (59.53%) have hits to chromosomal proteins with > 75% query coverage and > 75%

467    identity, and 1,623 (24.80%) have hits to chromosomal proteins with > 95% query coverage and

468    > 95% identity. This suggests that whereas the two sets of scaffolds encode many similar

469    (although not identical) proteins, the protein inventory of the extra chromosomal scaffolds only

470    partially overlaps with the gene inventory of the chromosomal scaffolds (we would expect them

471    to have a high level of overlap if they were haplotigs). Furthermore, the extra-chromosomal

472    scaffolds encode 12% of the total predicted genes but, when analyzed separately using BUSCO,

473    have only 1.9% of the Metazoa and 1.6% of the Eukaryota BUSCO genes recovered. This

474    conflict between the number of predicted genes in the scaffolds and the number of BUSCO

475    genes suggests that these scaffolds cannot be easily explained as simply unresolved haplotigs.

476    Finally, of the 3,896 proteins with top hits in the leniently filtered dataset (hit with > 75% query

477    coverage and > 75% identity), 2,748 (70.53%) were on scaffolds with other proteins with top hits

478    to different chromosomes. This suggests that the extra-chromosomal scaffolds have significant

479    structural differences when compared to the chromosomes. These results suggest that the extra-

480    chromosomal scaffolds do not comprise retained haplotigs however, given their significant size,

481    which increases the assembly size well above the estimated size, additional analyses will need to

482    be done to determine the placement of these sequences in the chromosomes and the genes that

483    they encode.

484

485    **Potential implications**

486    The substantial improvement in the contiguity and completeness of the assemblies and predicted

487    genes from the Hawaiian *M. capitata*, *Poc. meandrina*, *Poc. acuta*, and *Por. compressa* species

488    will enable many follow-up studies. The chromosome-level assembly of the *M. capitata* isolate

489    will not only serve as a key reference genome for future population studies focusing on this

490    species in Hawaii, but it will also enable more detailed studies on genome content (such as

491    repeats), gene content, and gene synteny with other species from reefs across the world. The *Poc.*

492    *acuta* genome, although not at chromosome-level resolution, is the most complete available for

493    this genus and will be a valuable model for not only comparative analysis, but for analysis of

494  ploidy in corals. As the first assembly ever generated from a non-diploid coral, this data will

495  open up new questions surrounding the role of ploidy in coral evolution and adaptation and how

496  this phenomenon is involved in the lifecycle of this species and potentially other *Pocillopora*

497  species, both in Hawaiʻi and other reefs across the world. These questions are critical, because an

498  understanding of how changes in ploidy evolve in these corals, particularly in response to stress,

499  will help us model the response of these ecosystems to anthropogenic climate change, and may

500  even provide a new avenue of research for the development of stress resistant "super" corals.

501

## Data availability

503  The Omni-C data generated from the Hawaiian *M. capitata* is available from the NCBI SRA

504  database, under BioProject PRJNA509219. The PacBio HiFi data generated from the Hawaiian

505  *Poc. meandrina* and *Poc. acuta* are also available from the SRA database under BioProject

506  PRJNA761443. The PacBio HiFi and Illumina data generated from the Hawaiian *Por. compressa*

507  are available from the SRA database under BioProject PRJNA663761. The genome assemblies

508  and predicted genes for the Hawaiian *M. capitata* is available from

509  http://cyanophora.rutgers.edu/montipora/ (Version 3), for *Poc. acuta* from

510  http://cyanophora.rutgers.edu/Pocillopora_acuta/ (Version 2), *Poc. meandrina* from

511  http://cyanophora.rutgers.edu/Pocillopora_meandrina/ (Version 1), *Por. compressa* from

512  http://cyanophora.rutgers.edu/Porites_compressa/ (Version 1). The data from the other

513  *Montipora*, *Pocillopora*, and *Porites* species used in this study are available from their respective

514  repositories listed in Supplementary Table S2.

