

Manuscript Number:	GIGA-D-18-00174	
Full Title:	The genome of common long-arm octopus <i>Octopus minor</i>	
Article Type:	Data Note	
Funding Information:	MABIK (2018M00900)	Dr. Hye Suck An
Abstract:	<p>Background: The common long-arm octopus (<i>Octopus minor</i>) is found in mudflats of subtidal zones and faces numerous environmental challenges. The ability to adapt its morphology and behavioural repertoire to diverse environmental conditions makes the species a promising model to understand genomic adaptation and evolution in cephalopods. Findings: The final genome assembly of <i>O. minor</i> is 5.09 Gb, with a contig N50 size of 197 kb and longest size of 3.027 Mb, from a total of 419 Gb raw reads generated using PacBio RS II platform. We identified 30,010 genes and 44.43% of the genome is composed of repeat elements. The genome-wide phylogenetic tree indicated the divergence time between <i>O. minor</i> and <i>O. bimaculoides</i> was estimated to be 43 million years ago (Mya) based on single-copy orthologous genes. In total, 178 gene families are expanded in <i>O. minor</i> in the 14 bilaterian animal species. Conclusion: We found that the <i>O. minor</i> genome was larger than that of closely related <i>O. bimaculoides</i>, and this difference could be explained by enlarged introns and recently diversified transposable elements. The high-quality <i>O. minor</i> genome assembly provides a valuable resource for understanding octopus genome evolution and the molecular basis of adaptations to mudflats.</p>	
Corresponding Author:	Hyun Park KOREA, REPUBLIC OF	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Bo-Mi Kim	
First Author Secondary Information:		
Order of Authors:	Bo-Mi Kim	
	Seunghyun Kang	
	Do-Hwan Ahn	
	Seung-Hyun Jung	
	Hwanseok Rhee	
	Jong Su Yoo	
	Jong-Eun Lee	
	SeungJae Lee	
	Yong-Hee Han	
	Kyoung-Bin Ryu	
	Sung-Jin Cho	
	Hyun Park	
	Hye Suck An	

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

The genome of common long-arm octopus *Octopus minor*

Bo-Mi Kim ^{a,†}, Seunghyun Kang ^{a,†}, Do-Hwan Ahn ^{a,†}, Seung-Hyun Jung ^{b,†}, Hwanseok Rhee ^{c,†}, Jong Su Yoo ^b, Jong-Eun Lee ^c, SeungJae Lee ^c, Yong-Hee Han ^d, Kyoung-Bin Ryu ^d, Sung-Jin Cho ^{d,*}, Hyun Park ^{a,e,*}, Hye Suck An ^{b,*}

Affiliations

^a Unit of Polar Genomics, Korea Polar Research Institute(KOPRI), Incheon 21990, Korea

^b Department of Genetic Resources Research, National Marine Biodiversity Institute of Korea (MABIK), Janghang-eup, Seochun-gun, Chungchungnam-do 33662, Korea

^c Genomics Lab, Cluster Center, DNA Link, Inc., 150, Bugahyeon-ro, Seodaemun-gu, Seoul 03759, Korea

^d School of Biological Sciences, College of Natural Sciences, Chungbuk National University, Cheongju, Chungbuk 28644, Korea

^e Polar Sciences, University of Science & Technology, Yuseong-gu, Daejeon 34113, Korea

*Co-corresponding author:

School of Biological Sciences, College of Natural Sciences, Chungbuk National University, Cheongju, Chungbuk 28644, Korea; E-mail address: sjchobio@chungbuk.ac.kr (S. Cho)

Unit of Polar Genomics, Korea Polar Research Institute, Incheon 21990, Korea; E-mail address: hpark@kopri.re.kr (H. Park)

Department of Genetic Resources Research, National Marine Biodiversity Institute of Korea (MABIK), Janghang-eup, Seochun-gun, Chungchungnam-do 33662, Korea; E-mail address: mgran@mabik.re.kr H.S. An)

[†]These authors contributed equally to this work.

1 **Abstract**

2
3 **Background:** The common long-arm octopus (*Octopus minor*) is found in mudflats of subtidal
4 zones and faces numerous environmental challenges. The ability to adapt its morphology and
5 behavioural repertoire to diverse environmental conditions makes the species a promising
6 model to understand genomic adaptation and evolution in cephalopods. **Findings:** The final
7 genome assembly of *O. minor* is 5.09 Gb, with a contig N50 size of 197 kb and longest size of
8 3.027 Mb, from a total of 419 Gb raw reads generated using PacBio RS II platform. We
9 identified 30,010 genes and 44.43% of the genome is composed of repeat elements. The
10 genome-wide phylogenetic tree indicated the divergence time between *O. minor* and *O.*
11 *bimaculoides* was estimated to be 43 million years ago (Mya) based on single-copy orthologous
12 genes. In total, 178 gene families are expanded in *O. minor* in the 14 bilaterian animal species.
13 **Conclusion:** We found that the *O. minor* genome was larger than that of closely related *O.*
14 *bimaculoides*, and this difference could be explained by enlarged introns and recently
15 diversified transposable elements. The high-quality *O. minor* genome assembly provides a
16 valuable resource for understanding octopus genome evolution and the molecular basis of
17 adaptations to mudflats.

