

Genome sequence of Japanese oak silk moth, *Antheraea yamamai*: the first draft genome in family Saturniidae --Manuscript Draft--

Manuscript Number:	GIGA-D-17-00085R3	
Full Title:	Genome sequence of Japanese oak silk moth, <i>Antheraea yamamai</i> : the first draft genome in family Saturniidae	
Article Type:	Data Note	
Funding Information:	Rural Development Administration (PJ010442)	Dr Seong-Ryul Kim
Abstract:	<p>Background <i>Antheraea yamamai</i>, also known as the Japanese oak silk moth, is a wild species of the silk moth. Silk produced by <i>A. yamamai</i>, referred to as tensan silk, shows different characteristics such as thickness, compressive elasticity and chemical resistance compared to the common silk produced from the domesticated silkworm, <i>Bombyx mori</i>. Its unique characteristics have led to its use in many research fields including biotechnology and medical science, and the scientific as well as economic importance of wild silk moth continues to gradually increase. However, no genomic information for wild silk moth, including <i>A. yamamai</i>, is currently available.</p> <p>Findings In order to construct the <i>A. yamamai</i> genome, a total of 147G base pairs using Illumina and Pacbio sequencing platforms were generated, providing 210-fold coverage based on the 700 Mb estimated genome size of <i>A. yamamai</i>. The assembled genome of <i>A. yamamai</i> was 656 Mb(>2kb) with 3,675 scaffolds and the N50 length of assembly was 739 Kb with 34.07% GC ratio. Identified repeat elements covered 37.33% of the total genome and the completeness of the constructed genome assembly was estimated to be 96.7% by BUSCO v2 analysis. A total of 15,481 genes were identified using Evidence Modeler based on the gene prediction results obtained from 3 different methods (ab initio, RNA-seq based, known-gene based) and manual curation.</p> <p>Conclusions Here we present the genome sequence of <i>A. yamamai</i>, the first genome sequence of wild silk moth. These results provide valuable genomic information which will help enrich our understanding of the molecular mechanisms related to not only specific phenotypes such as wild silk itself but also the genomic evolution of Saturniidae.</p>	
Corresponding Author:	Seung-Won Park KOREA, REPUBLIC OF	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Seong-Ryul Kim	
First Author Secondary Information:		
Order of Authors:	Seong-Ryul Kim	
	Woori Kwak	
	Hyaekang Kim	
	Kelsey Caetano-Anolles	
	Kee-Young Kim	
	Su-Bae Kim	
	Kwang-Ho Choi	

	Seong-Wan Kim
	Jae-Sam Hwang
	Min-Jee Kim
	Iksoo Kim
	Tae-Won Goo
	Seung-Won Park
Order of Authors Secondary Information:	
Response to Reviewers:	We attached two separated rebuttals for each reviewer.
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	Yes
Resources 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. Have you included the information requested as detailed in our Minimum Standards Reporting Checklist ?	Yes
Availability of data and materials 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	Yes

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?](#)

1 **Genome sequence of Japanese oak silk moth, *Antheraea yamamai*:**
2
3
4 **the first draft genome in family Saturniidae**
5
6
7

8 **Seong-Ryul Kim^{1†}, Woori Kwak^{2†}, Hyaekang Kim³, Kelsey Caetano-Anolles³, Kee-Young**
9 **Kim¹, Su-Bae Kim¹, Kwang-Ho Choi¹, Seong-Wan Kim¹, Jae-Sam Hwang¹, Min-Jee**
10 **Kim⁴, Iksoo Kim⁴, Tae-Won Goo⁵ and Seung-Won Park^{6*}**
11
12
13
14
15

16 ¹Department of Agricultural Biology, National Academy of Agricultural Science, Rural
17 Development Administration, Wanju-gun 55365, Republic of Korea; ²C&K Genomics, Main
18 Bldg. #420, SNU Research Park, Seoul 151-919, Republic of Korea; ³Department of
19 Agricultural Biotechnology and Research Institute of Agriculture and Life Sciences, Seoul
20 National University, Seoul 151-921, Republic of Korea; ⁴College of Agriculture & Life
21 Sciences, Chonnam National University, Gwangju, Republic of Korea; ⁵Department of
22 Biochemistry, Dongguk University College of Medicine, Gyeongju-si, Gyeongsangbuk-do
23 38066, Republic of Korea; ⁶Department of Biotechnology, Catholic University of Daegu,
24 Gyeongsan-si, Gyeongsangbuk-do 38430, Republic of Korea
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 Seong-Ryul Kim : ksr319@korea.kr; Woori Kwak : asleo@cnkgenomics.com; Hyaekang
43 Kim : hkim458@snu.ac.kr; Kelsey Caetano-Anolles : kelseyca@gmail.com, ORCID: 0000-
44 0001-6310-6679; Kee-Young Kim : applekky@korea.kr; Su-Bae Kim : subae@korea.kr;
45 Kwang-Ho Choi : ckh@korea.kr; Seong-Wan; Seong-Wan Kim : tarupa@korea.kr; Jae-Sam
46 Hwang : hwangjs@korea.kr; Min-Jae Kim : minjeekim3@gmail.com; Iksoo Kim :
47 ikkim81@chonnam.ac.kr; Tae-Won Goo : gootw@dongguk.ac.kr
48
49
50
51
52
53
54
55
56
57

58 † These authors equally contributed and should be regarded as co-first authors.
59
60
61
62
63
64
65

1 * Corresponding authors
2

3 Seung-Won Park
4

5
6 Department of Biotechnology,
7

8
9 Catholic University of Daegu,
10

11
12 Gyeongsan-si, Gyeongsangbuk-do 38430, Republic of Korea
13

14
15 Phone : +82-53-850-3176
16

17
18 Fax : +82-53-359-6846
19

20
21 E-mail: microsw@cu.ac.kr
22

23
24 ORCID: 0000-0002-2218-7748
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

Abstract

Background

Antheraea yamamai, also known as the Japanese oak silk moth, is a wild species of silk moth. Silk produced by *A. yamamai*, referred to as *tensan* silk, shows different characteristics such as thickness, compressive elasticity and chemical resistance compared to common silk produced from the domesticated silkworm, *Bombyx mori*. Its unique characteristics have led to its use in many research fields including biotechnology and medical science, and the scientific as well as economic importance of the wild silk moth continues to gradually increase. However, no genomic information for the wild silk moth, including *A. yamamai*, is currently available.

Findings

In order to construct the *A. yamamai* genome, a total of 147G base pairs using Illumina and Pacbio sequencing platforms were generated, providing 210-fold coverage based on the 700 Mb estimated genome size of *A. yamamai*. The assembled genome of *A. yamamai* was 656 Mb(>2kb) with 3,675 scaffolds and the N50 length of assembly was 739 Kb with a 34.07% GC ratio. Identified repeat elements covered 37.33% of the total genome and the completeness of the constructed genome assembly was estimated to be 96.7% by BUSCO v2 analysis. A total of 15,481 genes were identified using Evidence Modeler based on the gene prediction results obtained from 3 different methods (*ab initio*, RNA-seq based, known-gene based) and manual curation.

Conclusions

Here we present the genome sequence of *A. yamamai*, the first genome sequence of wild silk moth. These results provide valuable genomic information which will help enrich our

1 24 understanding of the molecular mechanisms relating to not only specific phenotypes such as
2
3 25 wild silk itself but also the genomic evolution of Saturniidae.
4
5

6 26 **Keywords**
7
8

9 27 *Antheraea yamamai*, genome assembly, Japanese silk moth, Japanese oak silk moth, wild
10
11 28 silkworm,
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Data description

Antheraea yamamai (NCBI Taxonomy ID: 7121), also known as the Japanese oak silk moth, is a wild silk moth species belonging to the Saturniidae family (Figure 1). Silk moths can be categorized into two families- Bombycidae and Saturniidae. Saturniidae has been estimated to contain approximately 1,861 species with 162 genera [1] and is known as the largest family in the Lepidoptera. Among the many species in family Saturniidae, only a few species, including *A. yamamai*, can be utilized for silk production. Previous phylogenetic studies have shown that family Saturniidae shares common ancestors with family Sphingidae, including the hawk moth (*Macroglossum stellatarum*) and Bombycidae family, including the most representative silkworm, *Bombyx mori* [2]. The estimated divergence time between *A. yamamai* and *B. mori* is 84 MYA (million years ago) making it similar to the 88 MYA estimated divergence time between human and mouse [3, 4].

A. yamamai produces specific silk, called tensan silk [5], which shows distinctive characteristics compared to common silk from *B. mori*, including characteristics such as thickness, bulkiness, compressive elasticity, and resistance to dyeing chemicals [6-8]. These characteristics have received the attention of researchers as a new biomaterial for use in various fields [9-11]. Additionally, it also has been studied for their applications to human health [12-15]. However, despite the potential importance of the wild silk moth in research and economic fields, no whole genomic information is currently available for this or any other species from the family Saturniidae.

In this study, we present the annotated genome sequence of *A. yamamai*, the first published genome in family Saturniidae, with transcriptome datasets collected from 10 different body organ tissues. This data will be a fundamental resource for future studies and provide more insight into the genome evolution and molecular phylogeny of the family Saturniidae.

53 Sequencing

54 For whole genome sequencing, we selected one male sample (Ay-7-male1) from a breeding
55 line (Ay-7) of *A. yamamai* raised at the National Academy of Agricultural Science, Rural
56 Development Administration, Korea. In lepidopterans, males are homogametic (ZZ) and
57 selecting a male sample can reduce the complexity of assembly from excessive repeats on the
58 W chromosome in females. For genomic library construction, we removed the guts of *A.*
59 *yamamai* to prevent contamination of genomes from other organisms such as gut microbes
60 and oak, the main food source of *A. yamamai*. Details of the sample preparation process used
61 in this study are presented in the supplementary information. Genomic DNA was extracted
62 using a DNeasy Animal Mini Kit (Qiagen, Hilden, Germany) and the quality of extracted
63 DNA was checked using trenean, picogreen assay and gel electrophoresis (1% agarose gel/
64 40ng loading). After quality control processing, we were left with a total of 61.5ug of *A.*
65 *yamamai* DNA for genome sequencing. Using standard Illumina whole genome shotgun
66 (WGS) sequencing protocol (paired-end and mate-pair), we added two long read sequencing
67 platforms, Moleclo (Illumina synthetic long read) and RS II (Pacific Bioscience). Table 1-3
68 shows a summary of generated data for each library used in this study. RNA-seq libraries
69 were also constructed for 10 different tissues (Hemocyte, Malpighian tube, Midgut, Fat Body,
70 Anterior-Middle/Silk gland, Posterior/Silk gland, Head, Integument, Testis, Ovary) with 3
71 biological replicates following standard manufacturer protocol (Illumina, San Diego, CA,
72 USA). For this, more than 100 individual *A. yamamai* samples in 5 instar stage from the same
73 breeding line were used for tissue anatomy and 3 samples from each tissue were selected
74 based on the quality of extracted RNA. Details of transcriptome library construction are
75 shown in the supplementary information. Information of libraries and generated data is
76 provided in Table 4, and a total of 147Gb of genomic data and 76Gb of transcriptomic data