515

## Additional Files

517  **Supplementary Figure S1.** GenomeScope2 (left) and Smudgeplot (right) profiles for (**A**)

518  Hawaiian *M. capitata* (this study), (**B**) Waiopae tide pools *M. capitata*, (**C**) *M. cactus*, (**D**) *M.*

519  *efflorescens*, (**E**) *Poc. meandrina* (this study), (**F**) Hawaiian *Poc. acuta* (this study), (**G**)

520  Indonesian *Poc. acuta*, (**H**) *Poc. verrucose*, (**I**) *Por. compressa* (this study), (**J**) *Por.*

521  *australiensis*, and (**K**) *Por. lutea*. The profiles were computed for each species using 21-mers

522  generated from the trimmed short-read data listed in Supplementary Table S2.

523

## Abbreviations

525 bp: base pairs

526 BUSCO: Benchmarking Universal Single-Copy Orthologs

527 Gbp: gigabase pairs

528 HM2: HaploMerger2

529 Kbp: Kilobase pairs

530 Mbp: megabase pairs

531 NCBI: National Center for Biotechnology Information

532 PacBio: Pacific BioSciences

533 SRA: Sequencing Read Archive

534

## Conflict of Interests

536 The authors declare that they have no other competing interests.

537

## Funding

552

## Author contributions

554 DB conceived the project with HMP and JL. TGS, JML, and YJJ did the bioinformatic analyses,

555 HSY provided sequencing resources, and HMP led the coral sample collection and processing

556     with EM. TGS wrote the manuscript draft with JML, and all authors commented on and

557     approved the submitted version.

558

## Acknowledgements

562

## References

564     1.      van Oppen MJH, Koolmees EM and Veron JEN. Patterns of evolution in the scleractinian
565             coral genus *Montipora* (Acroporidae). Marine Biology. 2004;144 1:9-18.
566             doi:10.1007/s00227-003-1188-3.
567     2.      Forsman ZH, Concepcion GT, Haverkort RD, Shaw RW, Maragos JE and Toonen RJ.
568             Ecomorph or endangered coral? DNA and microstructure reveal hawaiian species
569             complexes: *Montipora dilatata/flabellata/turgescens* & *M. patula/verrilli*. PLoS One.
570             2010;5 12:e15021. doi:10.1371/journal.pone.0015021.
571     3.      Schmidt-Roach S, Miller KJ, Lundgren P and Andreakis N. With eyes wide open: A
572             revision of species within and closely related to the *Pocillopora damicornis* species
573             complex (Scleractinia; Pocilloporidae) using morphology and genetics. Zoological
574             Journal of the Linnean Society. 2014;170 1:1-33. doi:doi.org/10.1111/zoj.12092.
575     4.      Cunning R, Bay RA, Gillette P, Baker AC and Traylor-Knowles N. Comparative analysis
576             of the *Pocillopora damicornis* genome highlights role of immune system in coral
577             evolution. Sci Rep. 2018;8 1:16134. doi:10.1038/s41598-018-34459-8.
578     5.      Williams A, Pathmanathan JS, Stephens TG, Su X, Chiles EN, Conetta D, et al. Multi-
579             omic characterization of the thermal stress phenome in the stony coral *Montipora
580             capitata*. PeerJ. 2021;9:e12335. doi:10.7717/peerj.12335.
581     6.      Mayfield AB, Chen YJ, Lu CY and Chen CS. The proteomic response of the reef coral
582             *Pocillopora acuta* to experimentally elevated temperatures. PLoS One. 2018;13
583             1:e0192001. doi:10.1371/journal.pone.0192001.
584     7.      Henley EM, Quinn M, Bouwmeester J, Daly J, Zuchowicz N, Lager C, et al.
585             Reproductive plasticity of Hawaiian *Montipora* corals following thermal stress. Sci Rep.
586             2021;11 1:12525. doi:10.1038/s41598-021-91030-8.
587     8.      Putnam HM, Davidson JM and Gates RD. Ocean acidification influences host DNA
588             methylation and phenotypic plasticity in environmentally susceptible corals. Evol Appl.
589             2016;9 9:1165-78. doi:10.1111/eva.12408.
590     9.      Jury CP, Delano MN and Toonen RJ. High heritability of coral calcification rates and
591             evolutionary potential under ocean acidification. Sci Rep. 2019;9 1:20419.
592             doi:10.1038/s41598-019-56313-1.
593     10.     Padilla-Gamino JL, Pochon X, Bird C, Concepcion GT and Gates RD. From parent to
594             gamete: vertical transmission of *Symbiodinium* (Dinophyceae) ITS2 sequence
595             assemblages in the reef building coral *Montipora capitata*. PLoS One. 2012;7 6:e38440.
596             doi:10.1371/journal.pone.0038440.