18
19 **Key words:**

20 Octopus genome, Cephalopods, adaptation and evolution, long-read sequencing

22 Introduction

23 Cephalopods (*e.g.* cuttlefish, nautilus, octopus, and squid) belong to the phylum Mollusca.
24 As advanced invertebrates, cephalopods have interesting biological characteristics, such as an
25 extraordinary life-history plasticity, rapid growth, short lifespan, large brain, and sophisticated
26 sense organs with a complex nervous system[1]. The ability to adapt their morphology and
27 behavioural repertoire to diverse environmental conditions and capacity for learning and
28 memory are common traits in cephalopods, but have rarely been observed in other
29 invertebrates[2]. Many cephalopod species have been considered for fisheries and are
30 promising candidates for aquaculture. There are an estimated 1,000 cephalopod species (~700
31 known marine-living species), and octopods are among the most well-known representatives
32 of the class, including over 150 species worldwide[3]. Studies have evaluated the biological
33 machinery underlying the fundamental nervous system functions, strong behavioural plasticity,
34 and learning ability in octopods[4, 5].

35 *Octopus minor* (Sasaki, 1920), also known as the common long-arm octopus, is a benthic
36 littoral species, and is a major commercial fishery product with a high annual yield[6]. *O. minor*
37 is relatively small and possesses a shorter life cycle (approximately 1 year), thinner arms, and
38 a lower ratio between head size and arm length compared to those of other octopus species
39 (**Fig. 1a**). The species is widely distributed in Northeast Asia, particularly in coastal regions of
40 South Korea, China, and Japan (**Fig. 1b**). Most *O. minor* habitats are mud and mud sand in
41 well-developed mudflats of coastal regions; they spawn in holes on the mudflat by digging
42 with the whole body. As an important economic cephalopod in South Korea, fishermen
43 normally catch *O. minor* by digging a hole in the mudflat with shovels. Thus, they are subjected
44 to the harsh environmental conditions of mudflats, including diurnal temperature changes,
45 steep salinity and pH gradients, desiccation, wave action and tides, oxygen availability, and
46 interrupted feeding. Owing to the ability of *O. minor* to tolerate environmental fluctuations, it
47 is a promising organism for studies of the molecular basis of plasticity and mechanisms
48 underlying adaptation to harsh environmental conditions, although relevant information is
49 scarce. To make full use of this emerging cephalopod model system and to understand the
50 interesting features of *O. minor*, including its plasticity in mudflats and genetic evolution, a
51 high-quality reference genome is required.

52 The published genome and multiple transcriptomes of the California two-spot octopus
53 *Octopus bimaculoides* have provided valuable information on genomic traits (*e.g.* gene family
54 expansion, genome rearrangements, and transposable element activity) related to the evolution

1 55 of neural complexity and morphological innovations[7]. In this study, we report a high-quality
2
3 56 genome assembly and annotation for *O. minor*. We compare the genomes of *O. minor* and *O.*
4
5 57 *bimaculoides* and provide evidence that the expansion of genes and/or gene families is related
6
7 58 to adaptation to the harsh environmental conditions of mudflats.
8
9 59

10 60 **Genome sequencing and annotation**

11 61 *O. minor* genomic DNA was extracted from leg muscle tissues. The average coverage of
12 62 SMRT sequences was ~76-fold using P6-C4 sequence chemistry from genomic DNA libraries.
13
14 63 The average subread length was 9.2 kb (Supplementary Table S1). For genome size estimation,
15
16 64 a k-mer analysis was performed using Jellyfish[8] with paired-end sequences of the genomic
17
18 65 DNA libraries. The *O. minor* genome was estimated to be 5.1 Gb (Supplementary Figs. S1 and
19
20 66 S2). The *de novo* assembly generated using FALCON-Unzip assembler was 5.09 Gb with
21
22 67 41,584 contigs[9]. Finally, evaluation of the genome completeness was checked using
23
24 68 BUSCO[10] (Supplementary Table S2).

25
26 69 Total RNA was extracted from thirteen tissues (brain, branchial heart, buccal mass, eye,
27
28 70 heart, kidney, liver, ovary, poison gland, siphon, skin, and suckers) using the RNeasy Mini Kit
29
30 71 (Qiagen, Hilden, Germany) according to the manufacturer's instructions. RNA quality was
31
32 72 confirmed using an Agilent Bioanalyzer™. Isoform sequencing was performed using pooled
33
34 73 RNA from thirteen organs. Library construction and sequencing were performed using Pacbio
35
36 74 RS II (Supplementary Table S3). The SMRTbell library for Iso-seq was sequenced using 16
37
38 75 SMRT cells (1–2 kb, three cells; 2–3 kb, six cells; and 3–6 kb, seven cells). Reads were
39
40 76 identified using the SMRT Analysis ver. 2.3 RS_IsoSeq.1 classification protocol. All full-
41
42 77 length reads derived from the same isoform were clustered and consensus sequences were
43
44 78 polished using the TOFU pipeline (isoseq-tofu)[11]. Additionally, chimeras of consensus
45
46 79 sequences were removed.

