

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

<b>Manuscript Number:</b>	GIGA-D-17-00085R2	
<b>Full Title:</b>	Genome sequence of Japanese oak silk moth, <i>Antheraea yamamai</i> : the first draft genome in family Saturniidae	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	Rural Development Administration (PJ010442)	Dr Seong-Ryul Kim
<b>Abstract:</b>	<p><b>Background</b>  <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><b>Findings</b>            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(&gt;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 21,124 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).</p> <p><b>Conclusions</b>            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>	
<b>Corresponding Author:</b>	Seung-Won Park  KOREA, REPUBLIC OF	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Seong-Ryul Kim	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	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
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	We appreciate suggestions of all the reviewers on our manuscript; all the comments made by the reviewers were quite valid. We have responded to all comments point-by-point in the separate rebuttal. We hope our revised manuscript can meet the quality standard of reviewers.
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . 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?	
<b>Resources</b>	Yes
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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.	
Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> ?	
<b>Availability of data and materials</b>	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](#)?

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<sup>1†</sup>, Woori Kwak<sup>2†</sup>, Hyaekang Kim<sup>3</sup>, Kelsey Caetano-Anolles<sup>3</sup>, Kee-Young**  
9 **Kim<sup>1</sup>, Su-Bae Kim<sup>1</sup>, Kwang-Ho Choi<sup>1</sup>, Seong-Wan Kim<sup>1</sup>, Jae-Sam Hwang<sup>1</sup>, Min-Jee Kim<sup>4</sup>,**  
10 **Iksoo Kim<sup>4</sup>, Tae-Won Goo<sup>5</sup> and Seung-Won Park<sup>6\*</sup>**  
11  
12  
13  
14  
15

16 <sup>1</sup>Department of Agricultural Biology, National Academy of Agricultural Science, Rural  
17 Development Administration, Wanju-gun 55365, Republic of Korea; <sup>2</sup>C&K Genomics, Main  
18 Bldg. #420, SNU Research Park, Seoul 151-919, Republic of Korea; <sup>3</sup>Department of  
19 Agricultural Biotechnology and Research Institute of Agriculture and Life Sciences, Seoul  
20 National University, Seoul 151-921, Republic of Korea; <sup>4</sup>College of Agriculture & Life  
21 Sciences, Chonnam National University, Gwangju, Republic of Korea; <sup>5</sup>Department of  
22 Biochemistry, Dongguk University College of Medicine, Gyeongju-si, Gyeongsangbuk-do  
23 38066, Republic of Korea; <sup>6</sup>Department 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](mailto:ksr319@korea.kr); Woori Kwak : [asleo@cnkgenomics.com](mailto:asleo@cnkgenomics.com); Hyaekang Kim :  
43 [hkim458@snu.ac.kr](mailto:hkim458@snu.ac.kr); Kelsey Caetano-Anolles : [kelseyca@gmail.com](mailto:kelseyca@gmail.com); Kee-Young Kim :  
44 [applekky@korea.kr](mailto:applekky@korea.kr); Su-Bae Kim : [subae@korea.kr](mailto:subae@korea.kr); Kwang-Ho Choi : [ckh@korea.kr](mailto:ckh@korea.kr); Seong-  
45 Wan; Seong-Wan Kim : [tarupa@korea.kr](mailto:tarupa@korea.kr); Jae-Sam Hwang : [hwangjs@korea.kr](mailto:hwangjs@korea.kr); Min-Jae Kim :  
46 [minjeekim3@gmail.com](mailto:minjeekim3@gmail.com); Iksoo Kim : [ikkim81@chonnam.ac.kr](mailto:ikkim81@chonnam.ac.kr); Tae-Won Goo :  
47 [gootw@dongguk.ac.kr](mailto: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](mailto:microsw@cu.ac.kr)  
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

# 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 the 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 wild silk moth continues to gradually increase. However, no genomic information for 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 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 21,124 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).

## 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 understanding of the molecular mechanisms related to not only specific phenotypes such as wild silk itself but also the genomic evolution of Saturniidae.

1 24 **Keywords**  
2  
3  
4 25 *Antheraea yamamai*, Japanese silk moth, Japanese oak silk moth, wild silkworm  
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

## Data description

*Antheraea yamamai* (Figure 1), also known as the Japanese oak silk moth, is a wild silk moth species belonging to the Saturniidae family. 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* was 84 MYA(million years ago) and it was similar to 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*, such as characteristics such as thickness, bulkiness, compressive elasticity, and resistance to dyeing chemicals[6-8]. These characteristics receive 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 importance of wild silk moth in research and economic fields, no whole genomic information is currently available for this wild silk moth or any other species from 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 family Saturniidae.



## 50 Sequencing

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

1 74

## 2 3 4 75 **Genome assembly and evaluation**

5  
6  
7  
8 76 Before conducting genome assembly, we conducted k-mer distribution analysis using a 350bp  
9  
10 77 paired-end library in order to estimate the size and characteristics of the *A. yamamai* genome.  
11  
12 78 The quality of our generated raw data was checked using FASTQC[16]( FastQC ,  
13  
14  
15 79 RRID:SCR\_014583). Sequencing artifacts such as adapter sequences and low-quality bases  
16  
17  
18 80 were removed using Trimmomatic[17]. Jellyfish[18] was used to count the k-mer frequency  
19  
20 81 for estimation of the genome size of *A. yamamai*. Figure 2 shows the 19-mer distribution of *A.*  
21  
22 82 *yamamai* genome using a 350bp paired-end library. In the 19-mer distribution, the second peak  
23  
24  
25 83 at approximately half the coverage value (x-axis) of the main peak indicates heterozygosity.  
26  
27 84 Although the inbred line used in this study was the single pair sib-mating maintained for more  
28  
29  
30 85 than 10 generations, high heterozygosity still remains. This phenomenon has been observed in  
31  
32 86 a previous genomic study of the Diamondback moth (*Plutella xylostella*), and sustained  
33  
34  
35 87 heterozygosity as an important genomic characteristic was hypothesized to be a result of  
36  
37 88 environmental adaption[19]. Based on the result of 19-mer distribution analysis, the genome  
38  
39  
40 89 size of *A. yamamai* was estimated to be 709Mb. However, this size might be larger than the  
41  
42 90 real genome size of *A.yamamai* because high heterozygosity could affect the estimation of  
43  
44  
45 91 genome size based on the K-mer distribution. Next, we conducted error correction on Illumina  
46  
47 92 paired-end libraries using the error correction module of Allpaths-LG[20] before the initial  
48  
49  
50 93 contig assembly process (ALLPATHS-LG , RRID:SCR\_010742). After error correction, initial  
51  
52 94 contig assembly with 350bp and 700bp libraries was conducted using SOAP denovo2[21] with  
53  
54  
55 95 the parameter option set at K=19; this approach showed the best assembly statistics compared  
56  
57 96 to other assemblers and parameters (SOAPdenovo2 , RRID:SCR\_014986). Quality control  
58  
59 97 processing for mate-pair libraries and scaffolding was conducted using Nxtrim[22] and  
60  
61  
62  
63  
64  
65