597    11.    Damjanovic K, Menendez P, Blackall LL and van Oppen MJH. Mixed-mode bacterial
598           transmission in the common brooding coral *Pocillopora acuta*. Environ Microbiol.
599           2020;22 1:397-412. doi:10.1111/1462-2920.14856.

600    12.    Cunha RL, Forsman ZH, Belderok R, Knapp ISS, Castilho R and Toonen RJ. Rare coral
601           under the genomic microscope: timing and relationships among Hawaiian *Montipora*.
602           BMC Evol Biol. 2019;19 1:153. doi:10.1186/s12862-019-1476-2.

603    13.    Johnston EC, Forsman ZH, Flot JF, Schmidt-Roach S, Pinzon JH, Knapp ISS, et al. A
604           genomic glance through the fog of plasticity and diversification in *Pocillopora*. Sci Rep.
605           2017;7 1:5991. doi:10.1038/s41598-017-06085-3.

606    14.    Aurelle D, Pratlong M, Oury N, Haguenauer A, Gélin P, Magalon H, et al. Population
607           genomics of *Pocillopora* corals: insights from RAD-sequencing. 2021-10-12 2021.

608    15.    Caruso C, de Souza MR, Ruiz-Jones L, Conetta D, Hancock J, Hobbs C, et al. Genetic
609           patterns in *Montipora capitata* across an environmental mosaic in Kāne'ohe Bay.
610           bioRxiv. 2021:2021.10.07.463582. doi:10.1101/2021.10.07.463582.

611    16.    Shumaker A, Putnam HM, Qiu H, Price DC, Zelzion E, Harel A, et al. Genome analysis
612           of the rice coral *Montipora capitata*. Sci Rep. 2019;9 1:2571. doi:10.1038/s41598-019-
613           39274-3.

614    17.    Shinzato C, Khalturin K, Inoue J, Zayasu Y, Kanda M, Kawamitsu M, et al. Eighteen
615           coral genomes reveal the evolutionary origin of *Acropora* strategies to accommodate
616           environmental changes. Molecular Biology and Evolution. 2021;38 1:16-30.
617           doi:10.1093/molbev/msaa216.

618    18.    Helmkampf M, Bellinger MR, Geib S, Sim SB and Takabayashi M. Draft genome of the
619           rice coral *Montipora capitata* obtained from linked-read sequencing. Genome Biol Evol.
620           2019;11 7:2045-54. doi:10.1093/gbe/evz135.

621    19.    Yuki Y, Go S, Yuna Z, Hiroshi Y and Chuya S. Comparative genomics highlight the
622           importance of lineage-specific gene families in evolutionary divergence of the coral
623           genus, *Montipora*. BMC Ecology and Evolution. 2021;  doi:10.21203/rs.3.rs-944849/v1.

624    20.    Parra G, Bradnam K and Korf I. CEGMA: A pipeline to accurately annotate core genes
625           in eukaryotic genomes. Bioinformatics. 2007;23 9:1061-7.
626           doi:10.1093/bioinformatics/btm071.

627    21.    Buitrago-Lopez C, Mariappan KG, Cardenas A, Gegner HM and Voolstra CR. The
628           genome of the cauliflower coral *Pocillopora verrucosa*. Genome Biol Evol. 2020;12
629           10:1911-7. doi:10.1093/gbe/evaa184.

