

Manuscript Number:	GIGA-D-16-00040R2	
Full Title:	The genome of Antarctic-endemic Copepod, <i>Tigriopus kingsejongensis</i>	
Article Type:	Data Note	
Funding Information:	Korea Polar Research Institute (KR) (PE16070)	Dr. Hyun Park
	Korea Polar Research Institute (PE14260)	Dr Sanghee Kim
Abstract:	<p>Background: The Antarctic intertidal zone is continuously subject to extreme fluctuations in biotic and abiotic stressors, and the West Antarctic Peninsula is the most rapidly warming region on earth. Organisms living in Antarctic intertidal pools are therefore of great interest for research on topics such as evolutionary adaptation to extreme environments and the effects of climate change.</p> <p>Findings: Here, we report the whole-genome sequence of the Antarctic endemic Harpacticoid copepod, <i>Tigriopus kingsejongensis</i> with a total of 37 Gb raw DNA sequence using Illumina Miseq platform and the libraries were prepared with 65-fold coverage with a total length of 295 Mb. The final assembly consists of 48,368 contigs with an N50 contig length of 17.5 kb and 27,823 scaffolds with N50 contig length of 159.2 kb and a total of 12,772 coding genes were inferred using the MAKER annotation pipeline approach. Comparative genome analysis revealed that <i>T. kingsejongensis</i> specific genes are enriched in transport and metabolism processes. Furthermore, rapidly evolving genes related to energy metabolism showed signatures of positive selection.</p> <p>Conclusions: The genome of <i>T. kingsejongensis</i> will provide an interesting example of an evolutionary strategy for Antarctic cold adaptation, and offers new genetic insights into Antarctic intertidal biota.</p>	
Corresponding Author:	Hyun Park KOREA, REPUBLIC OF	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Seunghyun Kang	
First Author Secondary Information:		
Order of Authors:	Seunghyun Kang	
	Do-Hwan Ahn	
	Jun Hyuck Lee	
	Sung Gu Lee	
	Seung Chul Shin	
	Jungeun Lee	
	Gi-Sik Min	
	Hyoungseok Lee	
	Hyun-Woo Kim	
	Sanghee Kim	
	Hyun Park	

Order of Authors Secondary Information:	
Response to Reviewers:	<p>GIGA-D-16-00040 The genome of Antarctic-endemic Copepod, <i>Tigriopus kingsejongensis</i> Seunghyun Kang; Do-Hwan Ahn; Jun Hyuck Lee; Sung Gu Lee; Seung Chul Shin; Jungeun Lee; Gi-Sik Min; Hyoungseok Lee; Hyun-Woo Kim; Sanghee Kim; Hyun Park GigaScience</p> <p>Reviewer reports: Reviewer #1: This paper is in reasonably good shape and the data will be useful for comparison to temperate <i>Tigriopus</i>. My only remaining comment is that the coverage of the CEGMA is relatively low and I don't know how that impacts the estimates of turnover in gene families (which the authors estimate to be quite high). Could this be a result of the fact that the assembly is missing a sizable proportion of core eukaryotic genes?</p> <p>Response) We have added BUSCO completeness results and two supplementary tables (Table S4 about CEGMA results and Table S5 about BUSCO results) in data description line 132.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<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>	
Resources	Yes
<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>	
Availability of data and materials	Yes

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.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

[Click here to view linked References](#)

1 **The genome of Antarctic-endemic Copepod, *Tigriopus***
2
3
4 ***kingsejongensis***
5
6
7
8

9
10 Seunghyun Kang^{1¶}, Do-Hwan Ahn^{1¶}, Jun Hyuck Lee^{1,2}, Sung Gu Lee^{1,2}, Seung Chul Shin¹,
11
12 Jungeun Lee¹, Gi-Sik Min³, Hyoungseok Lee¹, Hyun-Woo Kim^{4*&}, Sanghee Kim^{5*&} & Hyun
13
14 Park^{1,2*&}
15
16
17
18

19 ¹ Unit of Polar Genomics, Korea Polar Research Institute, Yeonsu-gu, Incheon, South Korea
20

21 ² Polar Sciences, University of Science & Technology, Yuseong-gu, Daejeon, South Korea
22

23 ³ Department of Biological Sciences, Inha University, Incheon, South Korea
24

25 ⁴ Department of Marine Biology, Pukyong National University, Busan, South Korea
26
27

28 ⁵ Division of Polar Life Sciences, Korea Polar Research Institute, Yeonsu-gu, Incheon, South
29
30
31 Korea
32

33
34 * Corresponding author
35

36 E-mail: kimhw@pknu.ac.kr; sangheekim@kopri.re.kr; hpark@kopri.re.kr
37
38

39 ¶These authors contributed equally to this work.
40

41 &These authors also contributed equally to this work.
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 **Abstract**

2
3
4 2 **Background:** The Antarctic intertidal zone is continuously subject to extreme fluctuations in
5
6
7 3 biotic and abiotic stressors, and the West Antarctic Peninsula is the most rapidly warming
8
9 4 region on earth. Organisms living in Antarctic intertidal pools are therefore of great interest
10
11 5 for research on topics such as evolutionary adaptation to extreme environments and the
12
13 6 effects of climate change.

14
15
16
17 7 **Findings:** Here, we report the whole-genome sequence of the Antarctic endemic
18
19 8 Harpacticoid copepod, *Tigriopus kingsejongensis* with a total of 37 Gb raw DNA sequence
20
21 9 using Illumina Miseq platform and the libraries were prepared with 65-fold coverage with a
22
23 10 total length of 295 Mb. The final assembly consists of 48,368 contigs with an N50 contig
24
25 11 length of 17.5 kb and 27,823 scaffolds with N50 contig length of 159.2 kb and a total of
26
27 12 12,772 coding genes were inferred using the MAKER annotation pipeline approach.
28
29 13 Comparative genome analysis revealed that *T. kingsejongensis* specific genes are enriched in
30
31 14 transport and metabolism processes. Furthermore, rapidly evolving genes related to energy
32
33 15 metabolism showed signatures of positive selection.

34
35
36
37 16 **Conclusions:** The genome of *T. kingsejongensis* will provide an interesting example of an
38
39 17 evolutionary strategy for Antarctic cold adaptation, and offers new genetic insights into
40
41 18 Antarctic intertidal biota.

42
43
44
45 19 **Keywords:** Copepoda, Genome, Antarctic, adaptation, *Tigriopus*
46
47
48
49
50
51
52
53
54
55 20
56
57
58
59
60
61
62
63
64
65

22 **Data description**

23 The subclass copepods are very diverse and approximately 12,000 copepod species have been
24 described [1, 2]. They dominate the zooplankton community contributing about 70% of total
25 zooplankton biomass [3] and play an important role in the marine meiobenthic food web
26 linking between the phytoplankton and higher trophic levels [4]. Harpacticoid copepods of
27 the genus *Tigriopus* Norman 1868 are dominant members of shallow supratidal rock pools
28 worldwide. They are distributed among habitats that vary widely in salinity, temperature,
29 desiccation risk, and UV radiation, and have been used as a model system to investigate
30 topics such as osmoregulation [5], temperature adaptation [6, 7] and environmental
31 toxicology [8]. As the genome resources of copepods has been publically available
32 (*Tigriopus californicus* [http://i5k.nal.usda.gov/Tigriopus_californicus], *Tigriopus japonicus*
33 [9], *Eurytemora affinis* [http://i5k.nal.usda.gov/Eurytemora_affinis] and salmon louse
34 *Lepeophtheirus salmonis* [<http://sealouse.imr.no/>]), now it is possible to explore their
35 fundamental biological processes and physiological responses to diverse environments.