1 98 SSPACE (SSPACE , RRID:SCR\_011848)[23], respectively. At each scaffolding step, SOAP  
2  
3 99 Gapcloser[21] with -l 155 and -p 31 parameters was repeatedly used to close the gaps within  
4  
5 100 each scaffold. In order to obtain a higher quality genome assembly of *A. yamamai*, we  
6  
7 101 employed several long read scaffolding strategies using SSPACE-LongRead[24]. First, we  
8  
9 102 used an Illumina synthetic long read sequencing platform called Moleclo which has been  
10  
11 103 proven valuable for the study of highly heterozygous genomes in previous studies[25, 26].  
12  
13 104 After scaffolding was performed using SSPACE-LongRead with Illumina synthetic long read  
14  
15 105 data, the total number of assembled scaffolds was effectively reduced from 398,446 to 24,558.  
16  
17 106 The average scaffold length was also extended from 1.7 Kb to 24.8 Kb. However, there was  
18  
19 107 no impressive improvement in N50 length (approximately 91 Kb to 112 Kb) of assembled  
20  
21 108 scaffolds. Therefore, we employed another type of long read data generated from 10 cells of  
22  
23 109 Pacbio RS II system with P6-C4 chemistry. After final scaffolding processing using Pacbio  
24  
25 110 long reads, the number of scaffolds was reduced to 3,675 and N50 length was effectively  
26  
27 111 extended from 112 Kb to 739 Kb. Summary statistics of the assembled *A. yamamai* genome is  
28  
29 112 provided in Table 5. Final assembly of the *A. yamamai* genome was 656 Mb(>2kb) long with  
30  
31 113 3,675 scaffolds and the N50 length of assembly was 739 Kb with a 34.07% GC ratio. To  
32  
33 114 evaluate the quality of the assembled genome, we conducted BUSCO (Benchmarking  
34  
35 115 Universal Single-Copy Orthologs) analysis[27] using BUSCO v2.0 with insecta\_odb9  
36  
37 116 including 1,658 BUSCOs from 42 species (BUSCO , RRID:SCR\_015008). From BUSCO  
38  
39 117 analysis, 96.7% of BUSCOs were completely detected in the assembled genome (1,576 :  
40  
41 118 complete and single-copy, 27 : complete and duplicated) among 1,658 tested BUSCOs. The  
42  
43 119 number of fragmented and missing BUSCOs was 21 and 34, respectively. Based on the result  
44  
45 120 of BUSCO analysis, the genome of *A.yamamai* presented here was considered properly  
46  
47 121 constructed for downstream analysis.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Repeat identification and comparative repeat analysis

To identify repeat elements of the *A. yamamai* genome, a custom repeat library was constructed using RepeatModeler with RECON[28], RepeatScout[29] and TRF[30]. The resulting constructed custom repeat library for *A. yamamai* was further curated using CENSOR[31] search and the curated library was employed in RepeatMasker[32] with Replibase[33]. RepeatMasker was conducted with RMBlast and 'no\_is' option for skipping bacterial insertion element check. Table 6 summarizes the proportion of identified mobile elements in the *A. yamamai* genome. The most prevalent repeat elements in the *A. yamamai* genome were LINE element (101 Mb, 15.31% of total genome) and total repeat elements accounted for 37.33% of the total genome. In order to compare the repeat elements of *A. yamamai* with that of other genomes, we conducted the same process for seven public genomes which are close neighbors of *A. yamamai* - *Aedes aegypti*[34], *Bombyx mori*[35], *Danaus plexippus*[36], *Drosophila melanogaster*[37], *Heliconius melpomene*[38], *Melitaea cinxia*[39] and *Plutella xylostella*[19]. Figure 3 displays the amount and proportion of identified repeat elements from the 8 species. Despite the small genome size of *B. mori*, the total amount of identified SINE element in the *B. mori* genome was 5.77 times larger than that of *A. yamamai*. The top 5 expanded repeat elements in *A. yamamai* genome were DNA/RC, LINE/L2, LINE/RTE-BovB, DNA/TcMar-Mariner and LINE/CR1. Among these, DNA/TcMar-Mariner was the specifically expanded repeat element in *A. yamamai* among 8 species. In *B. mori*, SINE/tRNA-CR1, LINE/Jockey, DNA/RC, LINE/CR1-Zenon and LINE/RTE-BovB were the top 5 expanded repeat elements. When comparing the repeat elements of *A. yamamai* with those of *B. mori*, which are both producers of the same type of silk, repeat elements showed family and species-specific patterns in the two silk moth lineages. Particularly, we found that the mariner repeat element, which was found specifically expanded in the *A. yamamai* genome, was also included in the fibroin gene.

1 147 A previous sequencing study also showed that the mariner repeat element was inserted in the  
2  
3 148 5'-end of fibroin gene of *A. yamamai*[40]. Fibroin is the core component of the silk protein  
4  
5  
6 149 found in silk moth, and the physical characteristics of silk mainly depend on the types and  
7  
8 150 unique repeat motif of the fibroin[41]. This gene is known to have hundreds of tandem repeat  
9  
10  
11 151 motifs and these kinds of tandem repeats can be derived through transposable elements. This  
12  
13 152 indicates that the mariner repeat element, specifically expanded in the *A. yamamai* genome,  
14  
15  
16 153 may play an important role in development of the unique silk of *A. yamamai*, and the lineage-  
17  
18 154 specific repeat elements may be one of the candidate evolution forces related to host-specific  
19  
20  
21 155 phenotype during genome evolution.  
22  
23  
24 156

## 27 157 **Gene prediction and annotation**

28  
29  
30  
31 158 Three different algorithms were used for gene prediction of the *A. yamamai* genome: *ab initio*,  
32  
33  
34 159 RNA-seq transcript based, and protein homology-based approaches. For *ab initio* gene  
35  
36 160 prediction, Augustus[42], Geneid[43] and GeneMarks-ET[44] were employed. Augustus was  
37  
38  
39 161 trained using known genes of *A. yamamai* in NCBI database and mapping information of RNA-  
40  
41 162 seq data obtained from Tophat[45]( TopHat , RRID:SCR\_013035) was also utilized for gene  
42  
43  
44 163 prediction. Geneid was used with predefined parameters for *Drosophila melanogaster*.  
45  
46 164 GeneMarks-ET was employed using junction information of genes from transcriptome data  
47  
48  
49 165 alignment. For RNA-seq transcript based prediction, generated transcriptome data from ten  
50  
51 166 organ tissues of *A. yamamai* were aligned to the assembled genome and gene information was  
52  
53  
54 167 predicted using Cufflinks[46](Cufflinks , RRID:SCR\_014597). The longest CDS sequences  
55  
56 168 were identified from Cufflinks results using Transdecoder. For the homology-based approach,  
57  
58  
59 169 all known genes of order Lepidoptera in the NCBI database were aligned using PASA[47].  
60  
61  
62  
63  
64  
65