1 77 was generated for this study.
2
3

4 78
5
6

7 79 **Genome assembly and evaluation** 8 9

10
11 80 Before conducting genome assembly, we conducted k-mer distribution analysis using a
12
13 81 350bp paired-end library in order to estimate the size and characteristics of the *A. yamamai*
14
15 82 genome. The quality of our generated raw data was checked using FASTQC [16] (FastQC ,
16
17 83 RRID:SCR_014583). Sequencing artifacts such as adapter sequences and low-quality bases
18
19 84 were removed using Trimmomatic (Trimmomatic, RRID:SCR_011848) [17]. Jellyfish [18]
20
21 85 was used to count the k-mer frequency for estimation of the genome size of *A. yamamai*.
22
23 86 Figure 2 shows the 19-mer distribution of *A. yamamai* genome using a 350bp paired-end
24
25 87 library. In the 19-mer distribution, the second peak at approximately half the coverage value
26
27 88 (x-axis) of the main peak indicates heterozygosity. Although the inbred line used in this study
28
29 89 was the single pair sib-mating maintained for more than 10 generations, high heterozygosity
30
31 90 still remains. This phenomenon has been observed in a previous genomic study of the
32
33 91 Diamondback moth (*Plutella xylostella*), and sustained heterozygosity as an important
34
35 92 genomic characteristic was hypothesized to be a result of environmental adaption [19]. The
36
37 93 underlying mechanism of these sustained heterozygosity is unclear, but associative
38
39 94 overdominance can be one of the candidate explanation of this phenomenon [20, 21]. Based
40
41 95 on the result of 19-mer distribution analysis, the genome size of *A. yamamai* was estimated to
42
43 96 be 709Mb. However, this size might be larger than the real genome size of *A.yamamai*
44
45 97 because high heterozygosity could affect the estimation of genome size based on the K-mer
46
47 98 distribution. Next, we conducted error correction on Illumina paired-end libraries using the
48
49 99 error correction module of Allpaths-LG [22] before the initial contig assembly process
50
51 100 (ALLPATHS-LG , RRID:SCR_010742). After error correction, initial contig assembly with
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 101 350bp and 700bp libraries was conducted using SOAP denovo2 [23] with the parameter
2
3 102 option set at K=19; this approach showed the best assembly statistics compared to other
4
5 103 assemblers and parameters (SOAPdenovo2 , RRID:SCR_014986). Quality control processing
6
7
8 104 for mate-pair libraries and scaffolding was conducted using Nxtrim [24] and SSPACE
9
10 105 (SSPACE , RRID:SCR_011848) [25], respectively. At each scaffolding step, SOAP
11
12 106 Gapcloser [23] with -l 155 and -p 31 parameters was repeatedly used to close the gaps within
13
14 107 each scaffold. In order to obtain a higher quality genome assembly of *A. yamamai*, we
15
16 108 employed several long read scaffolding strategies using SSPACE-LongRead [26]. First, we
17
18 109 used an Illumina synthetic long read sequencing platform called Moleculo which has been
19
20 110 proven valuable for the study of highly heterozygous genomes in previous studies [27, 28].
21
22 111 After scaffolding was performed using SSPACE-LongRead with Illumina synthetic long read
23
24 112 data, the total number of assembled scaffolds was effectively reduced from 398,446 to 24,558.
25
26 113 The average scaffold length was also extended from 1.7 Kb to 24.8 Kb. However, there was
27
28 114 no impressive improvement in N50 length (approximately 91 Kb to 112 Kb) of assembled
29
30 115 scaffolds. Therefore, we employed another type of long read data generated from 10 cells of
31
32 116 Pacbio RS II system with P6-C4 chemistry. After final scaffolding processing using Pacbio
33
34 117 long reads, the number of scaffolds was reduced to 3,675 and N50 length was effectively
35
36 118 extended from 112 Kb to 739 Kb. Summary statistics of the assembled *A. yamamai* genome
37
38 119 is provided in Table 5. Final assembly of the *A. yamamai* genome was 656 Mb (>2kb) long
39
40 120 with 3,675 scaffolds and the N50 length of assembly was 739 Kb with a 34.07% GC ratio. To
41
42 121 evaluate the quality of the assembled genome, we conducted BUSCO (Benchmarking
43
44 122 Universal Single-Copy Orthologs) analysis [29] using BUSCO v2.0 with insecta_odb9
45
46 123 including 1,658 BUSCOs from 42 species (BUSCO , RRID:SCR_015008). From BUSCO
47
48 124 analysis, 96.7% of BUSCOs were completely detected in the assembled genome (1,576 :
49
50 125 complete and single-copy, 27 : complete and duplicated) among 1,658 tested BUSCOs. The
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 126 number of fragmented and missing BUSCOs was 21 and 34, respectively. Based on the result
2
3 127 of BUSCO analysis, the genome of *A.yamamai* presented here was considered properly
4
5
6 128 constructed for downstream analysis.
7
8
9 129

130 **Repeat identification and comparative repeat analysis**

131 To identify repeat elements of the *A. yamamai* genome, a custom repeat library was
132 constructed using RepeatModeler with RECON [30], RepeatScout [31] and TRF [32]. The
133 resulting constructed custom repeat library for *A. yamamai* was further curated using
134 CENSOR [33] search and the curated library was employed in RepeatMasker [34] with
135 Repbase [35]. RepeatMasker (RepeatMasker, RRID:SCR_012954) was conducted with
136 RMBlast and ‘no_is’ option for skipping bacterial insertion element check. Table 6
137 summarizes the proportion of identified mobile elements in the *A. yamamai* genome. The
138 most prevalent repeat elements in the *A. yamamai* genome were LINE element (101 Mb,
139 15.31% of total genome) and total repeat elements accounted for 37.33% of the total genome.
140 In order to compare the repeat elements of *A. yamamai* with that of other genomes, we
141 conducted the same process for seven public genomes which are close neighbors of *A.*
142 *yamamai* - *Aedes aegypti* [36], *Bombyx mori* [37], *Danaus plexippus* [38], *Drosophila*
143 *melanogaster* [39], *Heliconius Melpomene* [40], *Melitaea cinxia* [41] and *Plutella xylostella*
144 [19]. Figure 3 displays the amount and proportion of identified repeat elements from the 8
145 species. Despite the small genome size of *B. mori*, the total amount of identified SINE
146 element in the *B. mori* genome was 5.77 times larger than that of *A. yamamai*. The top 5
147 expanded repeat elements in *A. yamamai* genome were DNA/RC, LINE/L2, LINE/RTE-
148 BovB, DNA/TcMar-Mariner and LINE/CR1. Among these, DNA/TcMar-Mariner was the
149 specifically expanded repeat element in *A. yamamai* among 8 species. In *B. mori*,

1 150 SINE/tRNA-CR1, LINE/Jockey, DNA/RC, LINE/CR1-Zenon and LINE/RTE-BovB were the
2
3 151 top 5 expanded repeat elements. When comparing the repeat elements of *A. yamamai* with
4
5
6 152 those of *B. mori*, which are both producers of the same type of silk, repeat elements showed
7
8 153 family and species-specific patterns in the two silk moth lineages. Particularly, we found that
9
10 154 the mariner repeat element, which was found specifically expanded in the *A. yamamai*
11
12
13 155 genome, was also included in the fibroin gene. A previous sequencing study also showed that
14
15
16 156 the mariner repeat element was inserted in the 5'-end of fibroin gene of *A. yamamai* [42].
17
18 157 Fibroin is the core component of the silk protein found in silk moth, and the physical
19
20 158 characteristics of silk mainly depend on the types and unique repeat motif of the fibroin [43].
21
22
23 159 This gene is known to have hundreds of tandem repeat motifs and these kinds of tandem
24
25 160 repeats can be derived through transposable elements. This indicates that the mariner repeat
26
27
28 161 element, specifically expanded in the *A. yamamai* genome, may play an important role in
29
30 162 development of the unique silk of *A. yamamai*, and the lineage-specific repeat elements may
31
32
33 163 be one of the candidate evolution forces related to host-specific phenotype during genome
34
35 164 evolution.

36 37 38 165 39 40 41 42 166 **Gene prediction and annotation**

43
44
45 167 Three different algorithms were used for gene prediction of the *A. yamamai* genome: *ab initio*,
46
47
48 168 RNA-seq transcript based, and protein homology-based approaches. For *ab initio* gene
49
50
51 169 prediction, Augustus (Augustus: Gene Prediction, RRID:SCR_008417)[44], Geneid [45] and
52
53
54 170 GeneMarks-ET [46] were employed. Augustus was trained using known genes of *A. yamamai*
55
56 171 in NCBI database and mapping information of RNA-seq data obtained from Tophat [47]
57
58 172 (TopHat , RRID:SCR_013035) was also utilized for gene prediction. Geneid was used with
59
60
61
62
63
64
65

1 173 predefined parameters for *Drosophila melanogaster*. GeneMarks-ET was employed using
2
3 174 junction information of genes from transcriptome data alignment. For RNA-seq transcript
4
5
6 175 based prediction, generated transcriptome data from ten organ tissues of *A. yamamai* were
7
8 176 aligned to the assembled genome and gene information was predicted using Cufflinks [48]
9
10
11 177 (Cufflinks , RRID:SCR_014597). The longest CDS sequences were identified from Cufflinks
12
13 178 results using Transdecoder. For the homology-based approach, all known genes of order
14
15
16 179 Lepidoptera in the NCBI database were aligned using PASA [49]. Table 7 shows the gene
17
18 180 prediction results from each method. Gene prediction results from different prediction
19
20
21 181 algorithms were combined using EVM (Evidence Modeler) [50] and a consensus gene set of
22
23 182 the *A. yamamai* genome was created. Manual curation was performed based on the 5
24
25 183 evidences (3 in-silico, known protein and RNA-seq) using IGV [51] and Blastp (BLASTP,
26
27
28 184 RRID:SCR_001010). Using IGV with each gene evidence and comparing results with known
29
30
31 185 genes via blastp, we mainly focused on the removing false positively predicted genes which
32
33 186 don't have enough evidences. And merged and spliced gene structured were corrected by
34
35 187 comparing the gene structure with known exon structure in NCBI NR database. In addition,
36
37
38 188 fibroin and sericin genes which couldn't be properly predicted because of its high repeat
39
40 189 motif were also manually identified with previously known sequences [42, 52] with RNA-seq
41
42
43 190 data. The final gene set of *A. yamamai* genome contains 15,481 genes. Summary statistics for
44
45 191 the consensus gene set is provided in Table 8. The average gene length was 11,016.34 bp with
46
47
48 192 a 34.38% GC ratio and the number of exons per gene was 5.64. In order to identify the
49
50 193 function of predicted genes in *A. yamamai*, three non-redundant sequence databases (Swiss-
51
52 194 Prot [53], Uniref100 [53], and NCBI NR [54]) as well as the gene information of two species
53
54
55 195 (*B. mori* and *D. melanogaster*) were used for target databases using Blastp. Additionally,
56
57 196 protein domain searches were conducted on the consensus gene set using InterproScan5
58
59
60 197 (InterProScan, RRID:SCR_005829) [55]. Figure S1 shows the top 20 identified terms from 7
61
62
63
64
65

1 198 different InterproScan5 analyses. Among the various analysis conducted using InterproScan5,
2
3 199 gene ontology analysis with Pfam database showed that a large proportion of genes in the
4
5 200 *A.yamamai* genome were related with the function of molecular binding, catalytic activity,
6
7
8 201 internal component of membrane, metabolic process, oxidation-reduction process and
9
10 202 transmembrane transport.
11
12
13
14 203
15
16
17

18 204 **Comparative genome analysis**

19
20
21

22 205 We used OrthoMCL [56] and RBH(Reciprocal Best Hit) within blastp for identification of
23
24 206 gene family clusters and 1:1 orthologous gene sets. Gene information of 7 taxa (*A. aegypti*, *B.*
25
26 207 *mori*, *D. plexippus*, *D. melanogaster*, *H. melpomene*, *M. cinxia* and *P. xylostella*), same taxa
27
28
29 208 used in repeat analysis, was employed for OrthoMCL with *A. yamamai*. A total of 17,406
30
31 209 gene family clusters were constructed and 3,586 1:1 orthologous genes were identified.
32
33
34 210 Before conducting comparative genome analysis, we constructed phylogenetic trees for the 8
35
36 211 species. In order to build the phylogenetic tree, multiple sequence alignment for the 1:1
37
38
39 212 orthologous genes of all 8 species was conducted using PRANK [57], and Gblocks [58] was
40
41 213 used to obtain conserved blocks for the phylogenetic tree. Conserved block sequences were
42
43
44 214 sequentially concatenated to obtain one consensus sequence for each species. MEGA [59]
45
46 215 was used for constructing Neighbor-Joining Trees (bootstrap 1000, maximum composite
47
48 216 likelihood, transitions + transversions, and gamma distributed option) and MrBayes
49
50
51 217 (MrBayes, RRID:SCR_012067) [60] was employed for the construction of Bayesian
52
53 218 inference trees. To select the best evolution model for our data, Modeltest [61] was conducted
54
55
56 219 and the GTR based invariant model was chosen based on the AIC value of Modeltest. Figure
57
58 220 4 shows the constructed phylogenetic tree of the 8 species using 3,586 orthologous genes.
59
60
61
62
63
64
65

1 221 The bootstrap value and Bayesian poster probability value of all nodes were 100 and 1,
2
3 222 respectively. The closest neighbor of *A. yamamai* was *B. mori*, which is included in
4
5
6 223 Bombycidae family; this result is consistent with that of previous studies. Three butterfly
7
8 224 species (*D. plexippus*, *M. cinxia* and *H. meplmene*) included in Nymphalidae family were
9
10
11 225 also shown to share a common ancestor with families Saturniidae and Bombycidae.