630    22.    Vidal-Dupiol J, Chaparro C, Pratlong M, Pontarotti P, Grunau C and Mitta G.
631           Sequencing, *de novo* assembly and annotation of the genome of the scleractinian coral,
632           *Pocillopora acuta*. bioRxiv. 2020:698688. doi:10.1101/698688.

633    23.    Robbins SJ, Singleton CM, Chan CX, Messer LF, Geers AU, Ying H, et al. A genomic
634           view of the reef-building coral *Porites lutea* and its microbial symbionts. Nat Microbiol.
635           2019;4 12:2090-100. doi:10.1038/s41564-019-0532-4.

636    24.    Shinzato C, Takeuchi T, Yoshioka Y, Tada I, Kanda M, Broussard C, et al. Whole-
637           genome sequencing highlights conservative genomic strategies of a stress-tolerant, long-
638           lived scleractinian coral, *Porites australiensis* Vaughan, 1918. Genome Biol Evol.
639           2021;13 12 doi:10.1093/gbe/evab270.

640    25.    Wong KH and Putnam HM. The Genome of the Mustard Hill Coral, *Porites astreoides*.
641           GIGAbyte (Under Review). 2022.

642 26. Kenyon JC. Models of reticulate evolution in the coral genus *Acropora* based on
643     chromosome numbers: Parallels with plants. Evolution. 1997;51 3:756-67.
644     doi:10.1111/j.1558-5646.1997.tb03659.x.
645 27. Nurk S, Walenz BP, Rhie A, Vollger MR, Logsdon GA, Grothe R, et al. HiCanu:
646     accurate assembly of segmental duplications, satellites, and allelic variants from high-
647     fidelity long reads. Genome Res. 2020;30 9:1291-305. doi:10.1101/gr.263566.120.
648 28. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an
649     integrated tool for comprehensive microbial variant detection and genome assembly
650     improvement. PLoS One. 2014;9 11:e112963. doi:10.1371/journal.pone.0112963.
651 29. Bolger AM, Lohse M and Usadel B. Trimmomatic: A flexible trimmer for Illumina
652     sequence data. Bioinformatics. 2014;30 15:2114-20. doi:10.1093/bioinformatics/btu170.
653 30. Marcais G and Kingsford C. A fast, lock-free approach for efficient parallel counting of
654     occurrences of *k*-mers. Bioinformatics. 2011;27 6:764-70.
655     doi:10.1093/bioinformatics/btr011.
656 31. Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids
657     Res. 1999;27 2:573-80. doi:10.1093/nar/27.2.573.
658 32. Bao Z and Eddy SR. Automated de novo identification of repeat sequence families in
659     sequenced genomes. Genome Res. 2002;12 8:1269-76. doi:10.1101/gr.88502.
660 33. Price AL, Jones NC and Pevzner PA. De novo identification of repeat families in large
661     genomes. Bioinformatics. 2005;21 Suppl 1:i351-8. doi:10.1093/bioinformatics/bti1018.
662 34. Huang S, Kang M and Xu A. HaploMerger2: Rebuilding both haploid sub-assemblies
663     from high-heterozygosity diploid genome assembly. Bioinformatics. 2017;33 16:2577-9.
664     doi:10.1093/bioinformatics/btx220.
665 35. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO:
666     assessing genome assembly and annotation completeness with single-copy orthologs.
667     Bioinformatics. 2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.
668 36. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. *De novo*
669     transcript sequence reconstruction from RNA-seq using the Trinity platform for reference
670     generation and analysis. Nat Protoc. 2013;8 8:1494-512. doi:10.1038/nprot.2013.084.
671 37. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length
672     transcriptome assembly from RNA-Seq data without a reference genome. Nat Biotechnol.
673     2011;29 7:644-52. doi:10.1038/nbt.1883.
674 38. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: Ultrafast
675     universal RNA-seq aligner. Bioinformatics. 2013;29 1:15-21.
676     doi:10.1093/bioinformatics/bts635.
677 39. Bruna T, Hoff KJ, Lomsadze A, Stanke M and Borodovsky M. BRAKER2: Automatic
678     eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a
679     protein database. NAR Genom Bioinform. 2021;3 1:lqaa108.
680     doi:10.1093/nargab/lqaa108.
681 40. Lomsadze A, Burns PD and Borodovsky M. Integration of mapped RNA-Seq reads into
682     automatic training of eukaryotic gene finding algorithm. Nucleic Acids Res. 2014;42
683     15:e119. doi:10.1093/nar/gku557.
684 41. Stanke M, Keller O, Gunduz I, Hayes A, Waack S and Morgenstern B. AUGUSTUS: *Ab*
685     *initio* prediction of alternative transcripts. Nucleic Acids Res. 2006;34 Web Server
686     issue:W435-9. doi:10.1093/nar/gkl200.