1 170 Table 7 shows the gene prediction results from each method. Gene prediction results from  
2  
3 171 different prediction algorithms were combined using EVM (Evidence Modeler)[48] and a  
4  
5  
6 172 consensus gene set of the *A. yamamai* genome was created. Manual curation was performed  
7  
8 173 based on the 5 evidences (3 in-silico, known protein and RNA-seq) using IGV[49] and Blastp.  
9  
10  
11 174 Using IGV with each gene evidence and comparing results with known genes via blastp, we  
12  
13 175 mainly focused on the removing false positively predicted genes which don't have enough  
14  
15  
16 176 evidences. And merged and spliced gene structured were corrected by comparing the gene  
17  
18 177 structure with known exon structure in NCBI NR database. In addition, fibroin and sericin  
19  
20  
21 178 genes which couldn't be properly predicted because of its high repeat motif were also manually  
22  
23 179 identified with previously known sequences[40, 50] with RNA-seq data. The final gene set of  
24  
25 180 *A. yamamai* genome contains 15,481 genes. Summary statistics for the consensus gene set is  
26  
27  
28 181 provided in Table 8. The average gene length was 11,016.34 bp with a 34.38% GC ratio and  
29  
30  
31 182 the number of exons per gene was 5.64. In order to identify the function of predicted genes in  
32  
33 183 *A. yamamai*, three non-redundant sequence databases (Swiss-Prot[51], Uniref100[51], and  
34  
35 184 NCBI NR[52]) as well as the gene information of two species (*B. mori* and *D. melanogaster*)  
36  
37  
38 185 were used for target databases using Blastp. Additionally, protein domain searches were  
39  
40 186 conducted on the consensus gene set using InterproScan5[53]. Figure S1 shows the top 20  
41  
42  
43 187 identified terms from 7 different InterproScan5 analyses. Among the various analysis  
44  
45 188 conducted using InterproScan5, gene ontology analysis with Pfam database showed that a large  
46  
47  
48 189 proportion of genes in the *A.yamamai* genome were related with the function of molecular  
49  
50 190 binding, catalytic activity, internal component of membrane, metabolic process, oxidation-  
51  
52 191 reduction process and transmembrane transport.

## 59 193 **Comparative genome analysis**

1 194 We used OrthoMCL[54] and RBH(Reciprocal Best Hit) within blastp for identification of gene  
2  
3 195 family clusters and 1:1 orthologous gene sets. Gene information of 7 taxa (*A. aegypti*, *B. mori*,  
4  
5 196 *D. plexippus*, *D. melanogaster*, *H. melpomene*, *M. cinxia* and *P. xylostella*), same taxa used in  
6  
7  
8 197 repeat analysis, was employed for OrthoMCL with *A. yamamai*. A total of 17,406 gene family  
9  
10 198 clusters were constructed and 3,586 1:1 orthologous genes were identified. Before conducting  
11  
12 199 comparative genome analysis, we constructed phylogenetic trees for the 8 species. In order to  
13  
14 200 build the phylogenetic tree, multiple sequence alignment for the 1:1 orthologous genes of all 8  
15  
16 201 species was conducted using PRANK[55], and Gblocks[56] was used to obtain conserved  
17  
18 202 blocks for the phylogenetic tree. Conserved block sequences were sequentially concatenated  
19  
20 203 to obtain one consensus sequence for each species. MEGA[57] was used for constructing  
21  
22 204 Neighbor-Joining Trees (bootstrap 1000, maximum composite likelihood, transitions +  
23  
24 205 transversions, and gamma distributed option) and MrBayes[58] was employed for the  
25  
26 206 construction of Bayesian inference trees. To select the best evolution model for our data,  
27  
28 207 Modeltest[59] was conducted and the GTR based invariant model was chosen based on the  
29  
30 208 AIC value of Modeltest. Figure 4 shows the constructed phylogenetic tree of the 8 species using  
31  
32 209 3,586 orthologous genes. The bootstrap value and Bayesian poster probability value of all  
33  
34 210 nodes were 100 and 1, respectively. The closest neighbor of *A. yamamai* was *B. mori*, which is  
35  
36 211 included in Bombycidae family; this result is consistent with that of previous studies. Three  
37  
38 212 butterfly species (*D. plexippus*, *M. cinxia* and *H. melpomene*) included in Nymphalidae family  
39  
40 213 were also shown to share a common ancestor with families Saturniidae and Bombycidae.  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50 214 Based on the constructed phylogenetic tree, gene family expansion and contraction analysis  
51  
52 215 was conducted using a 2 parameter model in CAFE[60] and the gene tree was constructed using  
53  
54 216 protein sequence via MEGA[57]. Figure 4 shows the result of gene family expansion and  
55  
56 217 contraction analysis of 8 species. 938 and 1,987 gene families of *A. yamamai* and 567 and 715  
57  
58 218 gene families of *B. mori* were estimated to be expanded and contracted from the common  
59  
60  
61  
62  
63  
64  
65

1 219 ancestors, respectively. Among these, 15 gene families in *A. yamamai* were estimated to be  
2  
3  
4 220 under rapid expansion during the evolution process. Functions of genes in rapidly expanded  
5  
6 221 gene families of *A. yamamai* were transposase, fatty acid synthase, zinc finger protein, chorion  
7  
8 222 (eggshell protein), reverse transcriptase, prostaglandin dehydrogenase, RNA-directed DNA  
9  
10 223 polymerase, gag like protein, juvenile hormone acid methyltransferase, facilitated trehalose  
11  
12  
13 224 transporter and glucose dehydrogenase. Figure 5 shows the gene tree of two chorion gene  
14  
15 225 (chorion class A and B) family clusters rapidly expanded in the *A. yamamai* genome. Chorion,  
16  
17  
18 226 called eggshell protein, composes the surface of egg and protects the embryo from  
19  
20  
21 227 environmental threats such as desiccation, flooding, freezing, infection of microorganisms, and  
22  
23 228 physical destruction. It also provides channels, such as aeropyle, which enables gas exchange  
24  
25 229 and maintains proper condition for diapause egg[61]. These diverse functions of eggshell are  
26  
27  
28 230 implemented by the specific eggshell structure and the surface structure of eggshell varies  
29  
30 231 between species for the adaptation in a different environment. The ancestor of *Antheraea* has  
31  
32  
33 232 the unique aeropyle structure called “aerophyle crown” on the eggshell surface[62]. This  
34  
35 233 unique structure is formed by the circular vertical projection of lamellar chorion from follicle  
36  
37  
38 234 cell and it surrounds the aeropyles near the end of oogenesis[63]. Acquiring this kind of *de*  
39  
40 235 *novo* complex structure requires numerous genetic changes and a previous study about  
41  
42 236 *Antheraea Polyphemus* has shown that over a hundred chorion specific polypeptides were  
43  
44  
45 237 involved for this unique ultra-structure[63]. Therefore, the specific rapid expansion of chorion  
46  
47 238 class A and B gene family in *A. yamamai* genome might be one of the convincing molecular  
48  
49  
50 239 explanation for acquiring this unique ultrastructure in the eggshell surface of *Antheraea* genus.  
51  
52 240 However, this unique ultra-structure tends to be reduced during current evolution process of  
53  
54  
55 241 the *Antheraea* genus. Types of eggshell structure in *Antheraea* genus can be categorized into  
56  
57 242 multiple classes based on the morphology and regional distribution of aeropyle[62]. The shape  
58  
59 243 of aeropyle in *A. yamamai* egg is known to be converted to mound shape from the crown shape  
60  
61  
62  
63  
64  
65