47 80 MAKER was used for genome annotation[12]. First, repetitive elements were identified
48
49 81 using RepeatMasker[13]. A *de novo* repeat library was constructed using RepeatModeler (ver.
50
51 82 1.0.3)[14], including RECON[15] and RepeatScout[16], with default parameters. Consensus
52
53 83 sequences and classification information for each repeat family were generated, and tandem
54
55 84 repeats, including simple repeats, satellites, and low-complexity repeats, were predicted using
56
57 85 Tandem Repeats Finder[11]. This masked genome sequence was used for *ab initio* gene
58
59 86 prediction with SNAP software[17]; subsequently, alignments of expressed sequence tags with
60
61 87 BLASTn and protein information from tBLASTx were included. The *de novo* repeat library of

1 88 *O. minor* from RepeatModeler was used for RepeatMasker; proteins from sequenced molluscs
2
3 89 (*L. gigantea*, *C. gigas*, and *Aplysia californica*) and an octopus species (*O. bimaculoides*) were
4
5 90 included in the analysis. Transcriptome assembly results were used for expressed sequence
6
7 91 tags. Next, MAKER polished the alignments using Exonerate, which provided integrated
8
9 92 information for SNAP annotation. Using MAKER, the final gene model was selected and
10
11 93 revised considering all information. A total of 30,010 *O. minor* genes were predicted using
12
13 94 MAKER. The Infernal software package (ver. 1.1)[18] and covariance models from the
14
15 95 Rfam[19] database were used to identify other non-coding RNAs in the *O. minor* scaffold.
16
17 96 Putative tRNA genes were identified using tRNAscan-SE[20]. tRNAscan-SE uses a covariance
18
19 97 model that scores candidates based on their sequence and predicted secondary structures.

20 98 The mean size of *O. minor* genes was 23.6 kb, with an average intron length of 5.4 kb (4.2
21
22 99 introns per gene) (Supplementary Table S4). The *O. minor* genome contained 30,010 protein-
23
24 100 coding genes (Table 1), of which 96% were annotated based on known proteins in public
25
26 101 databases, and 79% were similar to *O. bimaculoides* genes (Supplementary Table S5).

27 102

28 103 **Comparative genomic analyses and duplicate genes**

29 104 To resolve gene family evolution in the *O. minor* genome, we classified orthologous gene
30
31 105 clusters (Supplementary Table S6) from 14 species and found evidence for the recent expansion
32
33 106 of low-copy gene duplicates and the expansion of large gene families. Orthologous groups were
34
35 107 identified using both OrthoMCL[21] and Pfam[22] domain assignments. OrthoMCL generated
36
37 108 a graphical representation of sequence relationships, which was then divided into subgraphs
38
39 109 using the Markov Clustering Algorithm (MCL) from multiple eukaryotic genomes[21]. The
40
41 110 standard parameters and options of OrthoMCL were used for all steps, together with the
42
43 111 genomes of 14 species (Supplementary Table S6). For *O. minor*, the coding sequence from the
44
45 112 MAKER annotation pipeline was used. To construct a phylogenetic tree and estimate the
46
47 113 divergence time, 202 1:1:1 single-copy orthologous genes were used. Using the Probabilistic
48
49 114 Alignment Kit (PRANK)[23], protein-coding genes were aligned with the codon alignment
50
51 115 option, and poorly aligned regions with gaps were eliminated using Gblock[24] with a codon
52
53 116 model. A maximum-likelihood tree was built using RAxML[25] with 1,000 bootstrap
54
55 117 replicates, and the divergence time was calibrated using TimeTree[26]. The average gene gain-
56
57 118 loss was identified using CAFÉ 4.0[27].

58 119 Sequence divergence was estimated by calculating d_S values using the yn00 program from
59
60 120 the PAML package[28]. The Jukes–Cantor distances were adjusted using the Jukes–Cantor

1 121 formula $d_{XY} = -(3/4)\ln(1-4/3D)$, where D is the proportion of nucleotide differences between
2
3 122 the sequences. The time estimation was calibrated by assuming d_s of ~1 is 135 million years[7].

4
5 123 Gene family analyses of specific genes of interest were manually curated using manual
6
7 124 gene search methods. Gene or gene family targets identified in the genomes of *Octopus*
8
9 125 *bimaculoides*, *Crassostrea gigas*, *Lottia gigantea*, *Capitella teleta*, and *Homo sapiens* were
10 126 directly mapped to the *O. minor* genome database by a local BLAST analysis. Alignments were
11
12 127 generated using Clustal Omega (ClustalO)[29] and Multiple Sequence Comparison by Log-
13
14 128 Expectation (MUSCLE)[30], and phylogenetic trees were built using FastTree[31] or RAxML
15
16 129 with 1,000 bootstrap replicates.

17 130 A statistical analysis of the changes in gene family sizes indicated significantly greater gene
18
19 131 family expansion in *O. minor* (178 gene families) compared to other species, e.g. interleukin-
20
21 132 17, G protein-coupled receptor (GPCR) proteins, Zinc-finger of C2H2 type, heat shock protein
22
23 133 (HSP) 70 proteins, and cadherin-like domains (Supplementary Tables S7–S9). The divergence
24
25 134 time between *O. minor* and *O. bimaculoides* was estimated to be 43 million years ago (Mya)
26
27 135 based on single-copy orthologous genes (Fig. 2a) Further, Pfam domain and EggNOG
28
29 136 metazoan database searches consistently showed the expansion of gene families, including the
30
31 137 cadherin and protocadherin domains and interleukin-17 (Fig. 2b and Supplementary Tables
32 138 S10 and S11).