36 Antarctica provides not only an extreme habitat for extant organisms, but also a model for
37 research on evolutionary adaptations to cold environments [10, 11]. The Antarctic intertidal
38 zone, particularly in the Western Antarctic Peninsula region, is one of the most extreme
39 environments on earth. It also serves as a potential barometer for global climate changes,
40 since it is the fastest-warming region on earth [12]. Antarctic intertidal species that have
41 evolved stenothermal phenotypes through adaptation to a year-round climate of extreme cold
42 may now face extinction by global warming. The response of these species to further
43 warming in Western Antarctica is of serious concern; however, to date there have been few
44 studies focusing on species from the Antarctic intertidal zone.

45 *Tigriopus kingsejongensis* was first found and recognized as a new endemic species in a rock

1 46 pool in the Antarctic Peninsula, and is extremely cold-tolerant and can survive in frozen sea
2
3 47 water [13]. We observed the morphological differences, such as increased numbers of caudal
4
5 48 setae in nauplii, optimal growth temperature (ca. 8°C) and developmental characteristics have
6
7
8 49 been compared to those of the congener *Tigriopus japonicus*, which is found in the coastal
9
10 50 area of the Yellow Sea. *Tigriopus kingsejongensis* has evolved to overcome the unique
11
12 environmental constrains of Antarctica, and therefore provides an ideal experimental model
13 51
14
15 52 for all aspects of research on extreme habitats. This species may represent a case of rapid
16
17
18 53 speciation, since the intertidal zone on King George Island and surrounding areas did not
19
20 54 exist before 10,000 years ago [14]. *Tigriopus kingsejongensis* likely evolved as a distinct
21
22 species within this relatively short time period. Thus, inter- and intraspecies comparative
23 55
24
25 56 analyses of Antarctic *Tigriopus* species will help define the trajectory of adaptation to the
26
27
28 57 Antarctic environment and also provide insights into the genetic basis of *Tigriopus*
29
30 58 divergence and evolution.
31
32
33

34 59 **Library construction and sequencing**

35
36
37 60 *Tigriopus kingsejongensis* were collected from tidal pools in Potter Cove, near King
38
39 61 Sejong Station, on the northern Antarctic Peninsula (62°14'S, 58°47'W) (Fig. S1 and S2 in
40
41 additional file1) in January 2013 with a hand-nets. Water temperatures were $1.6 \pm 0.8^\circ\text{C}$
42 62
43
44 63 during this month. High-molecular-weight genomic DNA from pooled *T. kingsejongensis*
45
46
47 64 was extracted using the DNeasy Blood & Tissue Kit (Qiagen). For Illumina Miseq
48
49 65 sequencing, four library types were constructed with 350, 400, 450, and 500 bp for paired-
50
51
52 66 end libraries, and 3 kb and 8 kb for mate-pair libraries, prepared using the standard Illumina
53
54 67 sample preparation methods (Table 1). All sequencing processes were performed according
55
56
57 68 to the manufacturer's instructions (Illumina).
58
59
60
61
62
63
64
65

1 69 RNA was prepared from pooled *T. kingsejongensis* and *Tigriopus japonicus*
2
3
4 70 specimens from two different temperature experiments (4°C and 15°C) using the RNeasy
5
6 71 Mini Kit (Qiagen). For Illumina Miseq sequencing, subsequent experiments were carried out
7
8 72 under the manufacturer's instructions (Illumina). The *de novo* transcriptome assembly was
9
10
11 73 performed with CLC Genomics Workbench, setting the minimum allowed contig length to
12
13 74 200 nucleotides. The assembly process generated 40,172 contigs with a max length of 23,942
14
15
16 75 bp and an N50 value of 1,093 bp. These generated contigs were used as reference sequences
17
18 76 for mapping of trimmed reads, and fold changes in expression for each gene were calculated
19
20
21 77 with a significance threshold of $P \leq 0.05$ using CLC Genomics Workbench (Table 2 and 3).
22
23 78
24
25
26

27 79 **Genome assembly**

28
29
30 80 First, k-mer analysis was conducted using jellyfish 2.2.5 [15] to estimate the genome
31
32 81 size from DNA paired-end libraries. The estimated genome size was 298 Mb with main peak
33
34 82 at a depth of ~39x (Fig. 1). Then, assemblies were performed using a Celera Assembler with
35
36
37 83 Illumina short reads [16]. Prior to assembly, Illumina reads were trimmed using the FASTX-
38
39 84 Toolkit (http://hannonlab.cshl.edu/fastx_toolkit) with parameters -t 20, -l 70 and -Q 33, after
40
41
42 85 which a paired sequence from trimmed Illumina reads was selected. Finally, trimmed
43
44 86 Illumina reads with 65-fold coverage (insert sizes 350, 400, 450, and 500 bp) were obtained
45
46
47 87 and converted to the FRG file format (required by the Celera assembler) using FastqToCA.
48
49 88 Assembly was performed on a 96-processor workstation with Intel Xeon X7460 2.66 GHz
50
51
52 89 processors and 1 terabyte RAM with the following parameters: overlapper = ovl, unitigger =
53
54 90 bogart, utgGraphErrorRate = 0.03, utgGraphErrorLimit = 2.5, utgMergeErrorRate = 0.030,
55
56 91 utgMergeErrorLimit = 3.25, ovlErrorRate = 0.1, cnsErrorRate = 0.1, cgwErrorRate = 0.1,
57
58
59 92 merSize = 22, and doOverlapBasedTrimming = 1. The initial Celera assembly had a total size
60
61
62
63
64
65

1 93 of 305 Mb, N50 contig size of 17,566 bp, and max contig size of 349.5 kb. Scaffolding was
2
3 94 completed using the software SSPACE 2.0 scaffolder using mate-paired data [17].
4
5 95 Subsequently, we closed gaps using Gapfiller Ver.1.9 software with 65× trimmed Illumina
6
7 96 reads with default settings [18]. *De novo* assembly of 203 million reads from paired-end
8
9 97 libraries and mate-paired libraries yielded a draft assembly (65-fold coverage) with a total
10
11 97 length of 295 Mb, and contig and scaffold N50 sizes of 17.6 kb and 159.2 kb, respectively
12
13 98 (Table 4 and Fig. 2).
14
15
16
17
18
19