1 244 and these aeropyle mounds only exist in the narrow band surrounding the micropyle region[62].  
2  
3 245 Only a very few, small aeropyle crowns remained and it is entirely different with the ancestral  
4  
5  
6 246 form of eggshell surface mostly covered by aeropyle crowns. These regional differences were  
7  
8 247 known to be adjusted by regional difference of filler genes during choriogenesis[64] and the  
9  
10  
11 248 additional regulations of related genes for choriogenesis have to be considered. This indicates  
12  
13 249 that specifically expanded chorion gene families of *A. yamamai* may be one of the remaining  
14  
15  
16 250 evolutionary tracks in the genome of *Antheraea* genus. However, further functional studies  
17  
18 251 must be conducted to resolve the limited understanding about the relationship between these  
19  
20  
21 252 expanded chorion gene families and the current eggshell surface formation of *A. yamamai*.  
22  
23  
24 253 The constructed genome of *A.yamamai* presented here is the first announced genome in family  
25  
26 254 Saturniidae and the karyotyping analysis using gamete in metaphase showed that the genome  
27  
28  
29 255 of *A. yamamai* consists of 31 chromosomes (Figure 6). This constructed genome information  
30  
31 256 provides more insight into the genome evolution and phylogeny of family Saturniidae, which  
32  
33  
34 257 contains the largest number of species in Lepidoptera. For example, although two silk moths,  
35  
36 258 *A. yamamai* and *B. mori*, appear similar, comparative genome analysis showed the significant  
37  
38 259 differences in the genome size, specific expansion of repeat elements and gene families  
39  
40  
41 260 between families Saturniidae and Bombycidae. In case of molecular phylogeny, most previous  
42  
43 261 phylogenetic studies were limited to few genes due to the lack of genomic information on  
44  
45  
46 262 family Saturniidae. We expect our study and resulting constructed genome will resolve some  
47  
48 263 limitations of molecular phylogenetic and ecological research on Saturniidae species.  
49  
50  
51 264 Additionally, constructed genome information will help researchers better understand the  
52  
53 265 molecular background of wild silk and its production. Silk produced by *A. yamamai*, referred  
54  
55 266 to as *tensan* silk, shows unique characteristics which have made it valuable in various fields.  
56  
57  
58 267 However, *A. yamamai* has not been completely domesticated compared to *B. mori*, making  
59  
60 268 mass production of *tensan* silk infeasible. Understanding of the molecular mechanisms behind  
61  
62  
63  
64  
65

1 269 the tensan silk production process is essential for mass production using biotechnology, and  
2  
3  
4 270 this genome sequence with manually curated gene information is a fundamental resource for  
5  
6 271 related research and industrial improvement. Additionally, the transcriptome data of 10  
7  
8 272 different organ tissues with 3 biological replications presented here may be also useful  
9  
10  
11 273 resources for uncovering the molecular mechanisms related to specific phenotypes of  
12  
13 274 *A.yamamai* and family Saturniidae.  
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

## 275 **Availability of supporting data**

276 The generated genome sequence and gene information of *A. yamamai* are available in  
277 GigaDB[65] and generated raw data is available under project accession PRJNA383008 and  
278 PRJNA383025 of the NCBI database.

## 279 **Competing interests**

280 All authors report no competing interests.

## 281 **Abbreviation**

282 RBH – Reciprocal Best Hit

## 283 **Authors contributions**

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

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

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

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

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

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

290 Caetano-Anolles

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

292

1 **293 Acknowledgements**

2  
3  
4 294 This work was supported by a grant from the Rural Development Administration, Republic of  
5  
6 295 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*.

- Cambridge, UK: Babraham Institute, 2011.
17. Bolger, A.M., M. Lohse, and B. Usadel, *Trimmomatic: a flexible trimmer for Illumina sequence data*. *Bioinformatics*, 2014: p. btu170.
  18. Marçais, G. and C. Kingsford, *A fast, lock-free approach for efficient parallel counting of occurrences of k-mers*. *Bioinformatics*, 2011. **27**(6): p. 764-770.
  19. You, M., Z. Yue, W. He, et al., *A heterozygous moth genome provides insights into herbivory and detoxification*. *Nature genetics*, 2013. **45**(2): p. 220-225.
  20. Gnerre, S., I. MacCallum, D. Przybylski, et al., *High-quality draft assemblies of mammalian genomes from massively parallel sequence data*. *Proceedings of the National Academy of Sciences*, 2011. **108**(4): p. 1513-1518.
  21. Luo, R., B. Liu, Y. Xie, et al., *SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler*. *Gigascience*, 2012. **1**(1): p. 18.
  22. O'Connell, J., O. Schulz-Trieglaff, E. Carlson, et al., *NxTrim: optimized trimming of Illumina mate pair reads*. *Bioinformatics*, 2015. **31**(12): p. 2035-2037.
  23. Boetzer, M., C.V. Henkel, H.J. Jansen, et al., *Scaffolding pre-assembled contigs using SSPACE*. *Bioinformatics*, 2011. **27**(4): p. 578-579.
  24. Boetzer, M. and W. Pirovano, *SSPACE-LongRead: scaffolding bacterial draft genomes using long read sequence information*. *BMC bioinformatics*, 2014. **15**(1): p. 211.
  25. Voskoboynik, A., N.F. Neff, D. Sahoo, et al., *The genome sequence of the colonial chordate, Botryllus schlosseri*. *Elife*, 2013. **2**: p. e00569.
  26. McCoy, R.C., R.W. Taylor, T.A. Blauwkamp, et al., *Illumina TruSeq synthetic long-reads empower de novo assembly and resolve complex, highly-repetitive transposable elements*. *PloS one*, 2014. **9**(9): p. e106689.
  27. Simão, F.A., R.M. Waterhouse, P. Ioannidis, et al., *BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs*. *Bioinformatics*, 2015: p. btv351.
  28. Bao, Z. and S.R. Eddy, *Automated de novo identification of repeat sequence families in sequenced genomes*. *Genome Research*, 2002. **12**(8): p. 1269-1276.
  29. Price, A.L., N.C. Jones, and P.A. Pevzner, *De novo identification of repeat families in large genomes*. *Bioinformatics*, 2005. **21**(suppl 1): p. i351-i358.
  30. Benson, G., *Tandem repeats finder: a program to analyze DNA sequences*. *Nucleic acids research*, 1999. **27**(2): p. 573.
  31. Kohany, O., A.J. Gentles, L. Hankus, et al., *Annotation, submission and screening of repetitive elements in Repbase: RepbaseSubmitter and Censor*. *BMC bioinformatics*, 2006. **7**(1): p. 474.
  32. Tarailo-Graovac, M. and N. Chen, *Using RepeatMasker to identify repetitive elements in genomic sequences*. *Current Protocols in Bioinformatics*, 2009: p. 4.10. 1-4.10. 14.
  33. Bao, W., K.K. Kojima, and O. Kohany, *Repbase Update, a database of repetitive elements in eukaryotic genomes*. *Mobile DNA*, 2015. **6**(1): p. 11.
  34. Nene, V., J.R. Wortman, D. Lawson, et al., *Genome sequence of Aedes aegypti, a major*