12
13
14 226 Based on the constructed phylogenetic tree, gene family expansion and contraction analysis
15
16 227 was conducted using a 2 parameter model in CAFÉ [62] and the gene tree was constructed
17
18 228 using protein sequence via MEGA [59]. Figure 4 shows the result of gene family expansion
19
20
21 229 and contraction analysis of 8 species. 938 and 1,987 gene families of *A. yamamai* and 567
22
23 230 and 715 gene families of *B. mori* were estimated to be expanded and contracted from the
24
25
26 231 common ancestors, respectively. Among these, 15 gene families in *A. yamamai* were
27
28 232 estimated to be under rapid expansion during the evolution process. Functions of genes in
29
30
31 233 rapidly expanded gene families of *A. yamamai* were transposase, fatty acid synthase, zinc
32
33 234 finger protein, chorion (eggshell protein), reverse transcriptase, prostaglandin dehydrogenase,
34
35
36 235 RNA-directed DNA polymerase, gag like protein, juvenile hormone acid methyltransferase,
37
38 236 facilitated trehalose transporter and glucose dehydrogenase. Figure 5 shows the gene tree of
39
40
41 237 two chorion gene (chorion class A and B) family clusters rapidly expanded in the *A. yamamai*
42
43 238 genome. Chorion, called eggshell protein, composes the surface of egg and protects the
44
45
46 239 embryo from environmental threats such as desiccation, flooding, freezing, infection of
47
48 240 microorganisms, and physical destruction. It also provides channels, such as aeropyle, which
49
50
51 241 enables gas exchange and maintains proper condition for diapause egg [63]. These diverse
52
53 242 functions of eggshell are implemented by the specific eggshell structure and the surface
54
55 243 structure of eggshell varies between species for the adaptation in a different environment. The
56
57
58 244 ancestor of *Antheraea* has the unique aeropyle structure called “aerophyle crown” on the
59
60 245 eggshell surface [64]. This unique structure is formed by the circular vertical projection of
61
62
63
64
65

1 246 lamellar chorion from follicle cell and it surrounds the aeropyles near the end of oogenesis
2
3 247 [65]. Acquiring this kind of *de novo* complex structure requires numerous genetic changes
4
5
6 248 and a previous study about *Antheraea Polyphemus* has shown that over a hundred chorion
7
8 249 specific polypeptides were involved for this unique ultra-structure [65]. Therefore, the
9
10
11 250 specific rapid expansion of chorion class A and B gene family in *A. yamamai* genome might
12
13 251 be one of the convincing molecular explanation for acquiring this unique ultrastructure in the
14
15
16 252 eggshell surface of *Antheraea* genus. However, this unique ultra-structure tends to be reduced
17
18 253 during current evolution process of the *Antheraea* genus. Types of eggshell structure in
19
20
21 254 *Antheraea* genus can be categorized into multiple classes based on the morphology and
22
23 255 regional distribution of aeropyle [64]. The shape of aeropyle in *A. yamamai* egg is known to
24
25 256 be converted to mound shape from the crown shape and these aeropyle mounds only exist in
26
27
28 257 the narrow band surrounding the micropyle region [64]. Only a very few, small aeropyle
29
30 258 crowns remained and it is entirely different with the ancestral form of eggshell surface mostly
31
32
33 259 covered by aeropyle crowns. These regional differences were known to be adjusted by
34
35 260 regional difference of filler genes during choriogenesis [66] and the additional regulations of
36
37
38 261 related genes for choriogenesis have to be considered. This indicates that specifically
39
40 262 expanded chorion gene families of *A. yamamai* may be one of the remaining evolutionary
41
42 263 tracks in the genome of *Antheraea* genus. However, further functional studies must be
43
44
45 264 conducted to resolve the limited understanding about the relationship between these
46
47 265 expanded chorion gene families and the current eggshell surface formation of *A. yamamai*.

48
49
50 266 The constructed genome of *A.yamamai* presented here is the first announced genome in
51
52 267 family Saturniidae and the karyotyping analysis using gamete in metaphase showed that the
53
54
55 268 genome of *A. yamamai* consists of 31 chromosomes (Figure 6). This constructed genome
56
57
58 269 information provides more insight into the genome evolution and phylogeny of family
59
60 270 Saturniidae, which contains the largest number of species in Lepidoptera. For example,
61
62
63
64
65

1 271 although two silk moths, *A. yamamai* and *B. mori*, appear similar, comparative genome
2
3
4 272 analysis showed the significant differences in the genome size, specific expansion of repeat
5
6 273 elements and gene families between families Saturniidae and Bombycidae. In case of
7
8 274 molecular phylogeny, most previous phylogenetic studies were limited to few genes due to
9
10
11 275 the lack of genomic information on family Saturniidae. We expect our study and resulting
12
13 276 constructed genome will resolve some limitations of molecular phylogenetic and ecological
14
15 277 research on Saturniidae species. Additionally, constructed genome information will help
16
17
18 278 researchers better understand the molecular background of wild silk and its production. Silk
19
20
21 279 produced by *A. yamamai*, referred to as *tensan* silk, shows unique characteristics which have
22
23 280 made it valuable in various fields. However, *A. yamamai* has not been completely
24
25 281 domesticated compared to *B. mori*, making mass production of *tensan* silk infeasible.
26
27
28 282 Understanding of the molecular mechanisms behind the *tensan* silk production process is
29
30 283 essential for mass production using biotechnology, and this genome sequence with manually
31
32 284 curated gene information is a fundamental resource for related research and industrial
33
34
35 285 improvement. Additionally, the transcriptome data of 10 different organ tissues with 3
36
37 286 biological replications presented here may be also useful resources for uncovering the
38
39
40 287 molecular mechanisms related to specific phenotypes of *A.yamamai* and family Saturniidae.
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

288 **Availability of supporting data**

289 The generated genome sequence and gene information of *A. yamamai* are available in
290 GigaDB [67] and generated raw data is available under project accession PRJNA383008 and
291 PRJNA383025 of the NCBI database.

292 **Competing interests**

293 All authors report no competing interests.

294 **Abbreviation**

295 RBH – Reciprocal Best Hit

296 **Authors contributions**

297 Sampling - Kee-Young Kim, Su-Bae Kim

298 Sequencing - Kwang-Ho Choi, Seong-Wan Kim

299 Genome assembly - Seong-Ryul Kim, Woori Kwak, Jae-Sam Hwang, Seung-Won Park

300 Repeat element analysis - Seong-Ryul Kim, Woori Kwak, Seung-Won Park

301 Gene prediction - Seong-Ryul Kim, Woori Kwak, Hyaekang Kim, Jae-Sam Hwang

302 Comparative genome analysis - Seong-Ryul Kim, Woori Kwak, Min-Jae Kim, Kelsey

303 Caetano-Anolles

304 Funding and experimental design - Seong-Ryul Kim, Seung-Won Park

305

1 306 **Acknowledgements**

2
3
4 307 This work was supported by a grant from the Rural Development Administration, Republic of
5
6 308 Korea (grant no. PJ010442).

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

References

1. Regier, J.C., M.C. Grant, C. Mitter, et al., *Phylogenetic relationships of wild silkmoths (Lepidoptera: Saturniidae) inferred from four protein-coding nuclear genes*. Systematic Entomology, 2008. **33**(2): p. 219-228.
2. Regier, J.C., C. Mitter, A. Zwick, et al., *A large-scale, higher-level, molecular phylogenetic study of the insect order Lepidoptera (moths and butterflies)*. PLoS One, 2013. **8**(3): p. e58568.
3. Hedges, S.B., J. Dudley, and S. Kumar, *TimeTree: a public knowledge-base of divergence times among organisms*. Bioinformatics, 2006. **22**(23): p. 2971-2972.
4. Kawahara, A.Y. and J.R. Barber, *Tempo and mode of antibat ultrasound production and sonar jamming in the diverse hawkmoth radiation*. Proceedings of the National Academy of Sciences, 2015. **112**(20): p. 6407-6412.
5. Peigler, R.S., *Wild silks of the world*. American Entomologist, 1993. **39**(3): p. 151-162.
6. MATSUMOTO, Y.-I. and H. SAITO, *Load-extension characteristics of composite raw silk of *Antheraea yamamai* and *Bombyx mori**. The Journal of Sericultural Science of Japan, 1997. **66**(6): p. 497-501.
7. Nakamura, S., Y. Saegusa, Y. Yamaguchi, et al., *Physical properties and structure of silk. XI. Glass transition temperature of wild silk fibroins*. Journal of applied polymer science, 1986. **31**(3): p. 955-956.
8. Kweon, H. and Y. Park, *Structural characteristics and physical properties of wild silk fibres: *Antheraea pernyi* and *Antheraea yamamai**. Korean Journal of Sericultural Science (Korea Republic), 1994.
9. Zheng, Z., Y. Wei, S. Yan, et al., *Preparation of regenerated *Antheraea yamamai* silk fibroin film and controlled-molecular conformation changes by aqueous ethanol treatment*. Journal of applied polymer science, 2010. **116**(1): p. 461-467.
10. Omenetto, F., D. Kaplan, J. Amsden, et al., *Silk based biophotonic sensors*. 2011, Google Patents.
11. Takeda, S., *New field of insect science: Research on the use of insect properties*. Entomological Science, 2013. **16**(2): p. 125-135.
12. Omenetto, F. and D.L. Kaplan, *Silk-based multifunctional biomedical platform*. 2012, Google Patents.
13. Serban, M.A., *Silk medical devices*. 2016, Google Patents.
14. Jiang, G.-L., A.L. Collette, R.L. Horan, et al., *Drug delivery platforms comprising silk fibroin hydrogels and uses thereof*. 2010, Google Patents.
15. Kamiya, M., K. Oyauchi, Y. Sato, et al., *Structure-activity relationship of a novel pentapeptide with cancer cell growth-inhibitory activity*. Journal of Peptide Science, 2010. **16**(5): p. 242-248.
16. Bioinformatics, B., *FastQC A quality control tool for high throughput sequence data*.