20 100 **Annotation**

21
22 101 We used MAKER for genome annotation [19]. MAKER is a portable and easily
23
24 102 configurable genome annotation pipeline. MAKER first identified repetitive elements using
25
26 103 RepeatMasker [20]. This masked genome sequence was used for *ab initio* gene prediction
27
28 104 with SNAP software [21], after which alignment of expressed sequence tags with BLASTn
29
30 105 and protein information from tBLASTx were included. We used the *de novo* repeat library of
31
32 106 *T. kingsejongensis* from RepeatModeler for RepeatMasker; proteins from five species with
33
34 107 data from *D. melanogaster*, *D. pulex*, *T. japonicus*, and *Tigriopus californicus* were included
35
36 108 in the analysis. RNA-seq-based gene prediction was performed by aligning all RNA-seq data
37
38 109 against the assembled genome using TopHat [22], and Cufflinks [23] was used to predict
39
40 110 cDNAs from the resultant data. Next, MAKER polished the alignments using the program
41
42 111 Exonerate [24], which provided integrated information to synthesize SNAP annotation.
43
44 112 MAKER then selected and revised the final gene model considering all information. A total
45
46 113 of 12,772 genes were predicted using MAKER in *T. kingsejongensis*. Annotated genes
47
48 114 contained an average of 4.6 exons, with an average mRNA length of 1,090 bp. Additionally,
49
50 115 12,562 of 12,772 genes were assigned preliminary functions based on automated annotation
51
52 116 using Blast2GO (Ver. 2.6.0) [25] (Fig. S3 and S4 in additional file 1). The Infernal software
53
54
55
56
57
58
59
60
61
62
63
64
65

1 117 package (Ver. 1.1) [26] and covariance models (CMs) from the Rfam database [27] were
2
3 118 used to identify other non-coding RNAs in the *T. kingsejongensis* scaffold. We identified
4
5
6 119 putative tRNA genes using tRNAscan-SE [28] (Table S1 in additional file 2). tRNAscan-SE
7
8 120 uses a covariance model (CM) that scores candidates based on their sequence and predicted
9
10
11 121 secondary structures.

12
13 122 Non-gap sequences occupied 284.8 Mb (96.5%), and simple sequence repeats (SSRs)
14
15
16 123 were 1.2 Mb (0.4%) in total (Table S2 in additional file 2). Transposable elements (TEs)
17
18 124 comprised 6.5 Mb, which is roughly 2.3% of the assembled genome (Table S2 and S3 in
19
20
21 125 additional file 2). On the basis of homology and *ab initio* gene prediction, we found that the
22
23 126 genome of *T. kingsejongensis* contains 12,772 protein-coding genes (Table 5). By assessing
24
25 127 the quality of the annotated 12,772 gene models, we found that 11,686 protein-coding genes
26
27
28 128 (91.5%) were supported by the RNA-seq data, of which, 7,325 (63%) showed similarity to
29
30
31 129 proteins from other species. Analysis of Core Eukaryotic Genes Mapping Approach
32
33 130 (CEGMA) [29] showed that 179 of 248 CEGMA score genes were fully annotated (72.18 %
34
35 131 completeness) and 197 of 248 genes were partially annotated (79.44 % completeness) (Table
36
37 132 S4 in additional file 2). We also found that Benchmarking Universal Single-Copy Orthologs
38
39
40 133 (BUSCO) [30] analysis showed that the genome assembly contains 71 % of complete and 6
41
42 134 % of partial Metazoan orthologous gene set (Table S5 in additional file 2).

46 135 **Gene families**

47
48
49 136 The orthologous groups were identified from 11 species (*T. kingsejongensis*, *Aedes*
50
51 137 *aegypti*, *D. melanogaster*, *Ixodes scapularis*, *M. martensii*, *Strigamia martima*, *Tetranychus*
52
53
54 138 *urticae*, *D. pulex*, *Homo sapiens*, *Ciona intestinalis*, and *Caenorhabditis elegans*) (Table 6)
55
56
57 139 using OrthoMCL [31] with standard parameters and options, and transcript variants other
58
59 140 than the longest translation forms were removed. For *T. kingsejongensis*, the coding sequence
60
61
62
63
64
65

1 141 from the MAKER annotation pipeline was used. The 1:1:1 single-copy orthologous genes
2
3 142 were subjected to phylogenetic construction and divergence time estimation. Protein-coding
4
5 143 genes were aligned using PRANK with the codon alignment option [32], and poorly aligned
6
7
8 144 sequences with gaps were removed using Gblock under the codon model [33]. We
9
10 145 constructed a maximum-likelihood phylogenetic tree using RAxML with 1,000 bootstrap
11
12 146 values [34] and calibrated the divergence time between species with TimeTree [35]. Finally,
13
14 147 the average gene gain/loss rate along the given phylogeny was identified using the program
15
16 148 CAFÉ 3.1 [36]. We constructed orthologous gene clusters using four arthropod species
17
18 149 (Antarctic copepod, *T. kingsejongensis*; scorpion, *Mesobuthus martensii*; fruit fly, *Drosophila*
19
20 150 *melanogaster* and water flea, *Daphnia pulex*) to compare the genomic features and the
21
22 151 adaptive divergence in the arthropods. In total, 2,063 gene families are shared by all four
23
24 152 species, and 1,028 genes are specific to the Antarctic copepod. *Tigriopus kingsejongensis*
25
26 153 shares 4,559 (73.5%) gene families with *D. pulex*, which belongs to the same crustacean
27
28 154 lineage Vericrustacea, 3,531 (56.9%) with *D. melanogaster*, and 3,231 (52.1%) with *M.*
29
30 155 *martensii* (Fig. 3A). Gene ontology (GO) analysis revealed that the 1,028 *T. kingsejongensis*-
31
32 156 specific genes are enriched in transport (single-organism transport, GO: 0044765;
33
34 157 transmembrane transport, GO: 0055085; ion transport, GO: 0006811; cation transport, GO:
35
36 158 0006812) and single-organism metabolic processes (GO: 0044710) (Table S6 and S7 in
37
38 159 additional file 2). Subsequently, we performed gene gain-and-loss analysis on 11
39
40 160 representative species, and found that *T. kingsejongensis* gained 735 gene families and lost
41
42 161 4,401 gene families (Fig. 3B). Thus, this species exhibits a gene family turnover of 5,136, the
43
44 162 largest value among the eight arthropods. We also analyzed expansion and contraction of the
45
46 163 gene families (Table S8-S11 in additional file 2), and found 232 significantly expanded gene
47
48 164 families in *T. kingsejongensis*; these gene families are significantly overrepresented in amino
49
50 165 acid metabolism and carbohydrate metabolism in KEGG metabolic pathways.
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 166
2
3
4
5 167
6
7
8
9 168
10
11 169
12
13 170
14
15 171
16 172
17
18 173
19
20 174
21 175
22
23 176
24
25 177
26 178
27
28 179
29
30 180
31
32 181
33
34 182
35
36 183
37
38 184
39
40 185
41
42 186
43
44 187
45
46 188
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Genome evolution

Adaptive functional divergence caused by natural selection is commonly estimated based on the ratio of nonsynonymous (dN) to synonymous (dS) mutations. To estimate dN , dS , and average dN/dS ratio (w), and lineage-specific PSGs in *T. kingsejongensis* and *T. japonicus*, protein-coding genes from *T. japonicus* were added to define orthologous gene families among the four species (*T. kingsejongensis*, *T. japonicus*, *D. pulex*, and *D. melanogaster*) using the program OrthoMCL with the same conditions previously described. We identified 2,937 orthologous groups shared by all four species, and single-copy gene families were used to construct a phylogenetic tree and estimate the time since divergence using the same methods described above. Each of the identified orthologous genes was aligned using the PRANK, and poorly aligned sequences with gaps were removed using Gblock. Alignments showing less than 40% identity and genes shorter than 150 bp were eliminated in subsequent procedures. The values of dN , dS and w were estimated from each gene using the Codeml program implemented in the PAML package with the free-ratio model [37] under F3X4 codon frequencies, and orthologs with $w \leq 5$ and $dS \leq 3$ were retained [38]. To examine the accelerated nonsynonymous divergence in either *T. kingsejongensis* or *T. japonicus* lineage, a binomial test [39] was used to determine GO categories with at least 20 orthologous genes. To define PSGs in *T. kingsejongensis* and *T. japonicus*, we applied basic and branch-site models, and Likelihood Ratio Tests (LRTs) were used to remove genes under relaxation of selective pressure. To investigate which functional categories and pathways were enriched in the PSGs, we performed DAVID Functional Annotation [40] with Fisher's exact test (cutoff: $P \leq 0.05$).