- 1 *arbovirus vector*. Science, 2007. **316**(5832): p. 1718-1723.
- 2
- 3 35. Xia, Q., Z. Zhou, C. Lu, et al., *A draft sequence for the genome of the domesticated*
- 4 *silkworm (Bombyx mori)*. Science, 2004. **306**(5703): p. 1937-1940.
- 5
- 6 36. Zhan, S., C. Merlin, J.L. Boore, et al., *The monarch butterfly genome yields insights into*
- 7 *long-distance migration*. Cell, 2011. **147**(5): p. 1171-1185.
- 8
- 9 37. Adams, M.D., S.E. Celniker, R.A. Holt, et al., *The genome sequence of Drosophila*
- 10 *melanogaster*. Science, 2000. **287**(5461): p. 2185-2195.
- 11
- 12 38. Consortium, H.G., *Butterfly genome reveals promiscuous exchange of mimicry adaptations*
- 13 *among species*. Nature, 2012. **487**(7405): p. 94-98.
- 14
- 15 39. Ahola, V., R. Lehtonen, P. Somervuo, et al., *The Glanville fritillary genome retains an ancient*
- 16 *karyotype and reveals selective chromosomal fusions in Lepidoptera*. Nature
- 17 *communications*, 2014. **5**.
- 18
- 19 40. Hwang, J.-S., J.-S. Lee, T.-W. Goo, et al., *Cloning of the fibroin gene from the oak silkmoth,*
- 20 *Antheraea yamamai and its complete sequence*. Biotechnology letters, 2001. **23**(16): p.
- 21 1321-1326.
- 22
- 23
- 24 41. Malay, A.D., R. Sato, K. Yazawa, et al., *Relationships between physical properties and*
- 25 *sequence in silkmoth silks*. Scientific reports, 2016. **6**: p. 27573.
- 26
- 27 42. Stanke, M., M. Diekhans, R. Baertsch, et al., *Using native and syntenically mapped cDNA*
- 28 *alignments to improve de novo gene finding*. Bioinformatics, 2008. **24**(5): p. 637-644.
- 29
- 30 43. Blanco, E., G. Parra, and R. Guigó, *Using geneid to identify genes*. Current protocols in
- 31 *bioinformatics*, 2007: p. 4.3. 1-4.3. 28.
- 32
- 33 44. Lomsadze, A., P.D. Burns, and M. Borodovsky, *Integration of mapped RNA-Seq reads into*
- 34 *automatic training of eukaryotic gene finding algorithm*. Nucleic acids research, 2014: p.
- 35 gku557.
- 36
- 37
- 38 45. Trapnell, C., L. Pachter, and S.L. Salzberg, *TopHat: discovering splice junctions with RNA-*
- 39 *Seq*. Bioinformatics, 2009. **25**(9): p. 1105-1111.
- 40
- 41 46. Trapnell, C., A. Roberts, L. Goff, et al., *Differential gene and transcript expression analysis of*
- 42 *RNA-seq experiments with TopHat and Cufflinks*. Nature protocols, 2012. **7**(3): p. 562-578.
- 43
- 44 47. Campbell, M.A., B.J. Haas, J.P. Hamilton, et al., *Comprehensive analysis of alternative*
- 45 *splicing in rice and comparative analyses with Arabidopsis*. BMC genomics, 2006. **7**(1): p.
- 46 327.
- 47
- 48 48. Haas, B.J., S.L. Salzberg, W. Zhu, et al., *Automated eukaryotic gene structure annotation*
- 49 *using EVIDENCEModeler and the Program to Assemble Spliced Alignments*. Genome
- 50 *biology*, 2008. **9**(1): p. R7.
- 51
- 52
- 53 49. Robinson, J.T., H. Thorvaldsdóttir, W. Winckler, et al., *Integrative genomics viewer*. Nature
- 54 *biotechnology*, 2011. **29**(1): p. 24-26.
- 55
- 56 50. Zurovec, M., N. Yonemura, B. Kludkiewicz, et al., *Sericin Composition in the Silk of*
- 57 *Antheraea yamamai*. Biomacromolecules, 2016. **17**(5): p. 1776-1787.
- 58
- 59 51. Consortium, U., *Reorganizing the protein space at the Universal Protein Resource*
- 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
- (UniProt). *Nucleic acids research*, 2011: p. gkr981.
52. Pruitt, K.D., T. Tatusova, and D.R. Maglott, *NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins*. *Nucleic acids research*, 2007. **35**(suppl 1): p. D61-D65.
53. Jones, P., D. Binns, H.-Y. Chang, et al., *InterProScan 5: genome-scale protein function classification*. *Bioinformatics*, 2014. **30**(9): p. 1236-1240.
54. Li, L., C.J. Stoeckert, and D.S. Roos, *OrthoMCL: identification of ortholog groups for eukaryotic genomes*. *Genome research*, 2003. **13**(9): p. 2178-2189.
55. Löytynoja, A. and N. Goldman, *An algorithm for progressive multiple alignment of sequences with insertions*. *Proceedings of the National Academy of Sciences of the United States of America*, 2005. **102**(30): p. 10557.
56. Castresana, J., *Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis*. *Molecular biology and evolution*, 2000. **17**(4): p. 540-552.
57. Kumar, S., G. Stecher, and K. Tamura, *MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets*. *Molecular biology and evolution*, 2016. **33**(7): p. 1870-1874.
58. Ronquist, F. and J.P. Huelsenbeck, *MrBayes 3: Bayesian phylogenetic inference under mixed models*. *Bioinformatics*, 2003. **19**(12): p. 1572-1574.
59. Posada, D., *Using MODELTEST and PAUP\* to select a model of nucleotide substitution*. *Current protocols in bioinformatics*, 2003: p. 6.5. 1-6.5. 14.
60. De Bie, T., N. Cristianini, J.P. Demuth, et al., *CAFE: a computational tool for the study of gene family evolution*. *Bioinformatics*, 2006. **22**(10): p. 1269-1271.
61. Chapman, R.F., *The insects: structure and function*. 1998: Cambridge university press.
62. Regier, J.C., U. Paukstadt, L.H. Paukstadt, et al., *Phylogenetics of eggshell morphogenesis in Antheraea (Lepidoptera: Saturniidae): unique origin and repeated reduction of the aeropyle crown*. *Systematic biology*, 2005. **54**(2): p. 254-267.
63. Regier, J.C., G.D. Mazur, and F.C. Kafatos, *The silkmoth chorion: morphological and biochemical characterization of four surface regions*. *Developmental biology*, 1980. **76**(2): p. 286-304.
64. Hatzopoulos, A.K. and J.C. Regier, *Evolutionary changes in the developmental expression of silkmoth chorion genes and their morphological consequences*. *Proceedings of the National Academy of Sciences*, 1987. **84**(2): p. 479-483.
65. Sneddon, T.P., P. Li, and S.C. Edmunds, *GigaDB: announcing the GigaScience database*. *GigaScience*, 2012. **1**(1): p. 11.