687   42.   Marchler-Bauer A, Bo Y, Han L, He J, Lanczycki CJ, Lu S, et al. CDD/SPARCLE:
688         functional classification of proteins via subfamily domain architectures. Nucleic Acids
689         Res. 2017;45 D1:D200-D3. doi:10.1093/nar/gkw1129.
690   43.   Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, von Mering C, et al.
691         Fast Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-
692         Mapper. Molecular Biology and Evolution. 2017;34 8:2115-22.
693         doi:10.1093/molbev/msx148.
694   44.   Moriya Y, Itoh M, Okuda S, Yoshizawa AC and Kanehisa M. KAAS: an automatic
695         genome annotation and pathway reconstruction server. Nucleic Acids Research.
696         2007;35:W182-W5. doi:10.1093/nar/gkm321.
697   45.   Celis JS, Wibberg D, Ramirez-Portilla C, Rupp O, Sczyrba A, Winkler A, et al. Binning
698         enables efficient host genome reconstruction in cnidarian holobionts. Gigascience.
699         2018;7 7 doi:10.1093/gigascience/giy075.
700   46.   Pertea G and Pertea M. GFF Utilities: GffRead and GffCompare [version 1; peer review:
701         2 approved]. F1000Research. 2020;9 304 doi:10.12688/f1000research.23297.1.
702   47.   Ranallo-Benavidez TR, Jaron KS and Schatz MC. GenomeScope 2.0 and Smudgeplot for
703         reference-free profiling of polyploid genomes. Nat Commun. 2020;11 1:1432.
704         doi:10.1038/s41467-020-14998-3.
705   48.   Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
706         EMBnetjournal. 2011;17 1:3. doi:10.14806/ej.17.1.200.
707   49.   Weiss CL, Pais M, Cano LM, Kamoun S and Burbano HA. nQuire: A statistical
708         framework for ploidy estimation using next generation sequencing. BMC Bioinformatics.
709         2018;19 1:122. doi:10.1186/s12859-018-2128-z.
710   50.   Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years
711         of SAMtools and BCFtools. Gigascience. 2021;10 2 doi:10.1093/gigascience/giab008.
712   51.   Manni M, Berkeley MR, Seppey M, Simão FA and Zdobnov EM. BUSCO Update:
713         Novel and streamlined workflows along with broader and deeper phylogenetic coverage
714         for scoring of eukaryotic, prokaryotic, and viral genomes. Molecular Biology and
715         Evolution. 2021;  doi:10.1093/molbev/msab199.
716   52.   Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
717         Architecture and applications. BMC Bioinformatics. 2009;10:421. doi:10.1186/1471-
718         2105-10-421.
719   53.   Chen YB, Gonzalez-Pech RA, Stephens TG, Bhattacharya D and Chan CX. Evidence
720         that inconsistent gene prediction can mislead analysis of dinoflagellate genomes. J
721         Phycol. 2020;56 1:6-10. doi:10.1111/jpy.12947.

722

## Tables

**Table S1:** Summary of coral assemblies before and after haplotype merging.

**Table S2:** Metadata for the genome and gene models downloaded for the coral species used for comparative analysis.

**Table S3:** Comparison between the published *Montipora*, *Pocillopora*, and *Porites* genomes and those generated in this study. All statistics were calculated in this study using the available genome and gene models.