33 34 139 35 36 140 **Transposable element annotation and expansions**

37 141 The *O. minor* genome was larger than that of *O. bimaculoides* (2.7 Gb), with a high level
38
39 142 of repetitive sequences (44.43%) (Supplementary Tables S12–S14). Repeats accounted for 44%
40
41 143 (2.262 Gb) of the assembly, and were dominated by simple repeats (14.7%) and TEs, especially
42
43 144 DNA transposons and long interspersed elements (LINEs), which were more abundant in the
44
45 145 *O. minor* genome than in the *O. bimaculoides* genome. In an analysis of genes (i.e. exons and
46
47 146 introns) and intergenic sequences, TEs were highly distributed in the intergenic sequence
48
49 147 regions in both species (Supplementary Fig. S4). In particular, TE accumulation in intergenic
50
51 148 sequence regions was significantly greater in *O. minor* than in *O. bimaculoides*. The larger
52
53 149 gene size and higher repeat content may explain the larger genome of *O. minor* compared with
54 150 *O. bimaculoides*.

55
56 151 TEs are crucial components of animal genomes, with major roles in genome
57
58 152 rearrangements and evolution. Based on the mechanism of transposition, TEs are grouped into
59
60 153 two main classes, class I retrotransposons, which are subdivided into long terminal repeats

1 154 (LTRs) and non-LTR retrotransposons [*e.g.* LINEs and short interspersed elements (SINEs)],
2
3 155 and class II DNA transposons[32]. We detected more TEs in the larger genome of *O. minor*
4
5 156 than in the smaller genome of *O. bimaculoides*. Approximately half of the *O. minor* genome
6
7 157 was composed of TEs (11,547,325 TEs; 44% of the genome), while one-third of the *O.*
8
9 158 *bimaculoides* genome was composed of TEs (3,887,025 TEs; 35%) (Supplementary Table S12).
10 159 The majority of class I retrotransposons in the *O. minor* genome were LINEs (10%), as was
11
12 160 also the case in *O. bimaculoides* (9%), and the proportion of DNA transposons in *O. minor*
13
14 161 (13%) was comparable to that in *O. bimaculoides* (12%). Interestingly, the *O. minor* genome
15
16 162 had fewer SINEs (1,540 copies; 0.01%) and more rolling-circle (RC)-Helitrons (121,101
17
18 163 copies; 3.7%) than the *O. bimaculoides* genome (SINEs: 115,169 copies, 1.8%; RC-Helitron:
19 164 43,735 copies, 0.7%). A Kimura distance analysis revealed that the most frequent TE sequence
20
21 165 divergence relative to the TE consensus sequence was ~7–10%, with an additional peak at 3%
22
23 166 (Fig. 3a), compared to 16–17% in the *O. bimaculoides* genome (Fig. 3b and Supplementary
24
25 167 Table S12).

26 168 A more recent expansion of LINEs, without an increase in SINEs, was detected in the *O.*
27
28 169 *minor* genome, while ancient copies of all four types of TEs and an ancient transposition burst
29
30 170 of DNA transposons were observed in *O. bimaculoides*. Using the recent TE expansion in the
31
32 171 *O. minor* genome, we correlated Jukes–Cantor distance measures with d_s and identified two
33
34 172 unique expansion waves at 0.04 and 0.09 compared to the distribution of *O. bimaculoides* TEs
35
36 173 (Supplementary Figs. S5 and S6). This suggests that a major expansion of TEs in the *O. minor*
37
38 174 genome occurred 11 to 25 Mya, which is after the divergence of *O. minor* and *O. bimaculoides*.
39
40 175

41 176 **Conclusions**

42
43 177 *O. minor* has developed morphological and physiological adaptations to match their unique
44
45 178 mudflat habitats. In summary, we generated a high-quality sequence assembly for *O. minor* to
46
47 179 elucidate the molecular mechanisms underlying their adaptations. In a direct comparison
48
49 180 between the genomes of *O. minor* and *O. bimaculoides*, we discovered that they evolved
50
51 181 recently and independently from the octopus lineage during the successful transition from an
52
53 182 aquatic habitat to mudflats. We also found evidence suggesting that speciation in the genus
54
55 183 *Octopus* is closely related to the gene family expansion associated with environmental
56
57 184 adaptation. Finally, in addition to providing insights into the genome size increase via gene
58
59 185 family expansion, the *O. minor* genome sequence also provides an essential resource for studies
60
61 186 of Cephalopoda evolution.

1 187
2
3 188
4
5 189
6
7 190
8
9 191
10 192
11
12 193
13
14 194
15
16 195
17
18 196
19 197
20
21 198
22
23 199
24
25 200
26
27 201
28
29 202
30
31 203
32 204
33
34 205
35
36 206
37
38 207
39
40 208
41
42 209
43 210
44
45 211
46
47 212
48
49 213
50
51 214
52
53 215
54
55 216
56
57 217
58
59 218
60
61 219
62
63
64
65

Availability of supporting data

The octopus (*O. minor*) genome project was deposited at NCBI under BioProject number PRJNA421033. The whole-genome sequence was deposited in the Sequence Read Archive (SRA) database under accession number SRX3462978, and isoform sequence from PacBio sequencing data were deposited in the SRA database under accession numbers SRX3478495 and SRX3478496. Other supporting data, including annotations, alignments, and BUSCO results, are available in the GigaScience repository, GigaDB [---].

Additional files

Fig. S1. Estimation of genome size of *O. minor* based on distribution of 17 k-mer frequency in raw sequencing reads.

Fig. S2. Genome size determination by flow cytometry. The flow cytometry analysis provides as estimation of Propidium iodide (PI) staining. Accepting a haploid genome size estimate of 2.81 Gb for Mouse (Assembly; GRCm38.p6), we estimate the genome size of *O. minor* to be 5.38 Gb.

Fig. S3. Blast top hit distribution.