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
<b>Total</b>					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**

---

6 Size(1n)	656 Mb
8 GC level	34.07
10 No. scaffolds	3,675
12 N50 of scaffolds (bp)	739,388
14 N bases in scaffolds (%)	19,257,439 (2.93)
16 Longest(shortest) scaffolds (bp)	3,156,949 (2,003)
18 Average scaffold Length (bp)	178,657.53

---

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

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)

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

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

---



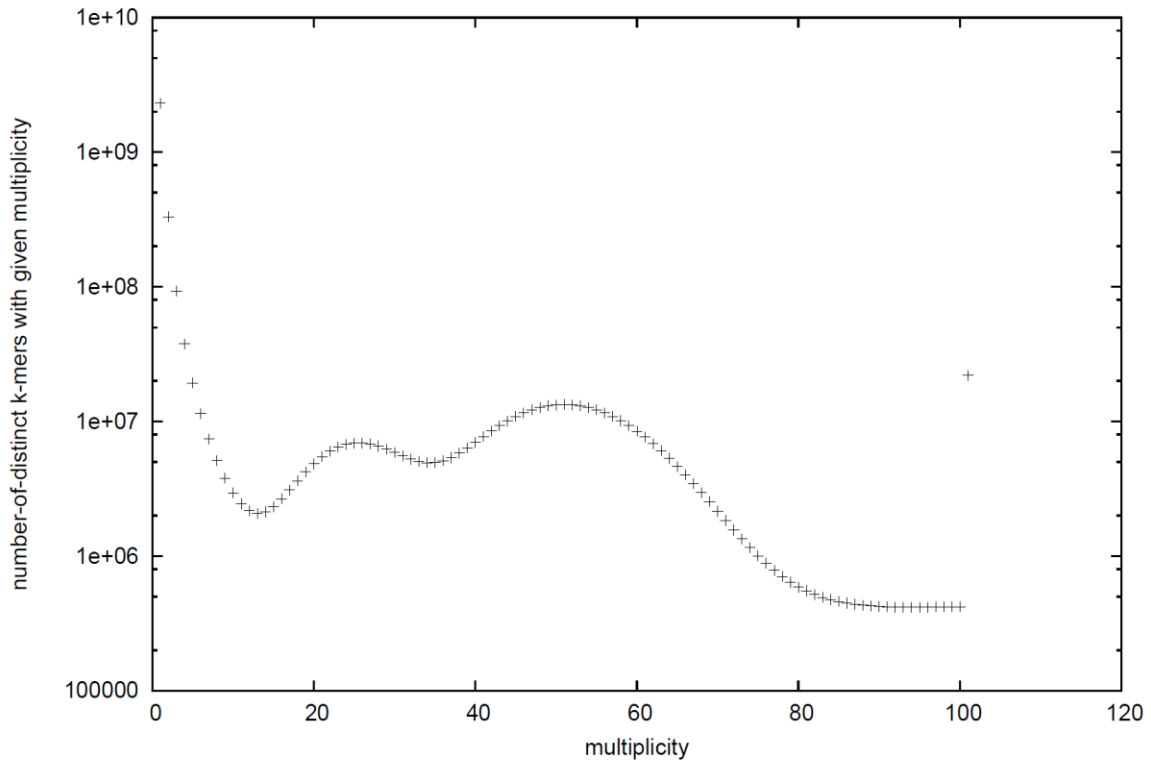
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