- 1 Cambridge, UK: Babraham Institute, 2011.
- 2
- 3 17. Bolger, A.M., M. Lohse, and B. Usadel, *Trimmomatic: a flexible trimmer for Illumina*
- 4 *sequence data*. Bioinformatics, 2014: p. btu170.
- 5
- 6 18. Marçais, G. and C. Kingsford, *A fast, lock-free approach for efficient parallel counting of*
- 7 *occurrences of k-mers*. Bioinformatics, 2011. **27**(6): p. 764-770.
- 8
- 9 19. You, M., Z. Yue, W. He, et al., *A heterozygous moth genome provides insights into*
- 10 *herbivory and detoxification*. Nature genetics, 2013. **45**(2): p. 220-225.
- 11
- 12 20. Maruyama, T. and M. Nei, *Genetic variability maintained by mutation and overdominant*
- 13 *selection in finite populations*. Genetics, 1981. **98**(2): p. 441-459.
- 14
- 15 21. Pamilo, P. and S. Pálsson, *Associative overdominance, heterozygosity and fitness*. Heredity,
- 16 1998. **81**(4): p. 381-389.
- 17
- 18 22. Gnerre, S., I. MacCallum, D. Przybylski, et al., *High-quality draft assemblies of mammalian*
- 19 *genomes from massively parallel sequence data*. Proceedings of the National Academy of
- 20 Sciences, 2011. **108**(4): p. 1513-1518.
- 21
- 22 23. Luo, R., B. Liu, Y. Xie, et al., *SOAPdenovo2: an empirically improved memory-efficient*
- 23 *short-read de novo assembler*. Gigascience, 2012. **1**(1): p. 18.
- 24
- 25 24. O'Connell, J., O. Schulz-Trieglaff, E. Carlson, et al., *NxTrim: optimized trimming of Illumina*
- 26 *mate pair reads*. Bioinformatics, 2015. **31**(12): p. 2035-2037.
- 27
- 28 25. Boetzer, M., C.V. Henkel, H.J. Jansen, et al., *Scaffolding pre-assembled contigs using*
- 29 *SSPACE*. Bioinformatics, 2011. **27**(4): p. 578-579.
- 30
- 31 26. Boetzer, M. and W. Pirovano, *SSPACE-LongRead: scaffolding bacterial draft genomes using*
- 32 *long read sequence information*. BMC bioinformatics, 2014. **15**(1): p. 211.
- 33
- 34 27. Voskoboynik, A., N.F. Neff, D. Sahoo, et al., *The genome sequence of the colonial chordate,*
- 35 *Botryllus schlosseri*. Elife, 2013. **2**: p. e00569.
- 36
- 37 28. McCoy, R.C., R.W. Taylor, T.A. Blauwkamp, et al., *Illumina TruSeq synthetic long-reads*
- 38 *empower de novo assembly and resolve complex, highly-repetitive transposable elements*.
- 39 *PloS one*, 2014. **9**(9): p. e106689.
- 40
- 41 29. Simão, F.A., R.M. Waterhouse, P. Ioannidis, et al., *BUSCO: assessing genome assembly and*
- 42 *annotation completeness with single-copy orthologs*. Bioinformatics, 2015: p. btv351.
- 43
- 44 30. Bao, Z. and S.R. Eddy, *Automated de novo identification of repeat sequence families in*
- 45 *sequenced genomes*. Genome Research, 2002. **12**(8): p. 1269-1276.
- 46
- 47 31. Price, A.L., N.C. Jones, and P.A. Pevzner, *De novo identification of repeat families in large*
- 48 *genomes*. Bioinformatics, 2005. **21**(suppl 1): p. i351-i358.
- 49
- 50 32. Benson, G., *Tandem repeats finder: a program to analyze DNA sequences*. Nucleic acids
- 51 *research*, 1999. **27**(2): p. 573.
- 52
- 53 33. Kohany, O., A.J. Gentles, L. Hankus, et al., *Annotation, submission and screening of*
- 54 *repetitive elements in Repbase: RepbaseSubmitter and Censor*. BMC bioinformatics, 2006.
- 55 **7**(1): p. 474.
- 56
- 57 34. Tarailo-Graovac, M. and N. Chen, *Using RepeatMasker to identify repetitive elements in*
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1 *genomic sequences*. Current Protocols in Bioinformatics, 2009: p. 4.10. 1-4.10. 14.
- 2
- 3 35. Bao, W., K.K. Kojima, and O. Kohany, *Rebase Update, a database of repetitive elements in*
- 4 *eukaryotic genomes*. Mobile DNA, 2015. **6**(1): p. 11.
- 5
- 6 36. Nene, V., J.R. Wortman, D. Lawson, et al., *Genome sequence of Aedes aegypti, a major*
- 7 *arbovirus vector*. Science, 2007. **316**(5832): p. 1718-1723.
- 8
- 9 37. Xia, Q., Z. Zhou, C. Lu, et al., *A draft sequence for the genome of the domesticated*
- 10 *silkworm (Bombyx mori)*. Science, 2004. **306**(5703): p. 1937-1940.
- 11
- 12 38. Zhan, S., C. Merlin, J.L. Boore, et al., *The monarch butterfly genome yields insights into*
- 13 *long-distance migration*. Cell, 2011. **147**(5): p. 1171-1185.
- 14
- 15 39. Adams, M.D., S.E. Celniker, R.A. Holt, et al., *The genome sequence of Drosophila*
- 16 *melanogaster*. Science, 2000. **287**(5461): p. 2185-2195.
- 17
- 18 40. Consortium, H.G., *Butterfly genome reveals promiscuous exchange of mimicry adaptations*
- 19 *among species*. Nature, 2012. **487**(7405): p. 94-98.
- 20
- 21 41. Ahola, V., R. Lehtonen, P. Somervuo, et al., *The Glanville fritillary genome retains an ancient*
- 22 *karyotype and reveals selective chromosomal fusions in Lepidoptera*. Nature
- 23 *communications*, 2014. **5**.
- 24
- 25 42. Hwang, J.-S., J.-S. Lee, T.-W. Goo, et al., *Cloning of the fibroin gene from the oak silkworm,*
- 26 *Antheraea yamamai and its complete sequence*. Biotechnology letters, 2001. **23**(16): p.
- 27 1321-1326.
- 28
- 29 43. Malay, A.D., R. Sato, K. Yazawa, et al., *Relationships between physical properties and*
- 30 *sequence in silkworm silks*. Scientific reports, 2016. **6**: p. 27573.
- 31
- 32 44. Stanke, M., M. Diekhans, R. Baertsch, et al., *Using native and syntenically mapped cDNA*
- 33 *alignments to improve de novo gene finding*. Bioinformatics, 2008. **24**(5): p. 637-644.
- 34
- 35 45. Blanco, E., G. Parra, and R. Guigó, *Using geneid to identify genes*. Current protocols in
- 36 *bioinformatics*, 2007: p. 4.3. 1-4.3. 28.
- 37
- 38 46. Lomsadze, A., P.D. Burns, and M. Borodovsky, *Integration of mapped RNA-Seq reads into*
- 39 *automatic training of eukaryotic gene finding algorithm*. Nucleic acids research, 2014: p.
- 40 gku557.
- 41
- 42 47. Trapnell, C., L. Pachter, and S.L. Salzberg, *TopHat: discovering splice junctions with RNA-*
- 43 *Seq*. Bioinformatics, 2009. **25**(9): p. 1105-1111.
- 44
- 45 48. Trapnell, C., A. Roberts, L. Goff, et al., *Differential gene and transcript expression analysis of*
- 46 *RNA-seq experiments with TopHat and Cufflinks*. Nature protocols, 2012. **7**(3): p. 562-578.
- 47
- 48 49. Campbell, M.A., B.J. Haas, J.P. Hamilton, et al., *Comprehensive analysis of alternative*
- 49 *splicing in rice and comparative analyses with Arabidopsis*. BMC genomics, 2006. **7**(1): p.
- 50 327.
- 51
- 52 50. Haas, B.J., S.L. Salzberg, W. Zhu, et al., *Automated eukaryotic gene structure annotation*
- 53 *using EVIDENCEModeler and the Program to Assemble Spliced Alignments*. Genome
- 54 *biology*, 2008. **9**(1): p. R7.
- 55
- 56 51. Robinson, J.T., H. Thorvaldsdóttir, W. Winckler, et al., *Integrative genomics viewer*. Nature
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1 biotechnology, 2011. **29**(1): p. 24-26.
- 2
- 3 52. Zurovec, M., N. Yonemura, B. Kludkiewicz, et al., *Sericin Composition in the Silk of*
- 4 *Antheraea yamamai*. *Biomacromolecules*, 2016. **17**(5): p. 1776-1787.
- 5
- 6 53. Consortium, U., *Reorganizing the protein space at the Universal Protein Resource (UniProt)*.
- 7 *Nucleic acids research*, 2011: p. gkr981.
- 8
- 9 54. Pruitt, K.D., T. Tatusova, and D.R. Maglott, *NCBI reference sequences (RefSeq): a curated*
- 10 *non-redundant sequence database of genomes, transcripts and proteins*. *Nucleic acids*
- 11 *research*, 2007. **35**(suppl 1): p. D61-D65.
- 12
- 13 55. Jones, P., D. Binns, H.-Y. Chang, et al., *InterProScan 5: genome-scale protein function*
- 14 *classification*. *Bioinformatics*, 2014. **30**(9): p. 1236-1240.
- 15
- 16 56. Li, L., C.J. Stoeckert, and D.S. Roos, *OrthoMCL: identification of ortholog groups for*
- 17 *eukaryotic genomes*. *Genome research*, 2003. **13**(9): p. 2178-2189.
- 18
- 19 57. Löytynoja, A. and N. Goldman, *An algorithm for progressive multiple alignment of*
- 20 *sequences with insertions*. *Proceedings of the National Academy of Sciences of the United*
- 21 *States of America*, 2005. **102**(30): p. 10557.
- 22
- 23 58. Castresana, J., *Selection of conserved blocks from multiple alignments for their use in*
- 24 *phylogenetic analysis*. *Molecular biology and evolution*, 2000. **17**(4): p. 540-552.
- 25
- 26 59. Kumar, S., G. Stecher, and K. Tamura, *MEGA7: Molecular Evolutionary Genetics Analysis*
- 27 *version 7.0 for bigger datasets*. *Molecular biology and evolution*, 2016. **33**(7): p. 1870-1874.
- 28
- 29 60. Ronquist, F. and J.P. Huelsenbeck, *MrBayes 3: Bayesian phylogenetic inference under mixed*
- 30 *models*. *Bioinformatics*, 2003. **19**(12): p. 1572-1574.
- 31
- 32 61. Posada, D., *Using MODELTEST and PAUP* to select a model of nucleotide substitution*.
- 33 *Current protocols in bioinformatics*, 2003: p. 6.5. 1-6.5. 14.
- 34
- 35 62. De Bie, T., N. Cristianini, J.P. Demuth, et al., *CAFE: a computational tool for the study of*
- 36 *gene family evolution*. *Bioinformatics*, 2006. **22**(10): p. 1269-1271.
- 37
- 38 63. Chapman, R.F., *The insects: structure and function*. 1998: Cambridge university press.
- 39
- 40 64. Regier, J.C., U. Paukstadt, L.H. Paukstadt, et al., *Phylogenetics of eggshell morphogenesis in*
- 41 *Antheraea (Lepidoptera: Saturniidae): unique origin and repeated reduction of the*
- 42 *aeropyle crown*. *Systematic biology*, 2005. **54**(2): p. 254-267.
- 43
- 44 65. Regier, J.C., G.D. Mazur, and F.C. Kafatos, *The silkmoth chorion: morphological and*
- 45 *biochemical characterization of four surface regions*. *Developmental biology*, 1980. **76**(2): p.
- 46 286-304.
- 47
- 48 66. Hatzopoulos, A.K. and J.C. Regier, *Evolutionary changes in the developmental expression of*
- 49 *silkmoth chorion genes and their morphological consequences*. *Proceedings of the*
- 50 *National Academy of Sciences*, 1987. **84**(2): p. 479-483.
- 51
- 52 67. Kim, S; Kwak, W; Kim, K; Kim, S; Choi, K; Kim, S; Hwang, J; Kim, I; Goo, T; Park, S (2017):
- 53 *The Japanese silk moth, Antheraea yamamai, draft genome sequence* GigaScience
- 54 *Database*. <http://dx.doi.org/10.5524/100382>
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

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

Tables

Table 1. Summary statistics of generated whole genome shotgun sequencing data using Illumina Nextseq 500.

Library Name	Library Type	Insert Size	Platform	Read Length	No. Reads	Total Base(bp)	Reads retained after trimming
350bp	Paired-end	350bp	Nextseq500	151	293,176,268	44,269,616,468	291,070,362
700bp	Paired-end	700bp	Nextseq500	151	246,945,900	37,288,830,900	244,698,580
3Kbp	Mate-pair	3Kbp	Nextseq500	76	284,204,762	21,599,561,912	195,095,164
6Kbp	Mate-pair	6Kbp	Nextseq500	76	246,238,370	18,714,116,120	152,496,372
9Kbp	Mate-pair	9Kbp	Nextseq500	76	239,919,538	18,233,884,888	148,612,724
Total					1,310,484,838	140,106,010,288	1,031,973,202

1 Table 2. Summary statistics of generated Illumina synthetic long read (Moleculo) library.
2

	500-1499bp	>= 1500bp
3		
4		
5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		
21		
22		
23		
24		
25		
26		
27		
28		
29		
30		
31		
32		
33		
34		
35		
36		
37		
38		
39		
40		
41		
42		
43		
44		
45		
46		
47		
48		
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		
61		
62		
63		
64		
65		

1 Table 3. Summary statistics of generated long reads data using Pacbio RS II system.
2

3		
4	Number of Reads	1,005,571
5		
6		
7		
8	Total Bases	5,836,969,225
9		
10		
11		
12	Length of longest (shortest) read	50,132(50)
13		
14		
15	Average read length	5,804.63
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

Table 4. Summary statistics of generated transcriptome data obtained from six organ tissues using Illumina platform.