Fig. S4. Composition of transposable elements in the regions of gene and intergenic sequence.

Fig. S5. Transposable elements Juke-cantor distance distribution.

Fig. S6. Transposable elements Juke-cantor distance distribution of *O. minor*.

Table S1. Statistics for SMRT sequencing for the *O. minor* genome sequencing.

Table S2. Benchmarking Universal Single-Copy Orthologs (BUSCO) evaluated for the completeness of the *O. minor* genome assembly.

Table S3. Isoform sequencing summary of transcriptome analysis of *O. minor* using PacBio RSII.

Table S4. Brief summary of gene statistics.

Table S5. Functional annotation statistics of transcriptome assembly.

Table S6. Summary of orthologous gene clusters analyzed in 14 species.

Table S7. CAFÉ gene family analysis results.

Table S8. Example of top 30 CAFÉ significantly expanded gene families.

Table S9. Example of top 30 CAFÉ significantly shrunked gene families.

Table S10. Top 30 expanded Pfam domains.

1 220 Table S11. Top 30 expanded EggNOG domains.
2
3 221 Table S12. Statistics of repeat analysis of the *O. minor* genome.
4
5 222 Table S13. Classifications and frequencies of transposable elements and other repeats.
6
7 223 Table S14. Classifications and frequencies of simple repeats.

8 224
9
10 225

11 226 **Acknowledgements**

12 227 We thank Jong Won Han and Ha Yeun Song of the National Marine Biodiversity Institute of
13
14 228 Korea (MABIK) for the sampling of 18 tissues used for transcriptome assembly, as well as
15
16 229 Keekwang Kim of Chungnam National University and Kun-Hee Kim of Chonnam National
17
18 230 University for their devotion to estimate the genome size of *O. minor* by flow cytometry. We
19
20 231 also thank Jeollanam-Do Oceans & Fisheries Science Institute for providing octopus embryos.
21
22 232

23 232

24 233 **Funding**

25 233
26 234 This work was supported by grants (2018M00900) from MABIK.
27
28 235

29 235

30 236 **Competing interests**

31 237 The authors declare that they have no competing interests.
32
33 238

34 238

35 239 **Author contributions**

36 239
37 240 H.S.A., H.P., and J.L. conceived the study. H.P., B.K., S.K., D.A., S.J., J.L., H.R., and S.L.
38
39 241 performed genome sequencing, assembly, and annotation. S.J., Y.H., K.R., and S.C. performed
40
41 242 experiments. J.S.Y., H.S.A., H.P., S.J., and J.L. advised and coordinated the study. B.K., S.K.,
42
43 243 D.A., and H.P. mainly wrote the paper. All authors contributed to writing and editing the
44
45 244 manuscript and supplementary information and producing the figures.
46
47 245

48 245

49 246 **References**

- 50 247 1. Boyle P and Rodhouse P. Cephalopods: ecology and fisheries. Oxford: Blackwe
51 248 ll Science Ltd; 2005.
52
53 249 2. Hanlon RT and Messenger JB. Cephalopod behaviour. Cambridge: Cambridge
54 250 University Press; 1998.
55
56 251 3. Guzik MT, Norman MD and Crozier RH. Molecular phylogeny of the benthic
57
58 252 shallow-water octopuses (Cephalopoda: Octopodinae). Mol Phylogen Evol. 2005;

59
60
61
62
63
64
65

1 253 37 1:235-48.

2

3 254 4. Hochner B, Shomrat T and Fiorito G. The octopus: a model for a comparative
4
5 255 analysis of the evolution of learning and memory mechanisms. *Biol Bull.* 200
6
7 256 6;210 3:308-17.

8 257 5. Mather JA. Cephalopod consciousness: behavioural evidence. *Conscious Cogn.*
9
10 258 2008;17 1:37-48.

11

12 259 6. MIFAFF. Food, Agriculture, Forestry and Fisheries statistical yearbook. Seoul:
13
14 260 Forestry and Fisheries (MIFAFF) Press; 2012.

15

16 261 7. Albertin CB, Simakov O, Mitros T, Wang ZY, Pungor JR, Edsinger-Gonzales
17
18 262 E, et al. The octopus genome and the evolution of cephalopod neural and mor
19
20 263 phological novelties. *Nature.* 2015;524 7564:220-4.

21 264 8. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel co
22
23 265 unting of occurrences of k-mers. *Bioinformatics.* 2011;27 6:764-70.

24

25 266 9. Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et a
26
27 267 l. Phased diploid genome assembly with single molecule real-time sequencing.
28
29 268 *Nat Methods.* 2016;13 12:1050.

30 269 10. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. B
31
32 270 USCO: assessing genome assembly and annotation completeness with single-cop
33
34 271 y orthologs. *Bioinformatics.* 2015;31 19:3210-2.

35

36 272 11. Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, et al. Widesprea
37
38 273 d polycistronic transcripts in fungi revealed by single-molecule mRNA sequenci
39
40 274 ng. *PLoS ONE.* 2015;10 7:e0132628.

41 275 12. Holt C and Yandell M. MAKER2: an annotation pipeline and genome-database
42
43 276 management tool for second-generation genome projects. *BMC bioinformatics.* 2
44
45 277 011;12 1:491.

46

47 278 13. Smit AFA HR, Green, P. RepeatMasker Open-3.0. 1996-2004 ([http://www.Repe
48
49 279 atMakser.org](http://www.RepeatMasker.org)).