1 189 The average w value from 2,937 co-orthologous genes of *T. kingsejongensis* (0.0027)
2
3
4 190 is higher than that of *T. japonicus* (0.0022). The GO categories that showing evidence of
5
6 191 accelerated evolution in *T. kingsejongensis* are energy metabolism (generation of precursor
7
8 192 metabolites and energy, GO: 0006091; cellular respiration, GO: 0045333) and carbohydrate
9
10 193 metabolism (monosaccharide metabolic process, GO: 0005996; hexose metabolic process,
11
12 194 GO: 0019318) (Figure 4A, Table S12 in Additional file 2). Branch-site model analysis
13
14 195 showed that the genes belonging to the functional categories above have undergone a
15
16 196 significant positive selection process by putative functional divergence in certain lineages.
17
18 197 There are 74 and 79 positively selected genes (PSGs) in *T. kingsejongensis* (Table S13 in
19
20 198 Additional file 2) and *T. japonicus* (Table S14 in Additional file 2), respectively. The
21
22 199 functional categories enriched in *T. kingsejongensis*, when compared to *T. japonicus*, support
23
24 200 the idea that the functional divergence in *T. kingsejongensis* is strongly related to energy
25
26 201 metabolism (oxidative phosphorylation, GO: 0006119; energy-coupled proton transport down
27
28 202 electrochemical gradient, GO: 0015985; ATP synthesis-coupled proton transport, GO:
29
30 203 0015986; generation of precursor metabolites and energy, GO: 0006091) (Figure 4B, Table
31
32 204 S15 and S16 in Additional file 2). In particular, three of the identified genes are involved in
33
34 205 the oxidative phosphorylation (OxPhos) pathway, which provides the primary cellular energy
35
36 206 source in the form of adenosine triphosphate (ATP). These three genes are nuclear-encoded
37
38 207 mitochondrial genes: the catalytic F1 ATP synthase subunit alpha (*ATP5A*) (Fig. S5 in
39
40 208 Additional file 1), a regulatory subunit acting as an electron transport chain such as
41
42 209 ubiquinol-cytochrome *c* reductase core protein (*UQCRC1*) (Fig. S6 in Additional file 1), and
43
44 210 an electron transfer flavoprotein alpha subunit (*ETFa*) (Fig. S7 in Additional file 1).
45
46
47
48
49
50
51
52
53
54
55

56 211 **Availability of supporting data**

57
58 212 The data for *T. kingsejongensis* genome and transcriptome has been deposited in the SRA as
59
60
61
62
63
64
65

1 213 BioProject PRJNA307207 and PRJNA307513, respectively.
2
3

4 214 **List of abbreviations**

5
6
7 215 simple sequence repeats, SSRs; Transposable elements, TEs; CEGMA, Core Eukaryotic
8
9
10 216 Genes Mapping Approach; Gene ontology, GO
11
12

13 217 **Competing interests**

14
15
16 218 The authors declare no competing interests.
17
18 219

19 20 220 **Funding**

21
22
23 221 This work was supported by an Antarctic organisms: Cold-adaptation mechanism and its
24
25 222 application grant (PE16070) and the basic research program (PE14260) funded by the Korea
26
27
28 223 Polar Research Institute (KOPRI).
29
30

31 224 **Author contributions**

32
33
34 225 H.P., Sanghee Kim and H.W.K. conceived and designed experiments and analyses;
35
36
37 226 Seunghyun Kang, D.-H.A., S.G.L., S.C.S., J.L., G.S.M. and H.L. performed experiments and
38
39 227 conducted bioinformatics. Seunghyun Kang, H.W.K., Sanghee Kim and H.P. wrote the paper.
40
41

42 228 **Acknowledgements**

43
44
45 229 We would like to thank Joseph A. Covi for comments and discussion.
46
47
48

49 230
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **References**

- 2
3
4 232 1. Huys R, Boxshall GA: *Copepod evolution*. Ray Society; 1991.
- 5
6 233 2. Humes AG: **How many copepods?** *Hydrobiologia* 1994, **292**:1-7.
- 7
8
9 234 3. Wells P, Persoone G, Jaspers E, C. C: *Marine ecotoxicological tests with zooplankton*.
10
11 235 *In: Persoone, G., Jaspers, E., Claus, C. (Eds.), Ecotoxicological Testing for the*
12
13 236 *Marine Environment*. Inst. Mar. Sci. Res., Bredene; 1984.
- 14
15
16 237 4. Ruppert E, Fox R, Barnes R: **Invertebrate Zoology, A Functional Evolutionary**
17
18 238 **Approach**. Brooks/Cole-Thomson Learning. Belmont, CA 2003.
- 19
20
21 239 5. Goolish E, Burton R: **Energetics of osmoregulation in an intertidal copepod:**
22
23 240 **Effects of anoxia and lipid reserves on the pattern of free amino accumulation.**
24
25 241 *Funct Ecol* 1989:81-89.
- 26
27
28 242 6. Lazzaretto I, Libertini A: **Karyological comparison among different**
29
30 243 **Mediterranean populations of the genus *Tigriopus* (Copepoda Harpacticoida).**
31
32 244 *Boll Zool* 2009, **53**:197-201.
- 33
34
35 245 7. Davenport J, Barnett P, McAllen R: **Environmental tolerances of three species of**
36
37 246 **the harpacticoid copepod genus *Tigriopus*. *J Mar Biol Assoc UK* 1997, **77**:3-16.**
- 38
39
40 247 8. Raisuddin S, Kwok KW, Leung KM, Schlenk D, Lee J-S: **The copepod *Tigriopus*: A**
41
42 248 **promising marine model organism for ecotoxicology and environmental**
43
44 249 **genomics. *Aquat Toxicol* 2007, **83**:161-173.**
- 45
46
47 250 9. Lee J-S, Rhee J-S, Kim R-O, Hwang D-S, Han J, Choi B-S, Park GS, Kim I-C, Park
48
49 251 HG, Lee Y-M: **The copepod *Tigriopus japonicus* genomic DNA information**
50
51 252 **(574Mb) and molecular anatomy. *Mar Environ Res* 2010, **69**:S21-S23.**
- 52
53
54 253 10. Thorne MAS, Kagoshima H, Clark MS, Marshall CJ, Wharton DA: **Molecular**
55
56 254 **analysis of the cold tolerant Antarctic Nematode, *Panagrolaimus davidi*. *PLOS***
57
58 255 ***one* 2014, **9**:e104526.**
- 59
60
61
62
63
64
65