Tissue	Sample Name	Read Length	Read Count	Total Base (bp)
Hemocyte	Hemocyte_1	76	20,815,674	1,581,991,224
	Hemocyte_2	76	26,704,666	2,029,554,616
	Hemocyte_2	76	53,068,562	4,033,210,712
Malpighian Tube	Malpighi_1	76	22,635,428	1,720,292,528
	Malpighi_2	76	24,893,788	1,891,927,888
	Malpighi_3	76	45,213,164	3,436,200,464
Midgut	Midgut_1	76	23,350,138	1,774,610,488
	Midgut_2	76	24,597,972	1,869,445,872
	Midgut_3	76	50,949,986	3,872,198,936
Head	Head_1	76	26,526,276	2,015,996,976
	Head_2	76	26,581,124	2,020,165,424
	Head_3	76	40,900,456	3,108,434,656
Integument	Skin_1	76	24,592,846	1,869,056,296
	Skin_2	76	42,775,430	3,250,932,680
	Skin_3	76	35,043,570	2,663,311,320
Fat Body	Fat Body_1	76	24,637,810	1,872,473,560
	Fat Body_2	76	24,037,494	1,826,849,544
	Fat Body_3	76	40,817,582	3,102,136,232
Anterior-Middle/Silk Gland	AM/Silk Gland_1	76	21,399,638	1,626,372,488
	AM/Silk Gland_2	76	24,292,386	1,846,221,336
	AM/Silk Gland_3	76	37,331,530	2,837,196,280
Posterior/Silk Gland	P/Silk Gland_1	76	27,359,580	2,079,328,080
	P/Silk Gland_2	76	23,300,962	1,770,873,112
	P/Silk Gland_3	76	39,421,430	2,996,028,680
Testis	Testis_1	76	40,890,404	3,107,670,704
	Testis_2	76	45,733,846	3,475,772,296
	Testis_3	76	44,985,224	3,418,877,024
Ovary	Ovary_1	76	40,797,628	3,100,619,728
	Ovary_2	76	40,409,752	3,071,141,152
	Ovary_3	76	42,417,892	3,223,759,792

1 Table 5. Summary statistics of the *A. yamamai* genome (>2kb).
2
3

4 **Assembled Genome**
5

6	Size(1n)	656 Mb
7		
8	GC level	34.07
9		
10	No. scaffolds	3,675
11		
12	N50 of scaffolds (bp)	739,388
13		
14	N bases in scaffolds (%)	19,257,439 (2.93)
15		
16	Longest(shortest) scaffolds (bp)	3,156,949 (2,003)
17		
18	Average scaffold Length (bp)	178,657.53
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
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Repeat Element	No. Element	Length (%)
SINE	59,968	8,615,338(1.30)
LINE	426,522	101,251,176(15.31)
LTR element	53,977	4,552,386(0.69)
DNA element	512,760	69,071,227(10.44)
Small RNA	43,645	6,691,619(1.01)
Simple repeat	135,989	6,256,839(0.95)
Low complexity	19,937	932,829(0.14)
Unclassified	294,190	54,552,009(8.25)

Table 6. Summary of identified repeat elements in the *A. yamamai* genome.

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

Table 7. Summary statistics of ab initio, RNA-seq based and homology-based gene prediction results.

Evidence Type	Programs	Element	Total count	Exon/Gene	Total length(bp)	Mean length(bp)	
<i>ab_initio</i>	Augustus	Gene	14,576	4.85	142,415,318	9,770.53	
		Exon	70,733		14,736,668	208.34	
	Geneid	Gene	10,946	2.25	46,119,402	4,213.35	
		Exon	24,686		3,925,563	159.01	
	GeneMarks-ET	Gene	27,754	5.50	273,745,951	9,863.29	
		Exon	152,660		30,847,503	202.06	
	RNA-seq	Cufflinks Transdecoder	Gene	36,213	7.03	840,429,061	23,207.94
			Exon	254,770		201,721,675	791.77
Known Gene (NCBI lepidoptera)	PASA (gmap)		44,561		22,484,151	504.57	

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

Table 8. Summary statistics for the consensus gene set of the *A. yamamai* genome.

Element	No. elements	Exon/Gene	Avg. length	Total length	Genome coverage(%)
Gene	15,481		11,016.34	170,543,958	25.78
		5.64			
Exon	87,346		1,346.23	20,840,925	3.31

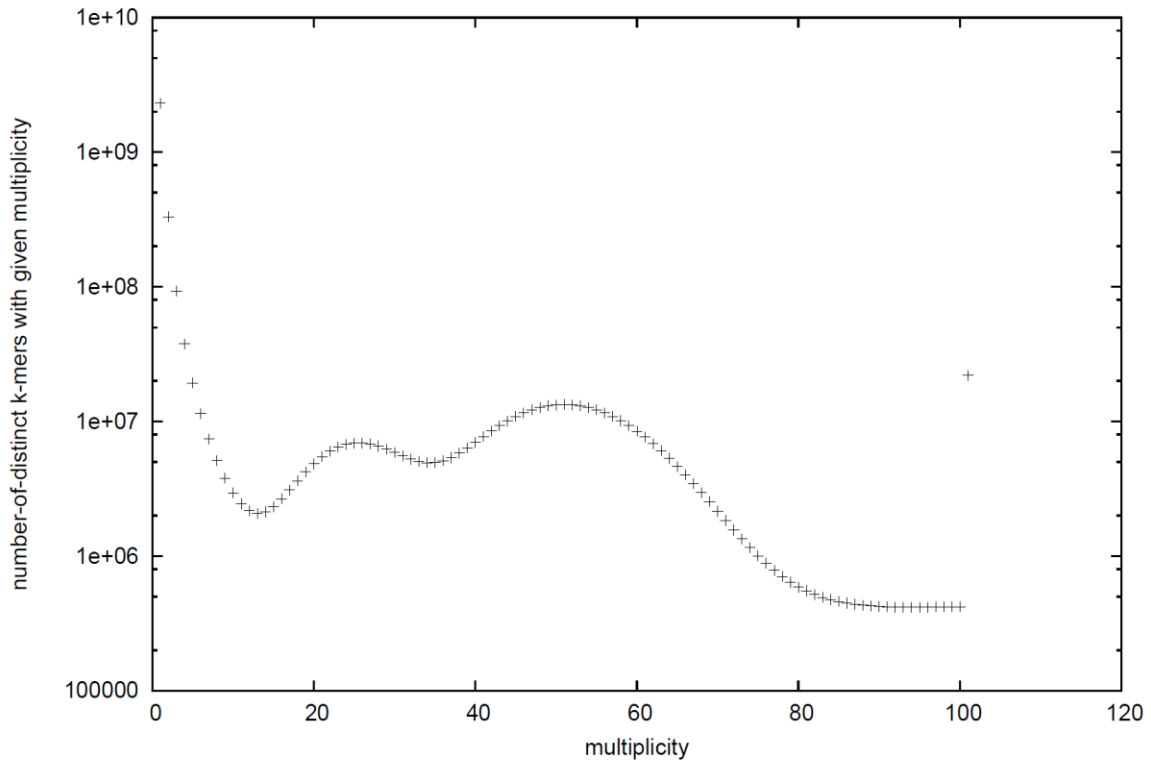
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figures

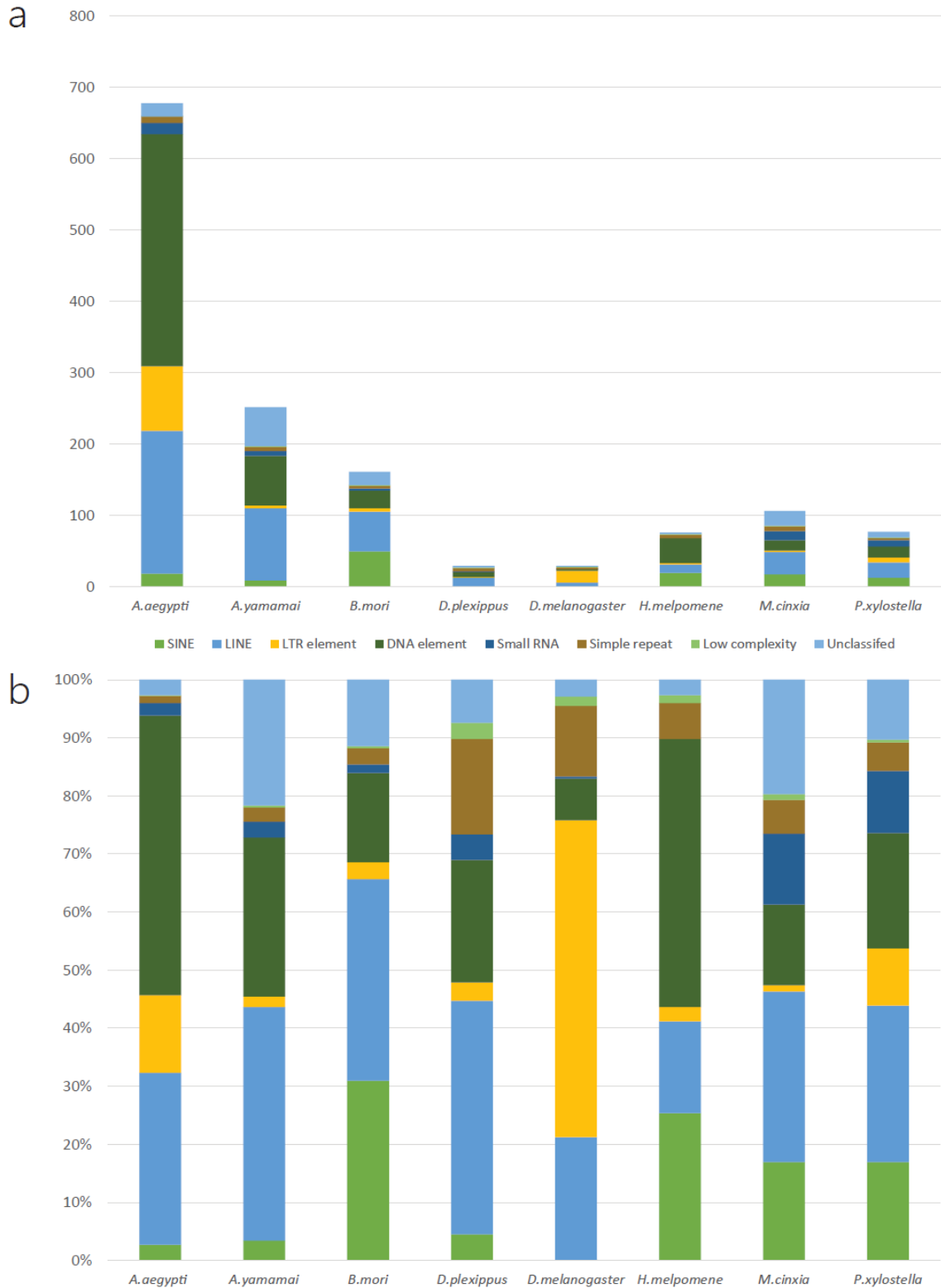
Figure 1. Photograph of *Antheraea Yamamai*. From left- larva, cocoon and adult *A. yamamai*, respectively. Green color is one of the representative characteristics of tensan silk.