50

51 280 14. Bao Z and Eddy SR. Automated de novo identification of repeat sequence fam
52
53 281 ilies in sequenced genomes. *Genome research.* 2002;12 8:1269-76.

54 282 15. Price AL, Jones NC and Pevzner PA. De novo identification of repeat families
55
56 283 in large genomes. *Bioinformatics.* 2005;21 suppl_1:i351-i8.

57

58 284 16. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucl
59
60 285 eic acids research.* 1999;27 2:573.

- 1 286 17. Korf I. Gene finding in novel genomes. *BMC bioinformatics*. 2004;5 1:59.
2
3 287 18. Nawrocki EP, Kolbe DL and Eddy SR. Infernal 1.0: inference of RNA alignm
4
5 288 ents. *Bioinformatics*. 2009;25 10:1335-7.
6
7 289 19. Gardner PP, Daub J, Tate J, Moore BL, Osuch IH, Griffiths-Jones S, et al. Rf
8
9 290 am: Wikipedia, clans and the “decimal” release. *Nucleic Acids Res*. 2010;39 su
10 291 ppl_1:D141-D5.
11
12 292 20. Lowe TM and Eddy SR. tRNAscan-SE: a program for improved detection of t
13
14 293 ransfer RNA genes in genomic sequence. *Nucleic acids research*. 1997;25 5:955
15
16 294 .
17
18 295 21. Li L, Stoeckert CJ, Jr. and Roos DS. OrthoMCL: identification of ortholog gro
19
20 296 ups for eukaryotic genomes. *Genome Res*. 2003;13 9:2178-89.
21
22 297 22. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al.
23
24 298 Pfam: the protein families database. *Nucleic Acids Res*. 2013;42 D1:D222-D30.
25
26 299 23. Löytynoja A and Goldman N. An algorithm for progressive multiple alignment
27
28 300 of sequences with insertions. *Proceedings of the National Academy of Sciences*
29
30 301 of the United States of America. 2005;102 30:10557-62.
31
32 302 24. Castresana J. Selection of conserved blocks from multiple alignments for their
33
34 303 use in phylogenetic analysis. *Molecular biology and evolution*. 2000;17 4:540-5
35
36 304 2.
37
38 305 25. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-ana
39
40 306 lysis of large phylogenies. *Bioinformatics*. 2014;30 9:1312-3.
41
42 307 26. Hedges SB, Dudley J and Kumar S. TimeTree: a public knowledge-base of div
43
44 308 ergence times among organisms. *Bioinformatics*. 2006;22 23:2971-2.
45
46 309 27. Han MV, Thomas GW, Lugo-Martinez J and Hahn MW. Estimating gene gain
47
48 310 and loss rates in the presence of error in genome assembly and annotation usi
49
50 311 ng CAFE 3. *Molecular biology and evolution*. 2013;30 8:1987-97.
51
52 312 28. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular bi*
53
54 313 ology and evolution. 2007;24 8:1586-91.
55
56 314 29. Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, et al. Fast, scala
57
58 315 ble generation of high-quality protein multiple sequence alignments using Clusta
59
60 316 l Omega. *Molecular systems biology*. 2011;7 1:539.
61
62 317 30. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and hig
63
64 318 h throughput. *Nucleic Acids Res*. 2004;32 doi:10.1093/nar/gkh340.
65

- 1 319 31. Price MN, Dehal PS and Arkin AP. FastTree 2—approximately maximum-likelih
2
3 320 ood trees for large alignments. PloS one. 2010;5 3:e9490.
4
5 321 32. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A
6
7 322 unified classification system for eukaryotic transposable elements. Nature Revie
8
9 323 ws Genetics. 2007;8 12:973-82.

10
11 324
12
13 325

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

1 326
2
3 327
4
5 328
6
7 329
8
9 330
10
11 331
12
13 332
14
15 333
16
17 334
18
19 335
20
21 336
22
23 337
24
25 338
26
27 339
28
29 340
30
31 341
32
33 342
34
35 343
36
37 344
38
39 345
40
41 346
42
43 347
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure legends

Figure 1: Common long-arm octopus (*Octopus minor*). **a** Habitat structure of mudflats and phenotypic differences between *Octopus minor* and *O. bimaculoides*. *O. minor* has a smaller body size and possesses longer, thinner arms than those of *O. bimaculoides*. **b** The distribution of *O. minor* is shown in dark red. The distribution map was updated from Roper *et al.* (1984).