- 1 256 11. Everatta MJ, Worlandb MR, Balea JS, Conveyb P, Hayward SAL: **Pre-adapted to**
2
3 257 **the maritime Antarctic? – Rapid cold hardening of the midge, *Eretmoptera***
4
5 258 ***murphyi*. *J Insect Physiol* 2012, **58**:1104–1111.**
6
7
8 259 12. Bromwich DH, Nicolas JP, Monaghan AJ, Lazzara MA, Keller LM, Weidner GA,
9
10 260 Wilson AB: **Central West Antarctica among the most rapidly warming regions**
11
12 **on Earth. *Nature Geoscience* 2013, **6**:139-145.**
13
14
15 261 13. Park E-O, Lee S, Cho M, Yoon SH, Lee Y, Lee W: **A new species of the genus**
16
17 ***Tigriopus* (Copepoda: Harpacticoida: Harpacticidae) from Antarctica. *Proc Biol***
18 263 ***Soc Wash* 2014, **127**:138-154.**
19
20 264
21
22 265 14. Birkenmajer K: **Geology of Admiralty Bay, King George Island (South Shetland**
23
24 **Islands). An outline. *Pol Polar Res* 1980, **1**:29-54.**
25 266
26
27 267 15. Marçais G, Kingsford C: **A fast, lock-free approach for efficient parallel counting**
28
29 **of occurrences of k-mers. *Bioinformatics* 2011, **27**:764-770.**
30 268
31
32 269 16. Myers EW, Sutton GG, Delcher AL, Dew IM, Fasulo DP, Flanigan MJ, Kravitz SA,
33
34 Mobarry CM, Reinert KH, Remington KA, et al: **A whole-genome assembly of**
35 270 ***Drosophila*. *Science* 2000, **287**:2196-2204.**
36
37 271
38
39 272 17. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W: **Scaffolding pre-**
40
41 **assembled contigs using SSPACE. *Bioinformatics* 2011, **27**:578-579.**
42 273
43
44 274 18. Nadalin F, Vezzi F, Policriti A: **GapFiller: a *de novo* assembly approach to fill the**
45
46 **gap within paired reads. *BMC Bioinformatics* 2012, **13**:S8.**
47 275
48
49 276 19. Holt C, Yandell M: **MAKER2: an annotation pipeline and genome-database**
50
51 **management tool for second-generation genome projects. *BMC Bioinformatics***
52 277
53 **2011, **12**:491.**
54 278
55
56 279 20. Smit AFA HR, Green, P.: **RepeatMasker Open-3.0. 1996-2004**
57
58 **(<http://www.RepeatMasker.org>).**
59 280
60
61
62
63
64
65

- 1 281 21. Korf I: **Gene finding in novel genomes.** *BMC Bioinformatics* 2004, **5**:59.
2
3 282 22. Trapnell C, Pachter L, Salzberg SL: **TopHat: discovering splice junctions with**
4
5 **RNA-Seq.** *Bioinformatics* 2009, **25**:1105-1111.
6 283
7
8 284 23. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg
9
10 SL, Wold BJ, Pachter L: **Transcript assembly and quantification by RNA-Seq**
11 285
12 **reveals unannotated transcripts and isoform switching during cell**
13 286
14 **differentiation.** *Nat Biotech* 2010, **28**:511-515.
15 287
16
17 288 24. Slater GS, Birney E: **Automated generation of heuristics for biological sequence**
18
19 **comparison.** *BMC Bioinformatics* 2005, **6**:31.
20 289
21
22 290 25. Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M: **Blast2GO: a**
23
24 **universal tool for annotation, visualization and analysis in functional genomics**
25 291
26 **research.** *Bioinformatics* 2005, **21**:3674-3676.
27 292
28
29 293 26. Nawrocki EP, Kolbe DL, Eddy SR: **Infernal 1.0: inference of RNA alignments.**
30
31 *Bioinformatics* 2009, **25**:1335-1337.
32 294
33
34 295 27. Gardner PP, Daub J, Tate J, Moore BL, Osuch IH, Griffiths-Jones S, Finn RD,
35
36 Nawrocki EP, Kolbe DL, Eddy SR, Bateman A: **Rfam: Wikipedia, clans and the**
37 296
38 **"decimal" release.** *Nucleic Acids Res* 2011, **39**:D141-145.
39 297
40
41 298 28. Lowe TM, Eddy SR: **tRNAscan-SE: a program for improved detection of transfer**
42
43 **RNA genes in genomic sequence.** *Nucleic Acids Res* 1997, **25**:955-964.
44 299
45
46 299 29. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate core**
47
48 **genes in eukaryotic genomes.** *Bioinformatics* 2007, **23**:1061-1067.
49 300
50 301
51 302 30. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO:**
52
53 **assessing genome assembly and annotation completeness with single-copy**
54 303
55 **orthologs.** *Bioinformatics* 2015:btv351.
56 304
57
58 305 31. Li L, Stoeckert CJ, Roos DS: **OrthoMCL: identification of ortholog groups for**
59
60
61
62
63
64
65

- 1 306 **eukaryotic genomes.** *Genome Res* 2003, **13**:2178-2189.
- 2
- 3 307 32. Löytynoja A, Goldman N: **An algorithm for progressive multiple alignment of**
- 4
- 5 **sequences with insertions.** *Proc Natl Acad Sci U S A* 2005, **102**:10557-10562.
- 6 308
- 7
- 8 309 33. Castresana J: **Selection of conserved blocks from multiple alignments for their use**
- 9
- 10 **in phylogenetic analysis.** *Mol Biol Evol* 2000, **17**:540-552.
- 11 310
- 12
- 13 311 34. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-**
- 14
- 15 **analysis of large phylogenies.** *Bioinformatics* 2014, **30**:1312-1313.
- 16 312
- 17
- 18 313 35. Hedges SB, Dudley J, Kumar S: **TimeTree: a public knowledge-base of divergence**
- 19
- 20 **times among organisms.** *Bioinformatics* 2006, **22**:2971-2972.
- 21 314
- 22
- 23 315 36. Han MV, Thomas GW, Lugo-Martinez J, Hahn MW: **Estimating gene gain and loss**
- 24
- 25 **rates in the presence of error in genome assembly and annotation using CAFE 3.**
- 26 316
- 27 *Mol Biol Evol* 2013, **30**:1987-1997.
- 28 317
- 29
- 30 318 37. Yang Z: **PAML 4: phylogenetic analysis by maximum likelihood.** *Mol Biol Evol*
- 31
- 32 2007, **24**:1586-1591.
- 33 319
- 34
- 35 320 38. Zhang G, Li C, Li Q, Li B, Larkin DM, Lee C, Storz JF, Antunes A, Greenwold MJ,
- 36
- 37 Meredith RW: **Comparative genomics reveals insights into avian genome**
- 38 321
- 39 **evolution and adaptation.** *Science* 2014, **346**:1311-1320.
- 40 322
- 41
- 42 323 39. Consortium TCSaA: **Initial sequence of the chimpanzee genome and comparison**
- 43
- 44 **with the human genome.** *Nature* 2005, **437**:69-87.
- 45 324
- 46
- 47 325 40. Huang DW, Sherman BT, Lempicki RA: **Systematic and integrative analysis of**
- 48
- 49 **large gene lists using DAVID bioinformatics resources.** *Nature protocols* 2008,
- 50 326
- 51 **4**:44-57.
- 52 327
- 53
- 54
- 55 328
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

1 329 **Figure legends**

2
3
4 330 **Figure 1** Estimation of the *T. kingsejongensis* genome size based on 33-mer analysis. The x-
5
6
7 331 axis represents the depth (peak at 39X) and the y-axis represents the proportion. The genome
8
9 332 size was estimated as 298 Mb (total k-mer number/volume peak).