1 Figure 2. 19-mer distribution of *A. yamamai* genome using jellyfish with 350bp paired-end
2
3 whole genome sequencing data.
4
5
6
7



1 Figure 3. Amount and proportion of identified repeat element from 8 species including *A.*
 2 *yamamai*. a. Absolute amount of repeat element classified into 8 different categories. b.
 3
 4 Proportion of each repeat element in identified total repeat element.
 5
 6
 7
 8
 9



1 Figure 4. Constructed phylogenetic tree and comparative gene family analysis. Nodes value
 2 indicate Bayesian posterior probability, bootstrap and gene expansion, contraction value.
 3
 4 Orange and blue color indicate expansion and contraction, respectively. Bar chart indicate the
 5 number of genes cauterized into 4 groups (Specific, 1:Multi, Multi:Multi and 1:1) using
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17
 18
 19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

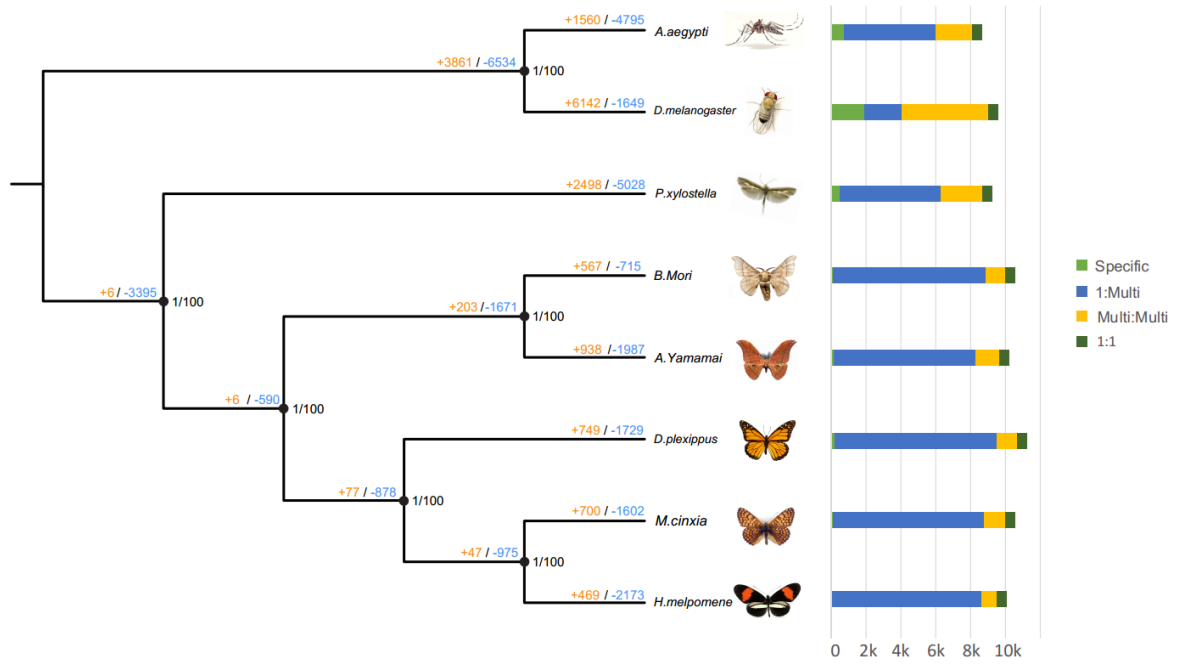
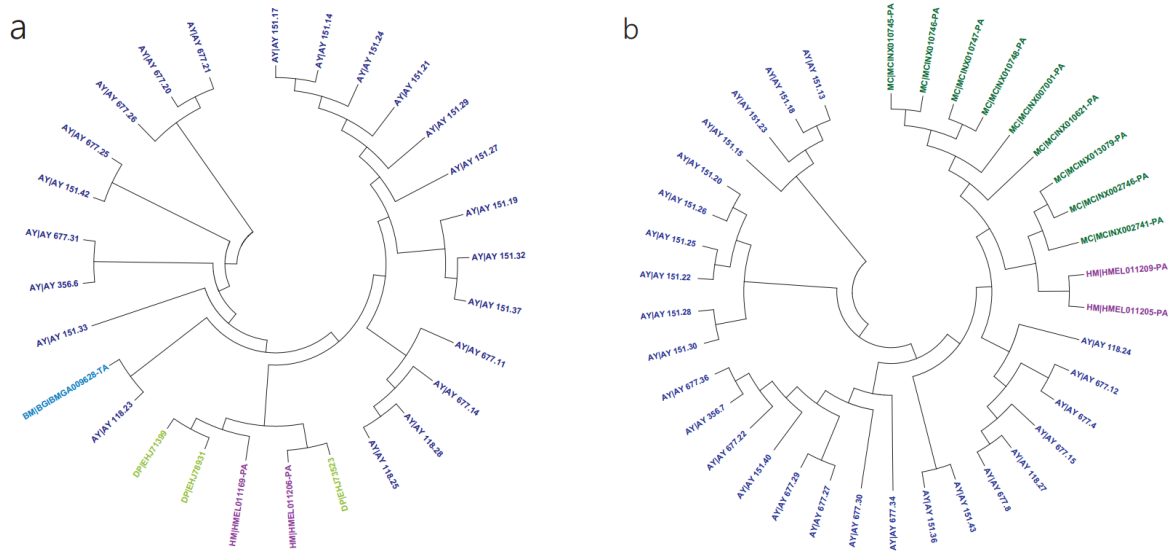
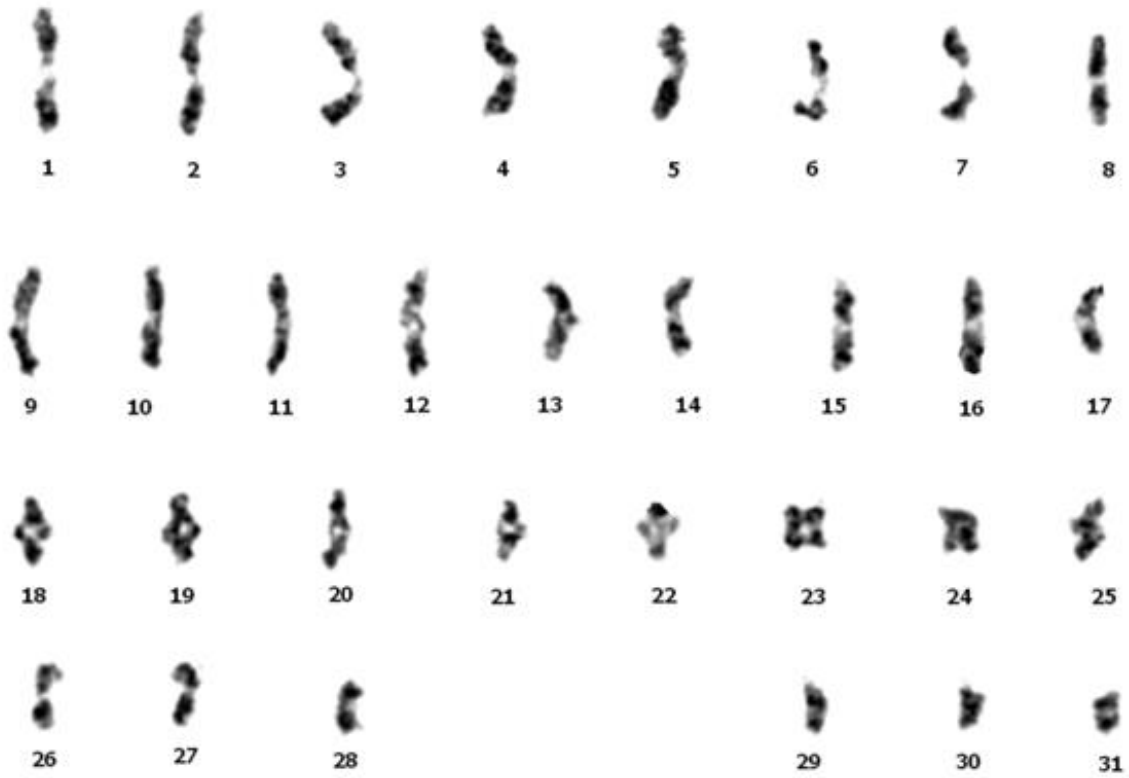


Figure 5. Expansion of chorion gene in *A.yamamai* genome. a and b shows the gene tree of chorion A and B in the rapid expanded gene family cluster, respectively. Color of terminal node indicates each taxon identified in the gene family cluster.



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

Figure 6. Karyotype of *A.yamamai* using a gamete of testis in metaphase.