**Table S4**: Results from nQuire lrdmodel ploidy estimation for the Hawaiian coral genomes analyzed in this study.

**Figure Legends**

**Figure 1:** (**A**) Cumulative and (**B**) individual length of scaffolds in the new Hawaiian *M.*

*capitata* genome assembly. Scaffolds were sorted by length in descending order; each point

along the x-axis of (**A**) and (**B**) represents a scaffold, with the longest scaffold being the first and

the shortest being the last on the x-axis of each plot. In (**A**) and (**B**) a zoomed-in section of the

larger plot is shown on the right highlighting the 40 largest scaffolds; a horizontal red line in (**A**)

shows the total assembled bases in the new genome and a vertical dashed line in (**A**) and (**B**) is

positioned after the 14th largest scaffold. GenomeScape2 linear *k*-mer distributions of the

Hawaiian (**C**) *M. capitata*, (**D**) *Poc. meandrina*, (**E**) *Poc. acuta*, and (**F**) *Por. compressa* species

with theoretical diploid (or triploid for *Poc. acuta*) models shown by the black lines. The

GenomeScape2 profiles were computed for each species using 21-mers generated from the

trimmed short-read data listed in Supplementary Table S2.

**Figure 2:** Results from BUSCO analysis run using the genomes and predicted genes from all

published (including this study) *Montipora*, *Pocillopora*, and *Porites* species, plus the old

version of the *M. capitata* genome that our group published in 2019 [16]. BUSCO results for

each species using the (**A**) Metazoa dataset (genome mode), (**B**) Eukaryota dataset (genome

mode), (**C**) Metazoa dataset (protein mode), and (**D**) Eukaryota dataset (protein mode).