10
11 333
12
13
14 334 **Figure 2** Scaffold and contig size distributions of *T. kingsejongensis*. The percentage of the
15
16 335 assembly included (y-axis) in contigs or scaffolds of a minimum size (x-axis, log scale) is
17
18 336 shown for the contig (red) and scaffold (blue).

19 337
20
21
22
23
24 338 **Figure 3** Comparative genome analyses of the *T. kingsejongensis* genome. **a** Venn diagram
25
26 339 of orthologous gene clusters between the four arthropod lineages. **b** Gene family gain-and-
27
28 340 loss analysis. The number of gained gene families (red), lost gene families (blue) and orphan
29
30 341 gene families (black) are indicated for each species. Time lines specify divergence times
31
32 342 between the lineages.

33
34
35
36
37 343
38
39 344 **Table legends**

40
41
42 345 **Table 1** Statistics for each DNA library.

43
44
45 346 **Table 2** Sequencing and assembly results of transcriptome analysis of *T. japonicus*.

46
47
48 347 **Table 3** Sequencing statistics of RNA-seq analysis of *T. kingsejongensis*.

49
50
51 348 **Table 4** Statistics of genome assembly.

52
53
54 349 **Table 5** General statistics of genes in *T. kingsejongensis*.

55
56
57 350 **Table 6** Summary of orthologous gene clusters in the 11 representative species.

58
59
60
61
62
63
64
65

Table 1 Statistics for each DNA library.

Library		Reads (n)	Ave. length	Sequences (bp) (n)	Reads (trimmed) (n)	Ave. length	Sequences (trimmed) (n)
Paired-end	Sum	99,710,266		29,271,916,613	65,644,374		14,668,956,871
	350S1	6,661,392	300	2,005,078,992	4,446,394	233	1,034,231,244
	350S2	4,933,058	265	1,311,700,122	4,618,711	211	975,471,763
	400S1	65,668,598	300	19,766,247,998	36,863,154	228	8,397,426,481
	450S1	3,418,988	300	1,029,115,388	2,812,455	230	646,302,159
	450S2	8,009,162	245	1,968,652,020	7,660,814	199	1,527,566,312
	500S1	11,019,068	289	3,191,122,093	9,242,846	226	2,087,958,911
Mate-Paired	Sum	103,373,998		7,753,049,850	73,515,391		5,169,006,268
	3KS1	8,374,238	75	628,067,850	6,745,546	73	493,099,413
	3KS2	9,250,994	75	693,824,550	5,281,513	65	344,618,723
	3KS3	51,349,594	75	3,851,219,550	39,147,167	72	2,816,638,666
	3KS4	3,063,232	75	229,742,400	1,740,986	65	112,554,745
	8KS1	9,847,636	75	738,572,700	7,887,612	73	572,246,251
	8KS2	16,322,038	75	1,224,152,850	9,653,293	65	630,842,698
	8KS3	5,166,266	75	387,469,950	3,059,274	65	199,005,774
Total	203,084,264		37,024,966,463	139,159,765		19,837,963,139	
Coverage (folds)			120.7		64.7		

Table 2 Sequencing and assembly results of transcriptome analysis of *T. japonicus*.

Sequencing	
Total reads (n)	37,956,160
Total bases (n)	7,714,415,316
Trimmed reads (n)	35,577,636
Trimmed bases (n)	5,989,188,343
Assembly	
Contigs (n)	40,172
Total contig length (bases)	28,850,726
N50 contig length (bases)	1,093
Max scaffold length (bases)	23,942
Annotation	
With blast results	20,392
Without blast hits	7,090
With mapping results	8,172
Annotated sequences	4,518

Table 3 Sequencing statistics of RNA-seq analysis of *T. kingsejongensis*.

	4°C	15°C
Total reads (n)	15,786,118	16,417,072
Total bases (n)	3,567,662,668	3,763,295,032
Trimmed reads (n)	14,845,103	15,388,513
Trimmed bases (n)	2,761,189,158	2,833,805,442

Table 4 Statistics of genome assembly.

Celera assembler (Version : 8.0)

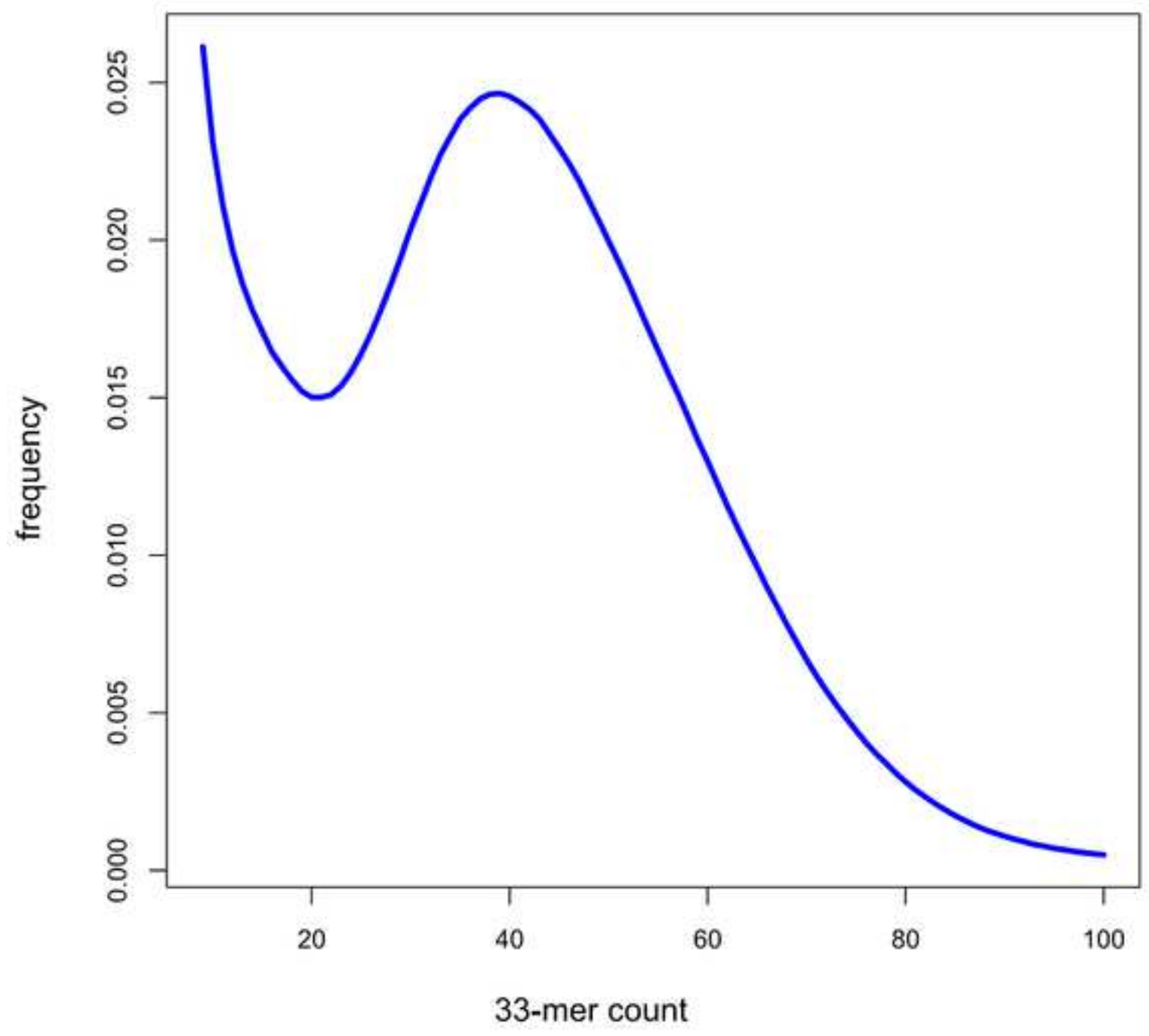
Scaffold	Total scaffold length (bases)	295,233,602
	Gap size (bases)	10,474,460
	Scaffolds (n)	11,558
	N50 scaffold length (bases)	159,218
	Max scaffold length (bases)	3,401,446
Contig	Total contig length (bases)	305,712,242
	Contigs (n)	48,368
	N50 contig length (bases)	17,566
	Max contig length (bases)	349,507