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



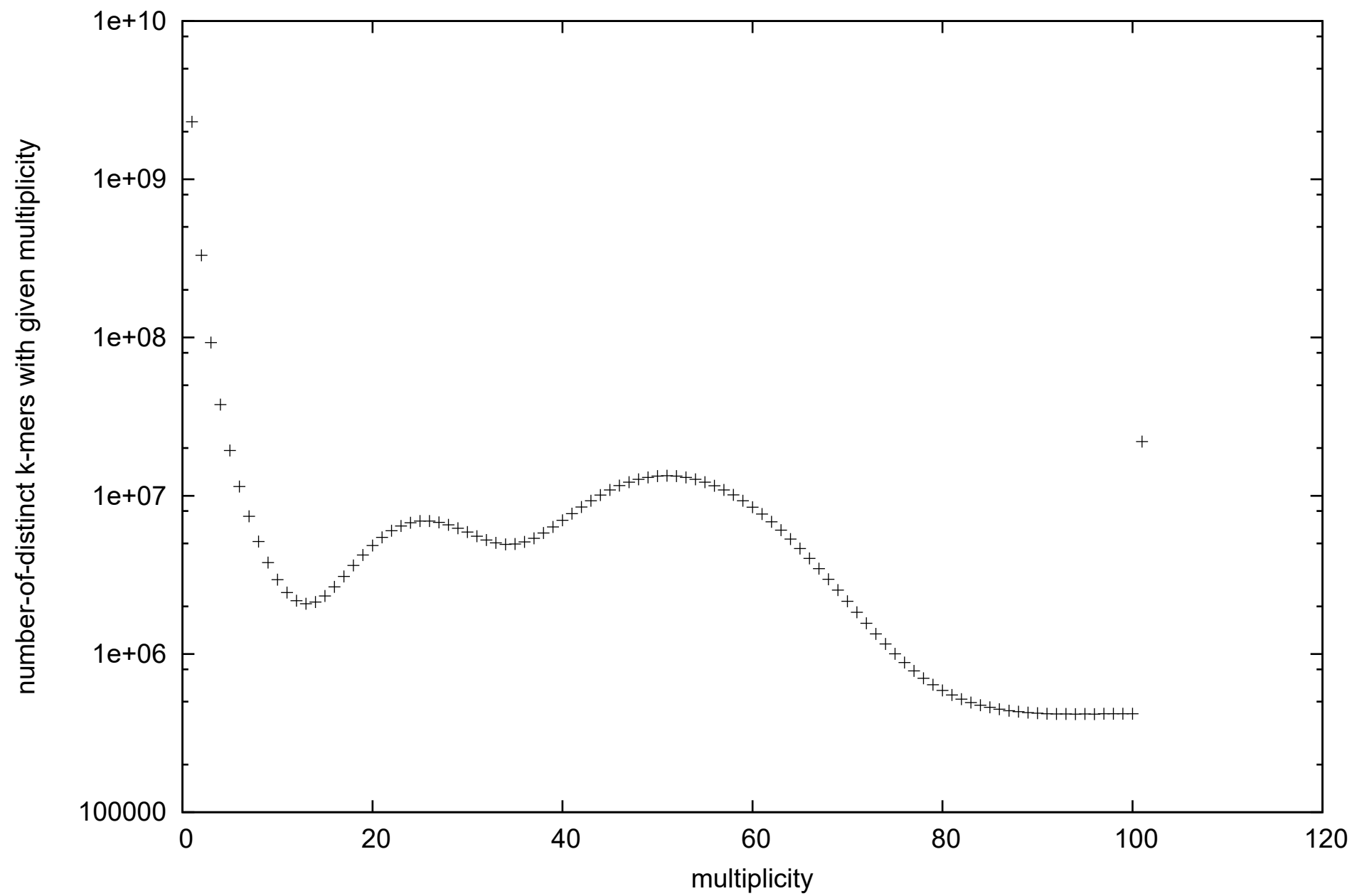
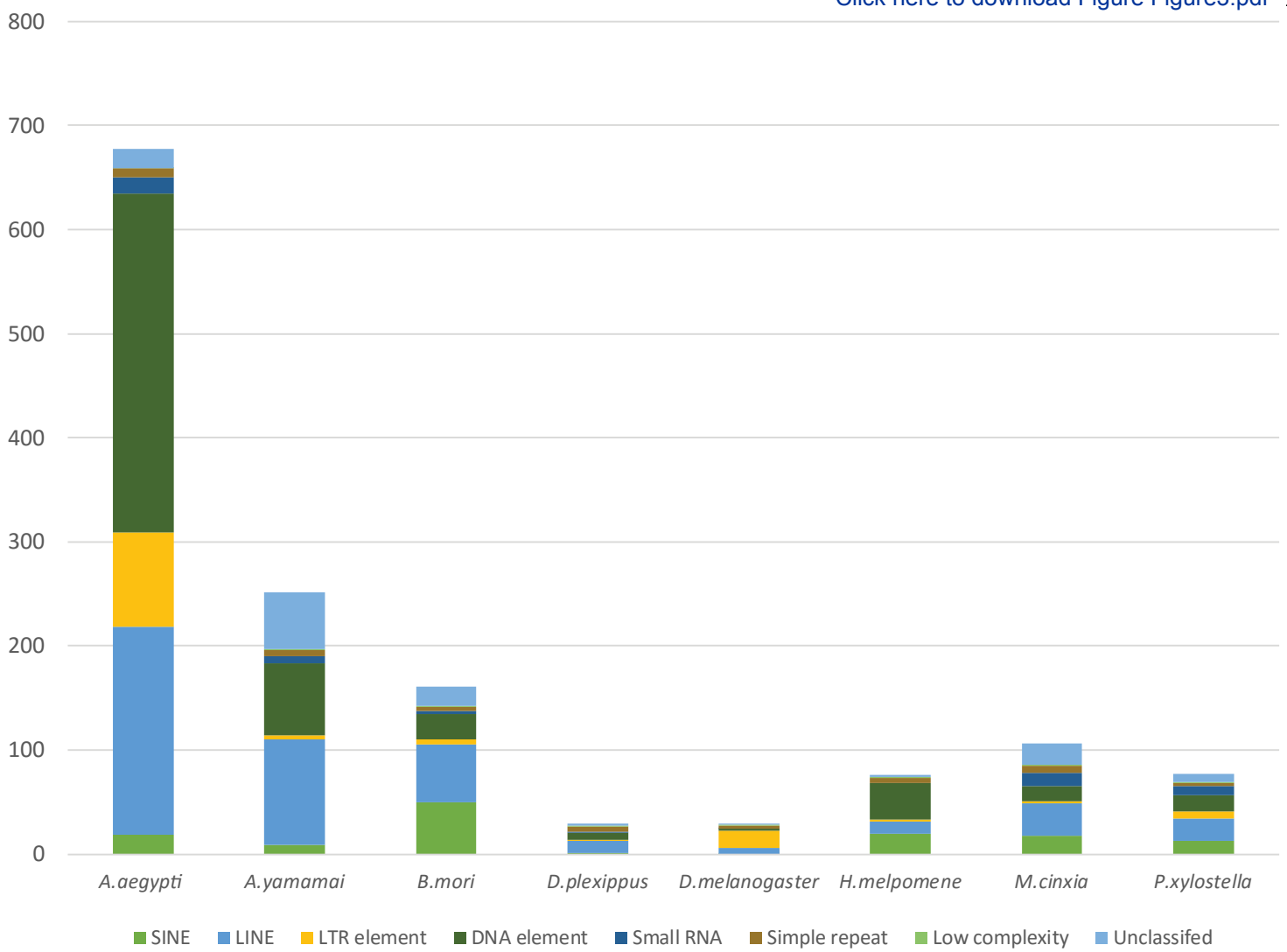
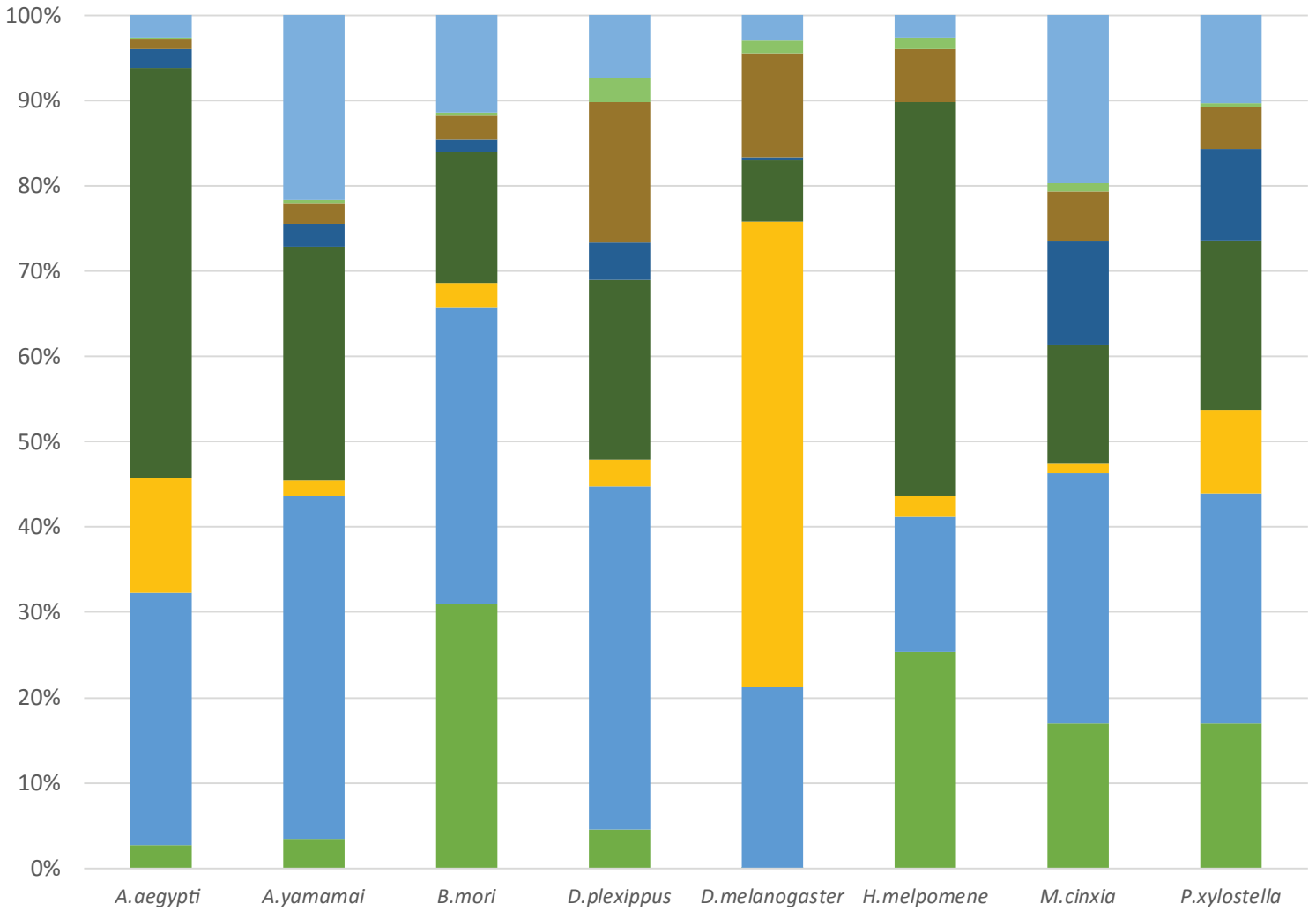
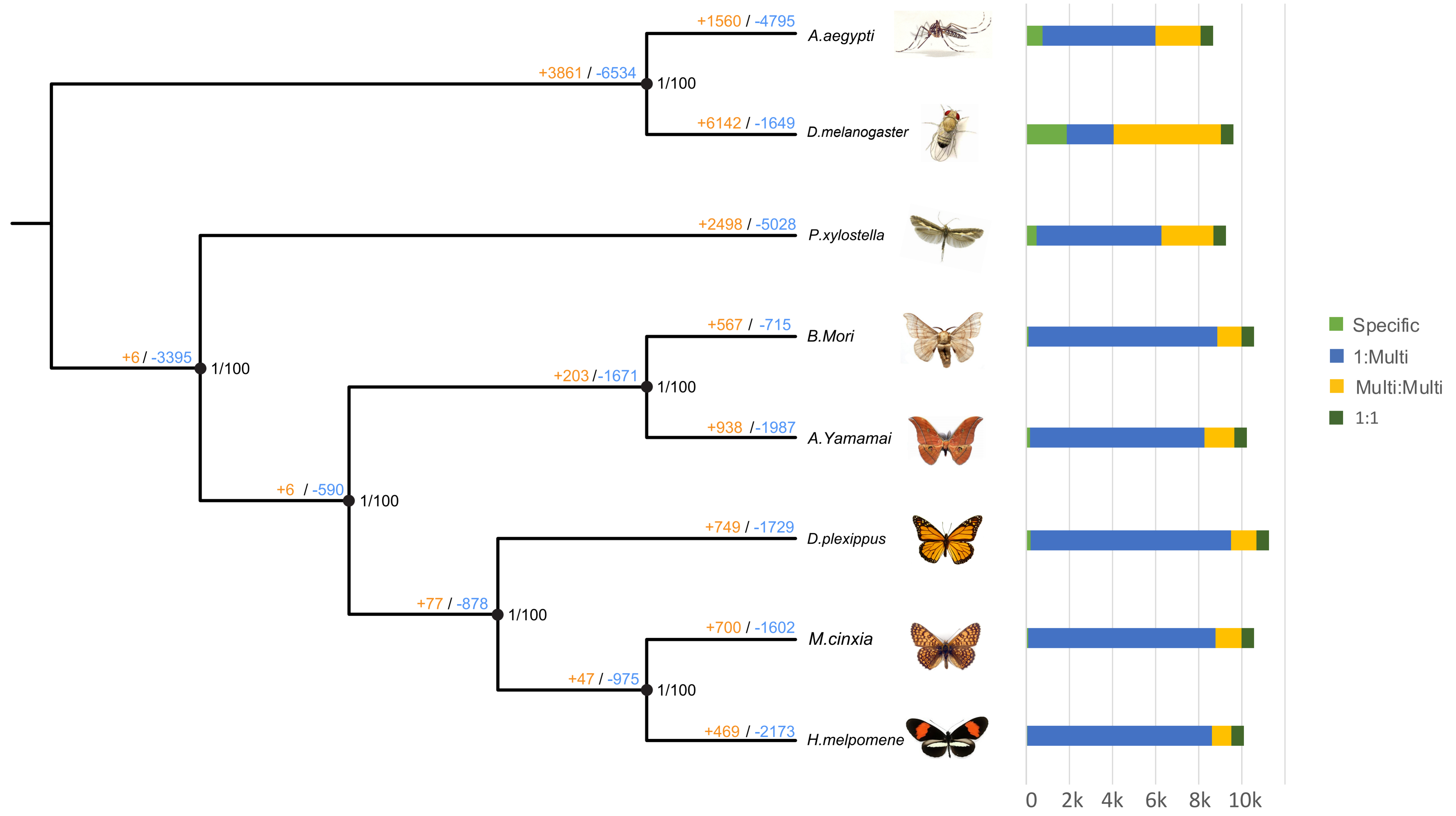