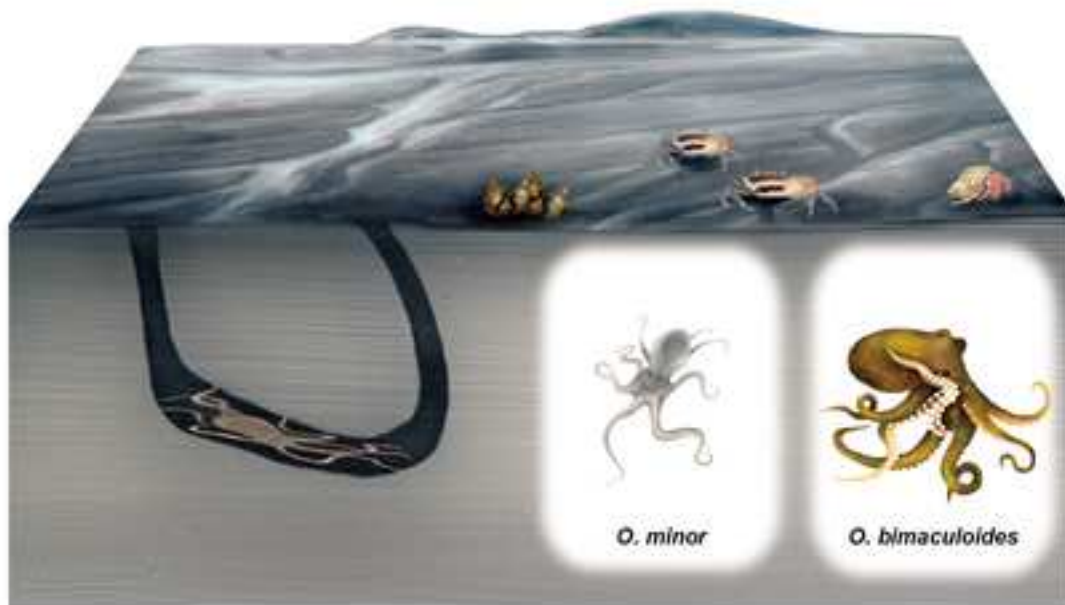
Figure 2: Gene family analysis for 14 bilaterian animal species. **a** Divergence times and gene family gain-and-loss analysis of 14 bilaterian animal species. **b** Heat map of expanded Pfam domains in the *O. minor* genome. OM, *Octopus minor*; OB, *Octopus bimaculoides*; LG, *Lottia gigantea*; CG, *Crassostrea gigas*; PF, *Pinctada fucata*; LA, *Lingula anatina*; CT, *Capitella teleta*; HR, *Helobdella robusta*; CE, *Caenorhabditis elegans*; DM, *Drosophila melanogaster*; DP, *Daphnia pulex*; SP, *Strongylocentrotus purpuratus*; MM, *Mus musculus*; HS, *Homo sapiens*.

Figure 3: Transposable element (TE) accumulation history in the *Octopus* genomes. Kimura distance-based copy divergence analysis of TEs for **a**, *O. minor* and **b**, *O. bimaculoides*. *x*-axis, K-value; *y*-axis, genome coverage for each type of TE.

Table 1 Overview of the assembly and annotation of the *Octopus minor* genome.

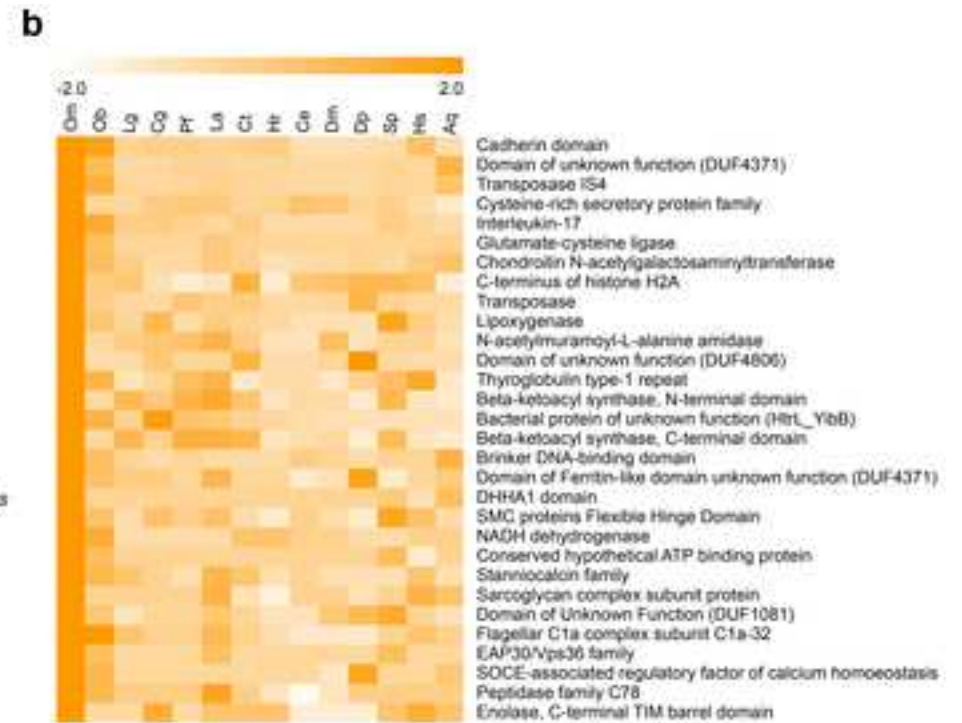
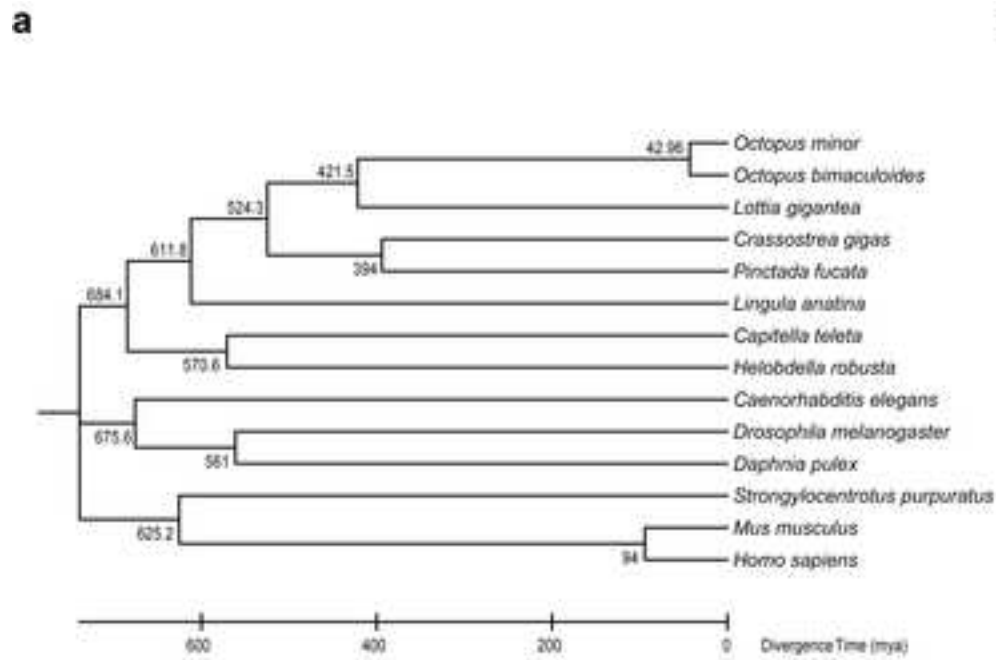
Total length (bp)	5,090,349,614
Number of contigs	41,584
Contig N50 (bp)	196,941
Largest contigs (bp)	3,027,443
GC content (%)	36.33
Number of protein-coding genes	30,010