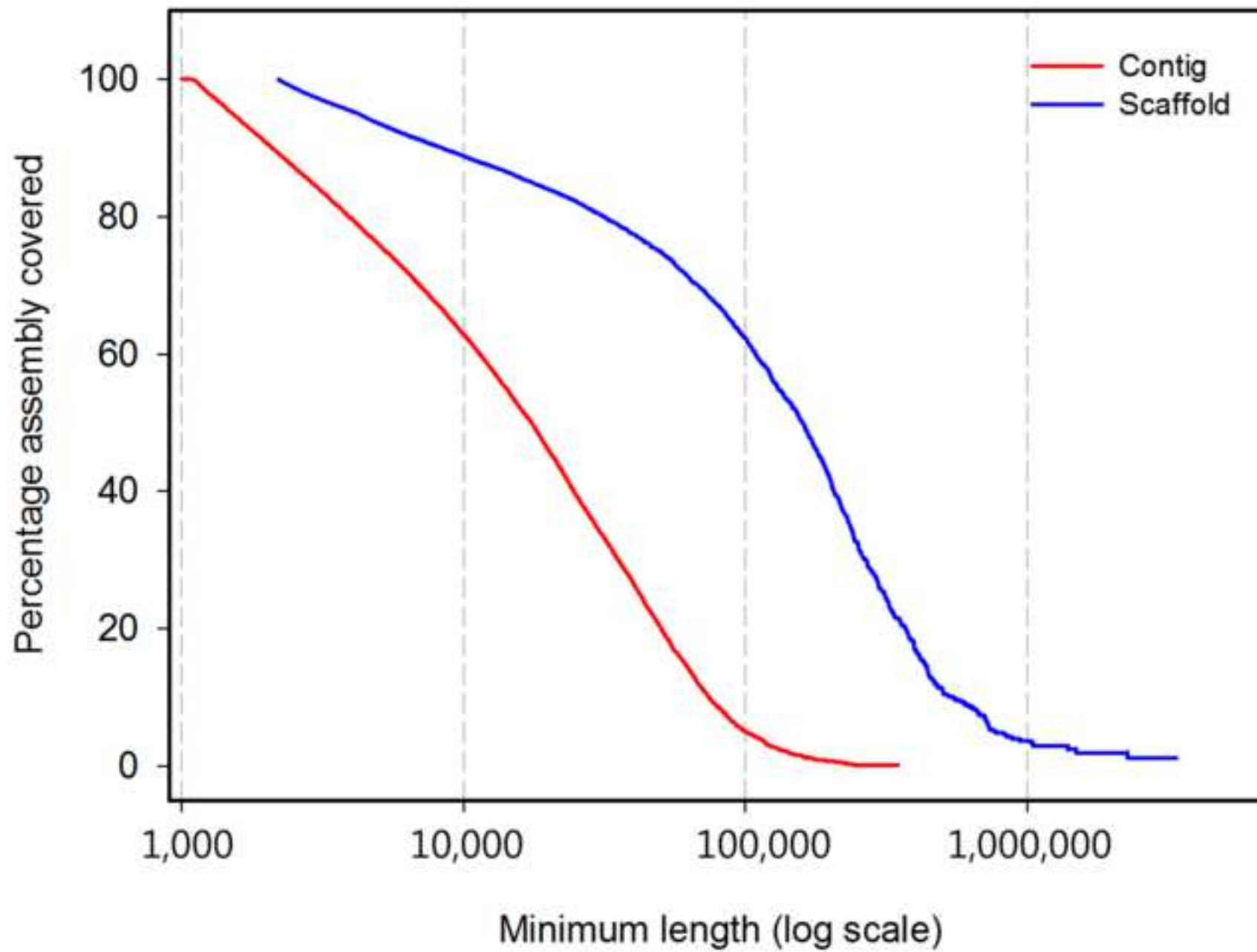
Table 5 General statistics of genes in *T. kingsejongensis*.

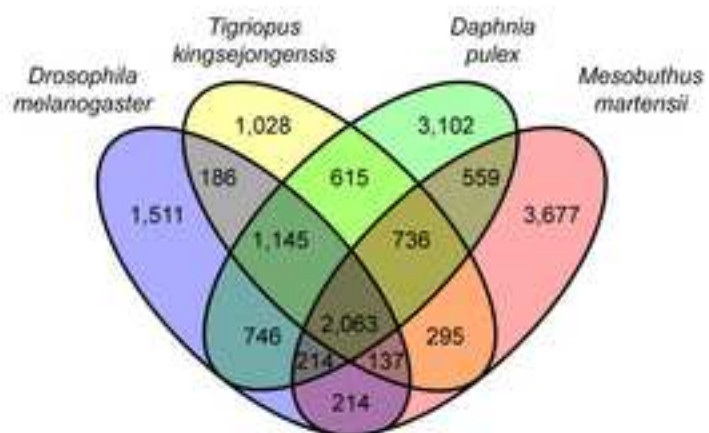
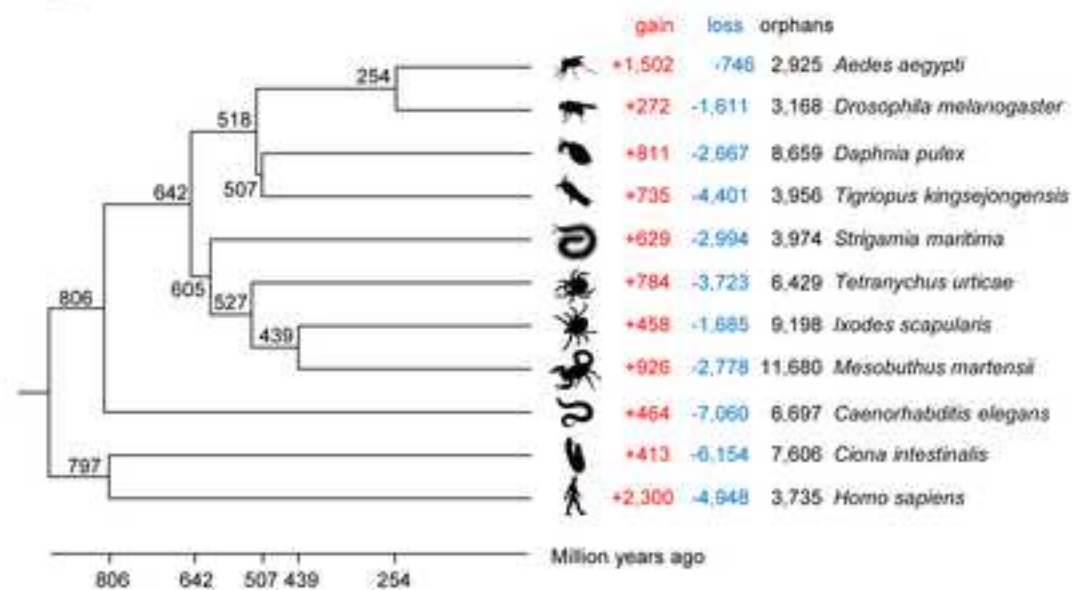
Genes (n)	12,772
Gene Length Sum (bp)	82,293,116
Exons per genes (n)	4.6
mRNA Length Sum (bp)	43,306,342
Average mRNA length (bp)	1,090
Number of tRNA	1,393
Number of rRNA	215