Figure 1

Figure 1

Figure 2

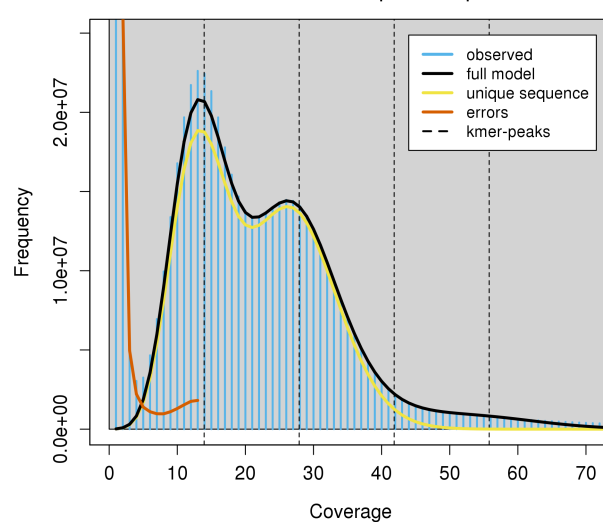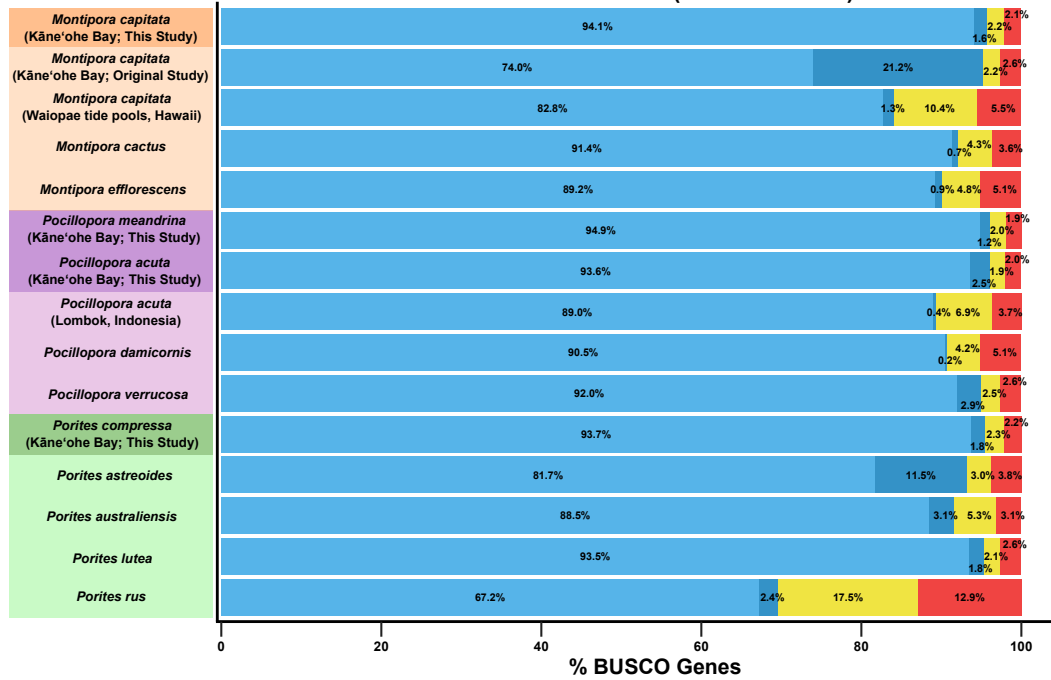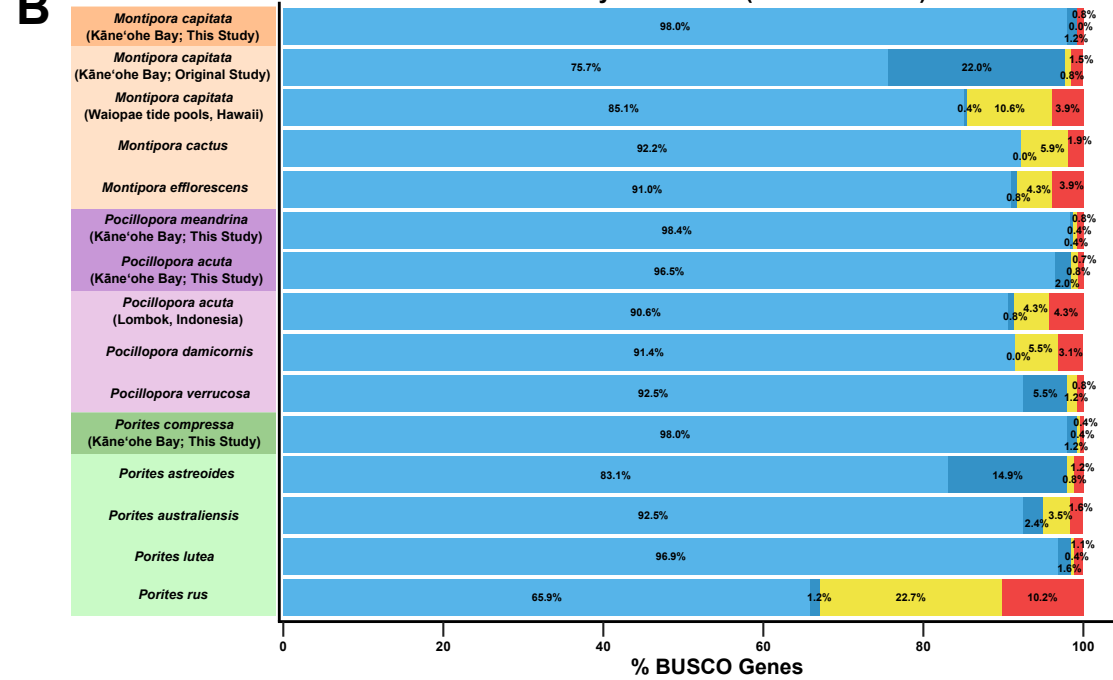**A** — BUSCO Metazoa dataset (Genome mode)

**B** — BUSCO Eukaryota dataset (Genome mode)

**C** — BUSCO Metazoa dataset (Protein mode)

**D** — BUSCO Eukaryota dataset (Protein mode)

Legend: Complete (C) and single-copy (S); Complete (C) and duplicated (D); Fragmented (F); Missing (M)

Figure S1

Click here to access/download
**Supplementary Material**
Figure_S1.pdf

Click here to access/download
**Supplementary Material**
Supplementary_Tables.xlsx