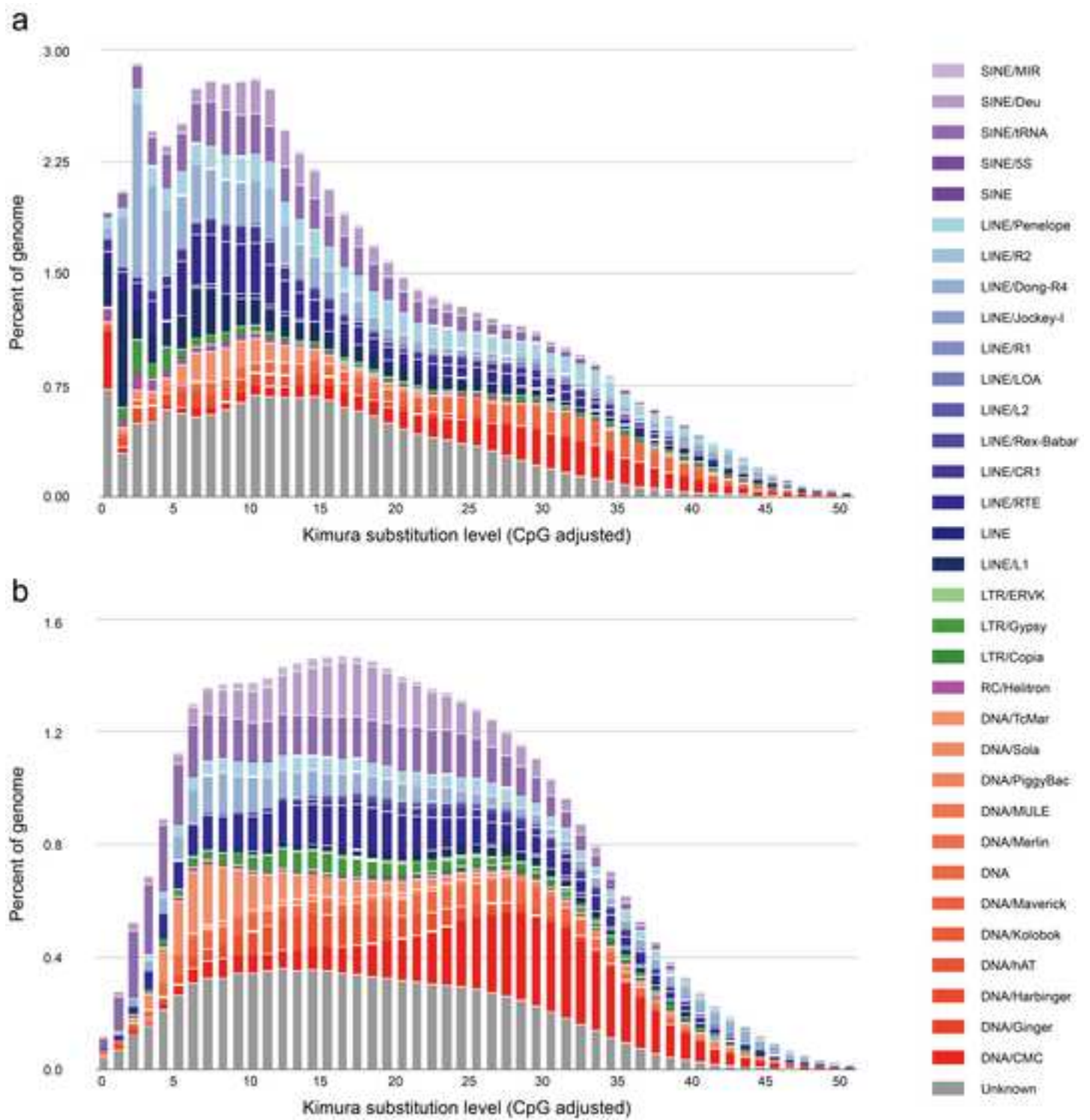
a



b









[Click here to access/download](#)

Supplementary Material

GIGA_Additional file 1_figure.docx





Click here to access/download
Supplementary Material
GIGA_Additional file 1_Table.docx