# 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

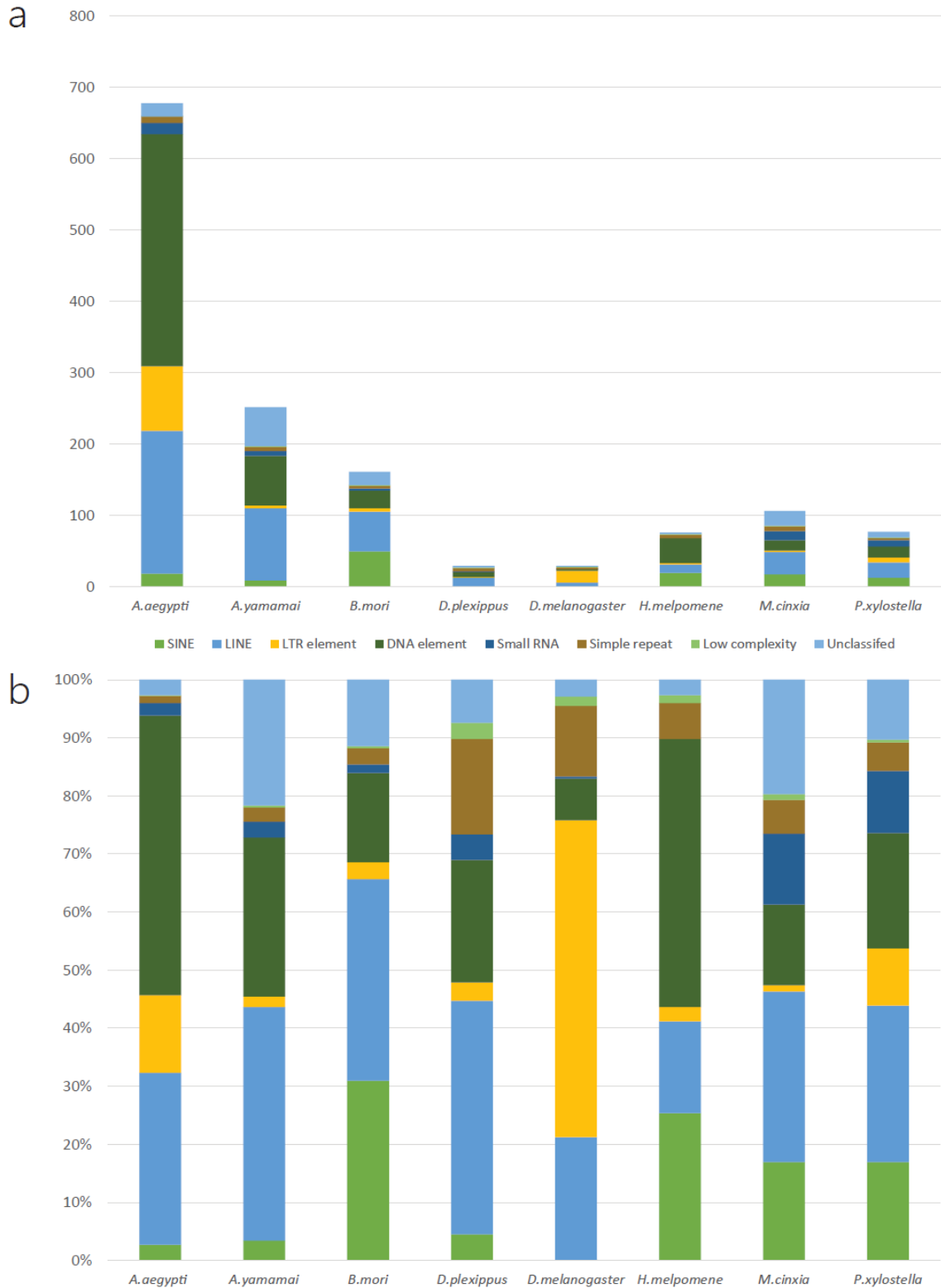


Figure 4. Constructed phylogenetic tree and comparative gene family analysis. Nodes value indicate Bayesian posterior probability, bootstrap and gene expansion, contraction value. Orange and blue color indicate expansion and contraction, respectively. Bar chart indicate the number of genes cauterized into 4 groups (Specific, 1:Multi, Multi:Multi and 1:1) using OrthoMCL.

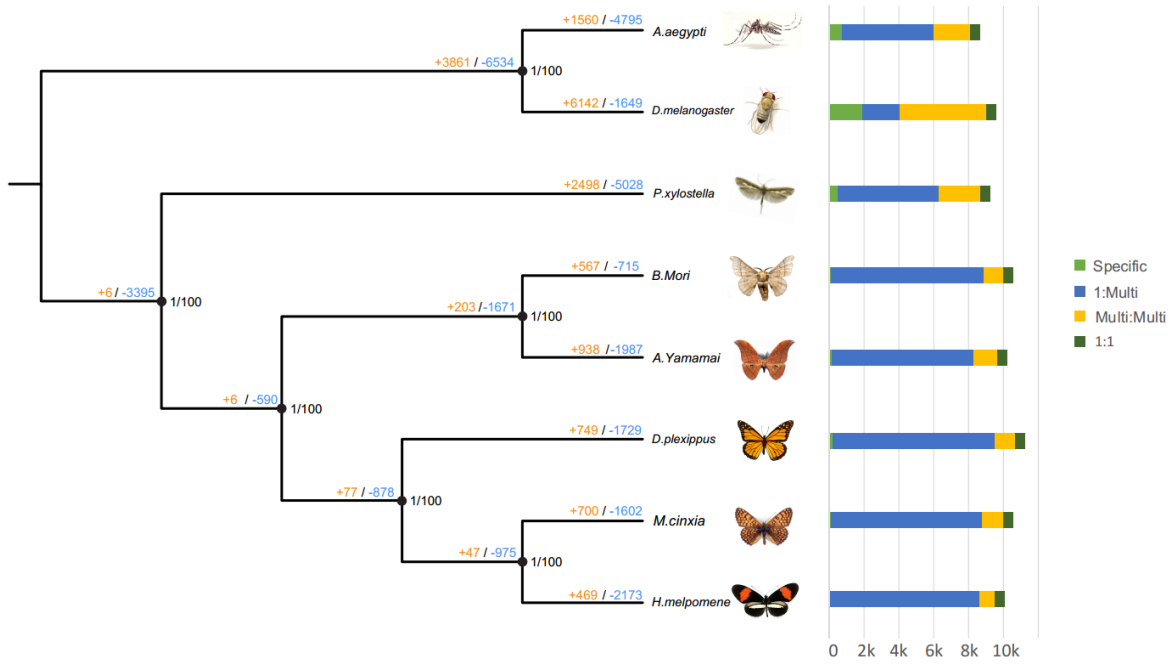
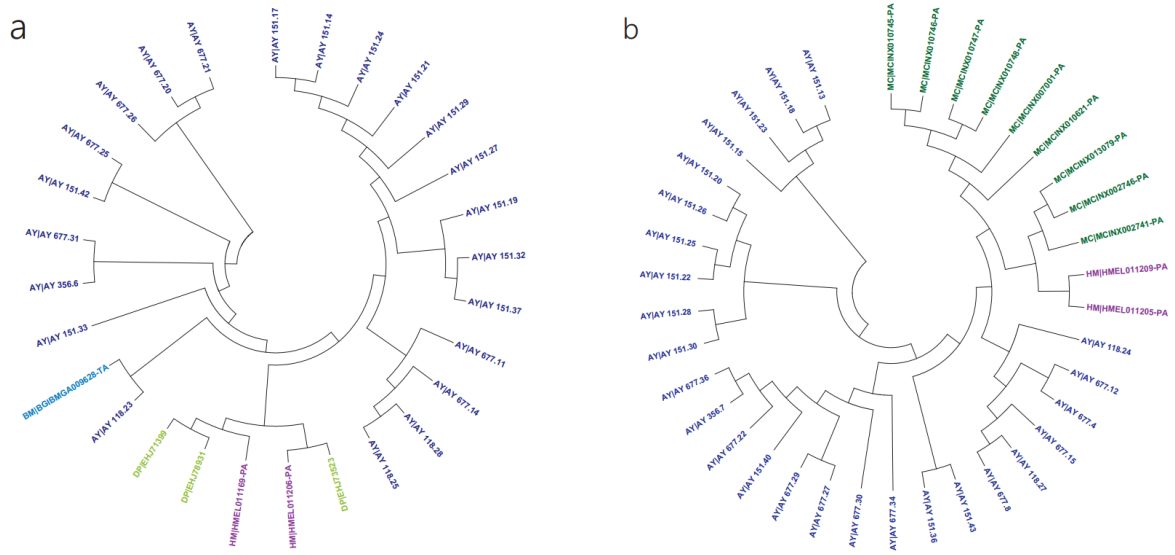
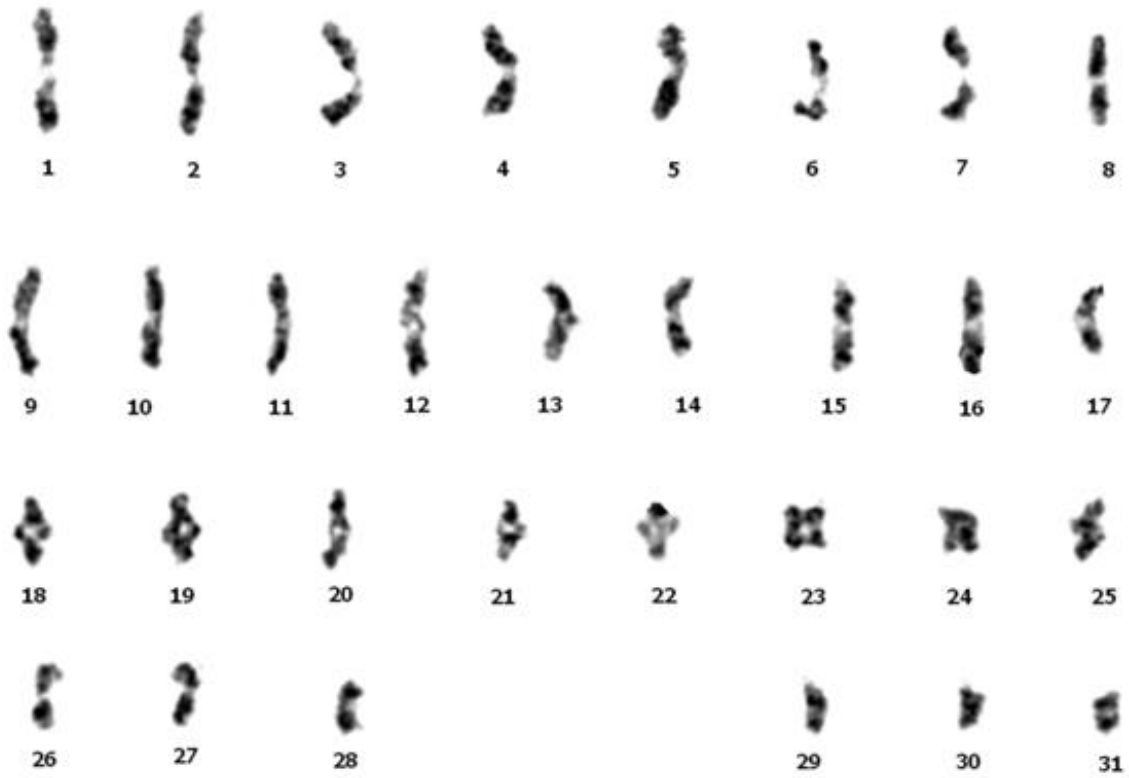