Figure 3
a

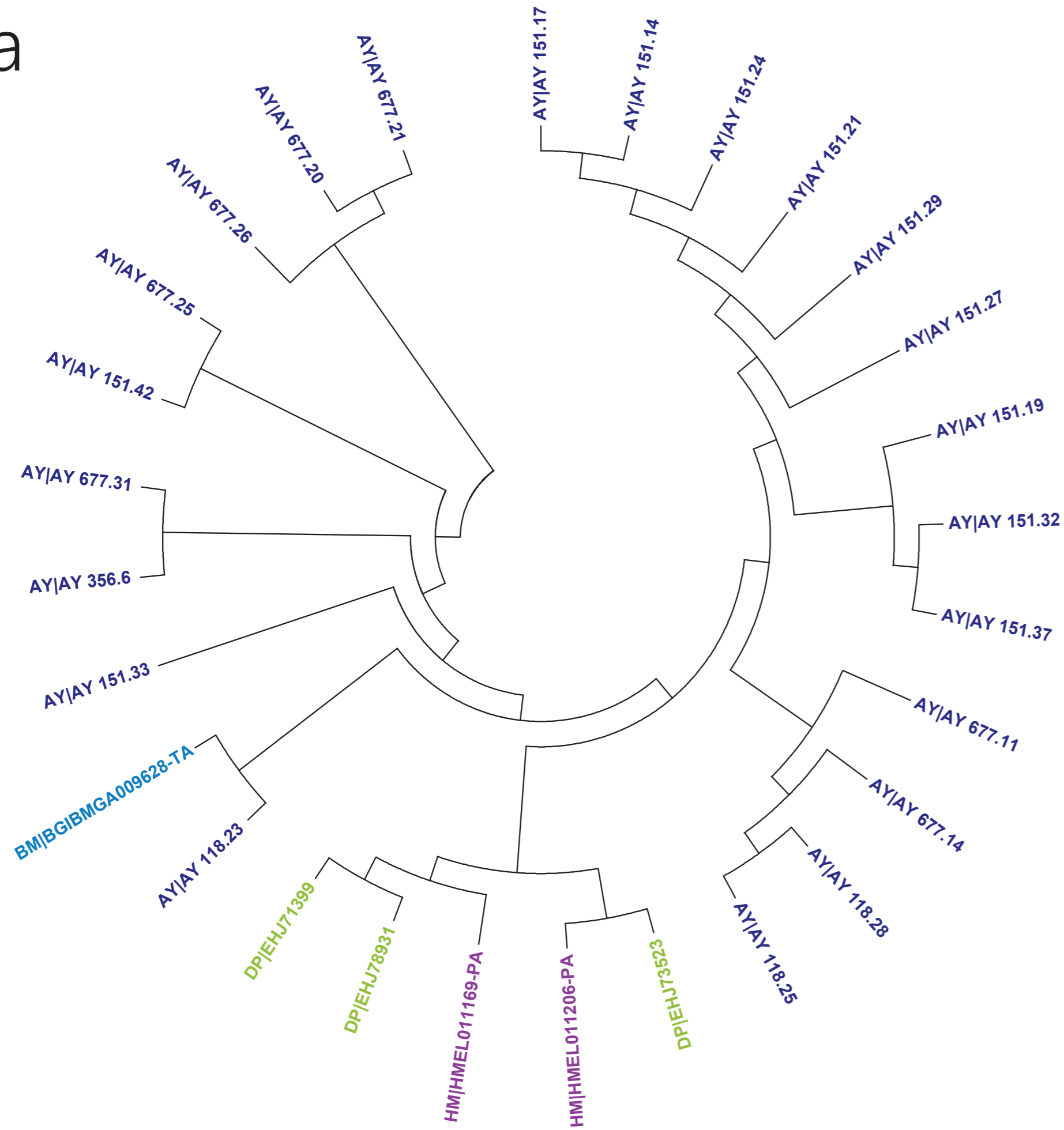


b

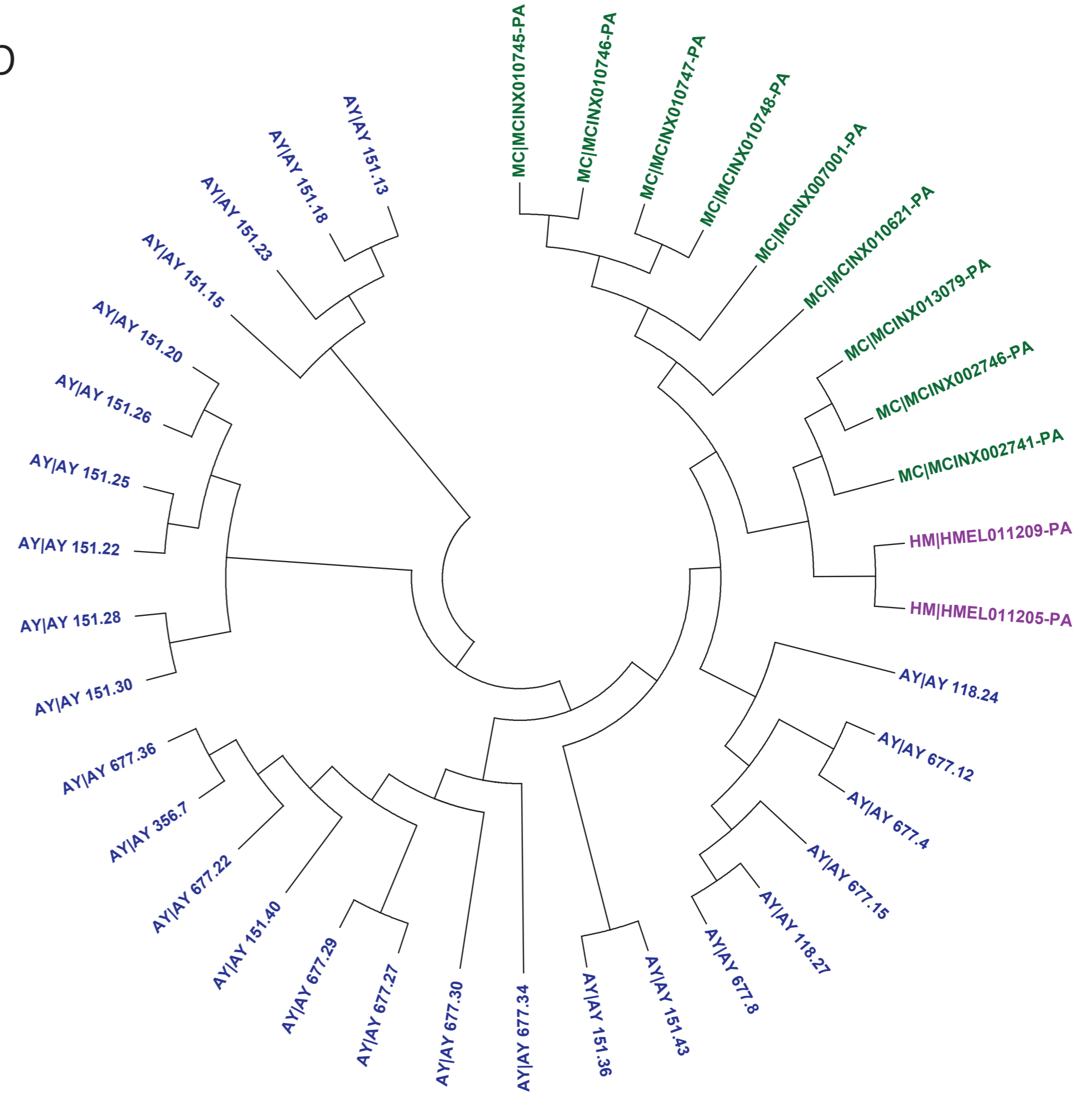


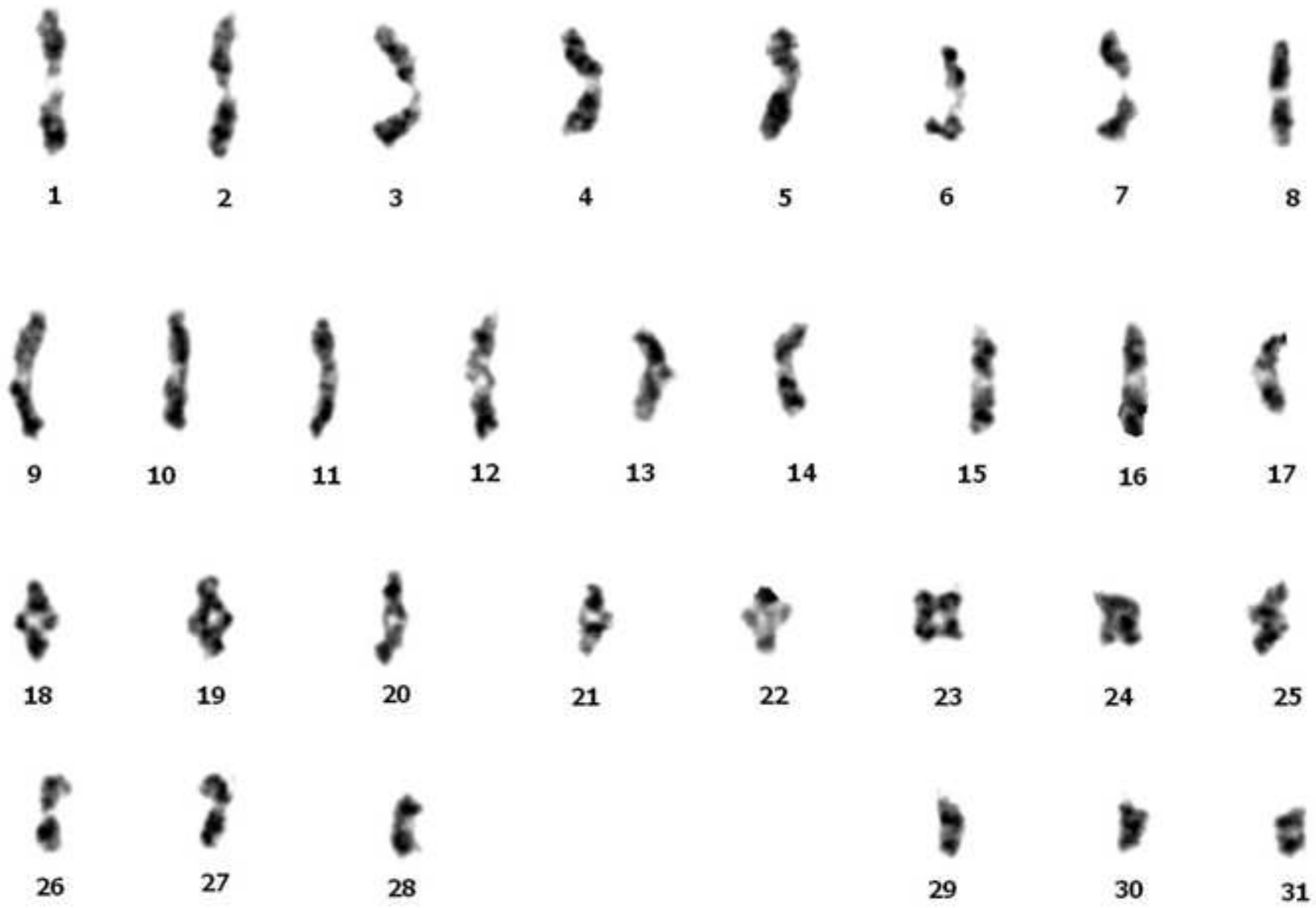


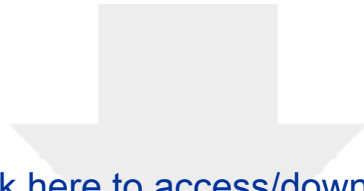
a



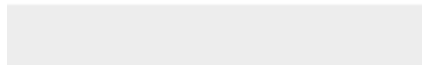
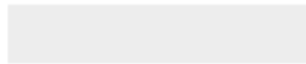
b







Click here to access/download
Supplementary Material
Supplementary_Information.docx



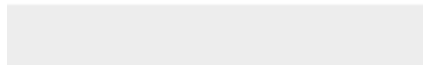


Click here to access/download
Supplementary Material
Supplementary_information2.xlsx



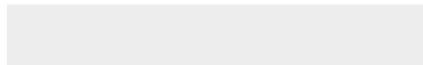
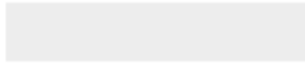


Click here to access/download
Supplementary Material
Response_to_Reviewer 1.docx





Click here to access/download
Supplementary Material
Response_to_Reviewer 2.docx



Oct 17, 2017

Dear Editor of *Gigascience*,

We thank for sending our manuscript out for 2nd review and for obtaining constructive suggestions from two expert referees. Especially, we want to appreciate the editor who give this great opportunity for publication in *Gigascience* of great reputation.

Remained minor concerns were fully addressed and genome data was also sent to the administrator of Lepbase followed the suggestion of reviewer 2.

We hope that our revised manuscript and supporting analysis can convince the reviewers and meet the high-quality standard of *Gigascience*.

Looking forward to hearing from you again.

Thank you.

With best regards,

Prof. Seung-Won Park

Department of Biotechnology,

Catholic University of Daegu, Gyeongsan-si, Gyeongsangbuk-do 38430, Republic of Korea,

Tel: +82-53-850-3176, E-mail: microsw@cu.ac.kr