Table 6 Summary of orthologous gene clusters in the 11 representative species.

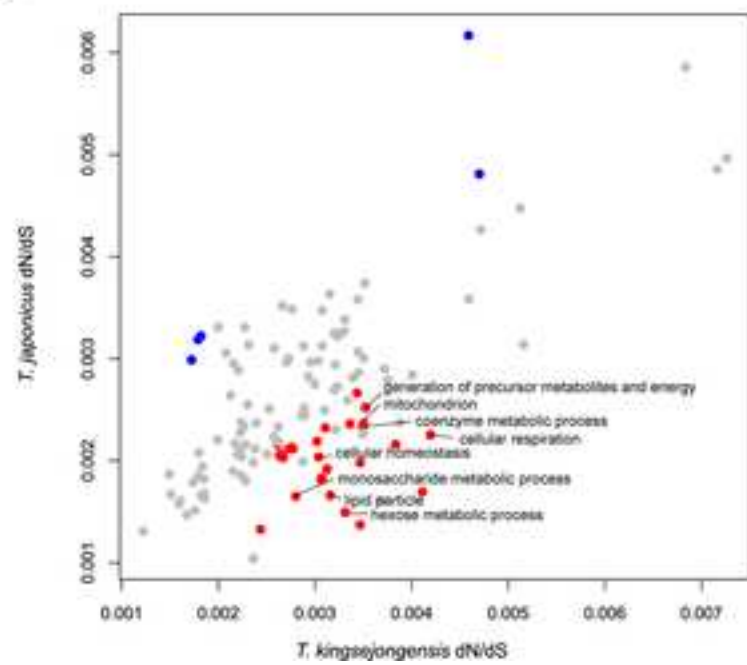
Species	Source of data	No. of coding genes	No. of gene families	No. of genes in gene families	No. of orphan genes	No. of unique gene families	Average No. of genes in gene families
<i>Aedes aegypti</i>	Ensembl genome 25	15,797	7,958	12,792	7,839	854	1.61
<i>Caenorhabditis elegans</i>	Ensembl gene 78	20,447	6,536	13,737	13,911	1,528	2.10
<i>Ciona intestinalis</i>	Ensembl gene 78	16,671	7,017	9,058	9,654	503	1.29
<i>Daphnia pulex</i>	Ensembl genome 25	30,590	6,710	8,362	7,208	368	1.25
<i>Drosophila melanogaster</i>	Ensembl gene 78	13,918	9,673	21,917	20,917	2,408	2.27
<i>Homo sapiens</i>	Ensembl gene 78	20,300	8,696	17,186	11,604	1,065	1.98
<i>Ixodes scapularis</i>	Ensembl genome 25	20,486	8,097	11,277	12,389	873	1.39
<i>Mesobuthus martensii</i>	http://lifecenter.sgst.cn/main/en/scorpion.jsp	32,016	8,389	19,961	23,627	2,276	2.38
<i>Strigamia maritima</i>	Ensembl genome 25	14,992	7,727	11,012	7,265	583	1.43
<i>Tetranychus urticae</i>	Ensembl genome 25	18,224	6,602	11,788	11,622	939	1.79
<i>T. kingsejongensis</i>	this study	12,772	6,205	8,813	6,567	649	1.42



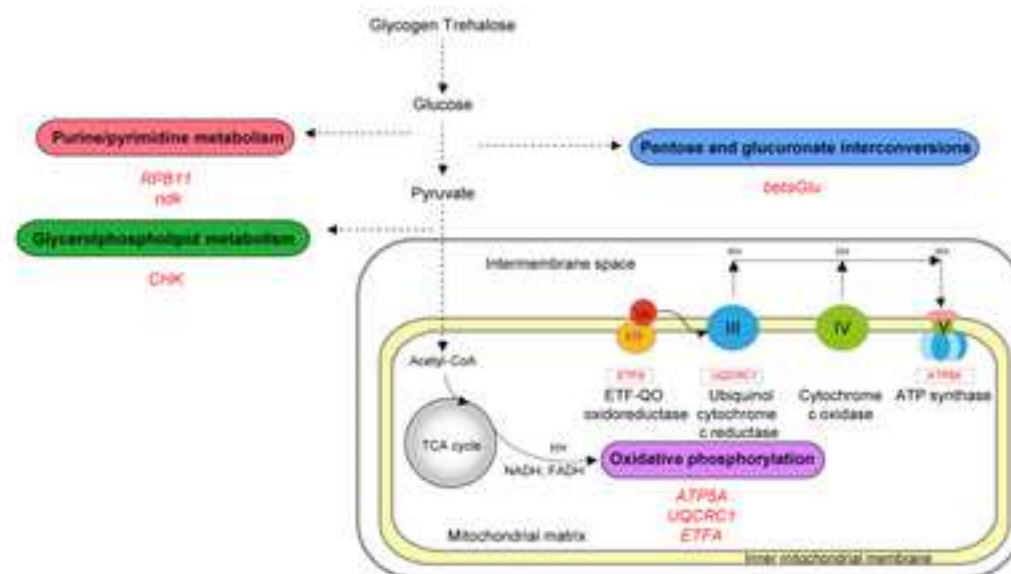


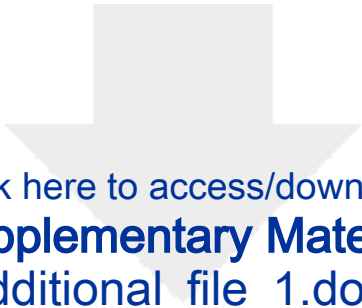
A**B**

A




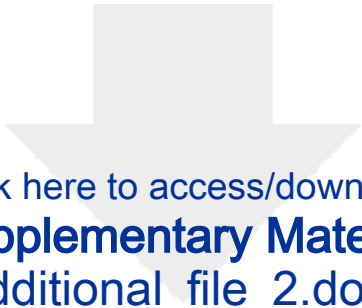
B





Click here to access/download
Supplementary Material
Additional_file_1.docx





Click here to access/download
Supplementary Material
Additional_file_2.docx