Figure 5. Expansion of chorion gene in *A.yamamai* genome. a and b shows the gene tree of chorion A and B in the rapid expanded gene family cluster, respectively. Color of terminal node indicates each taxon identified in the gene family cluster.



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



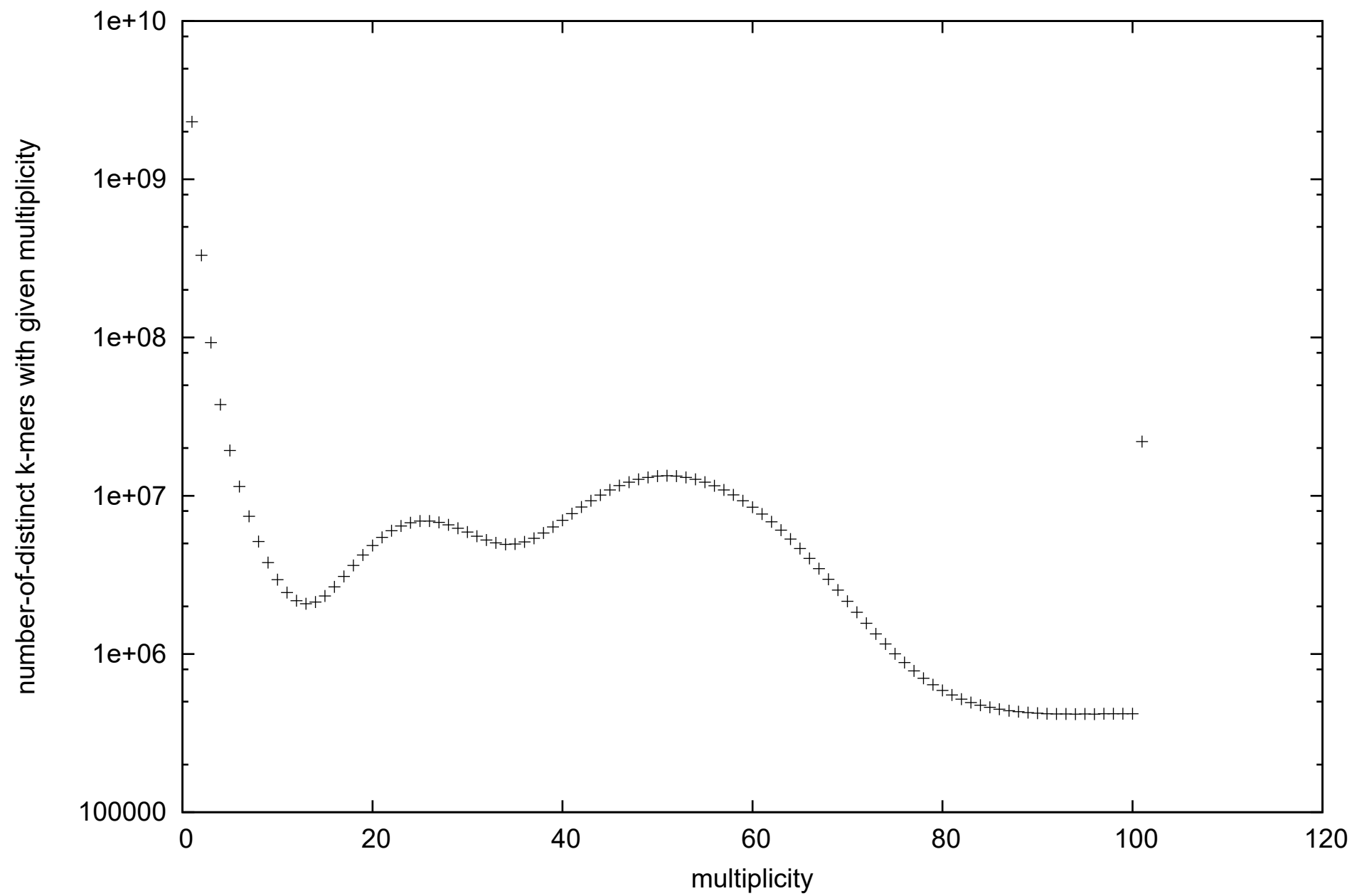
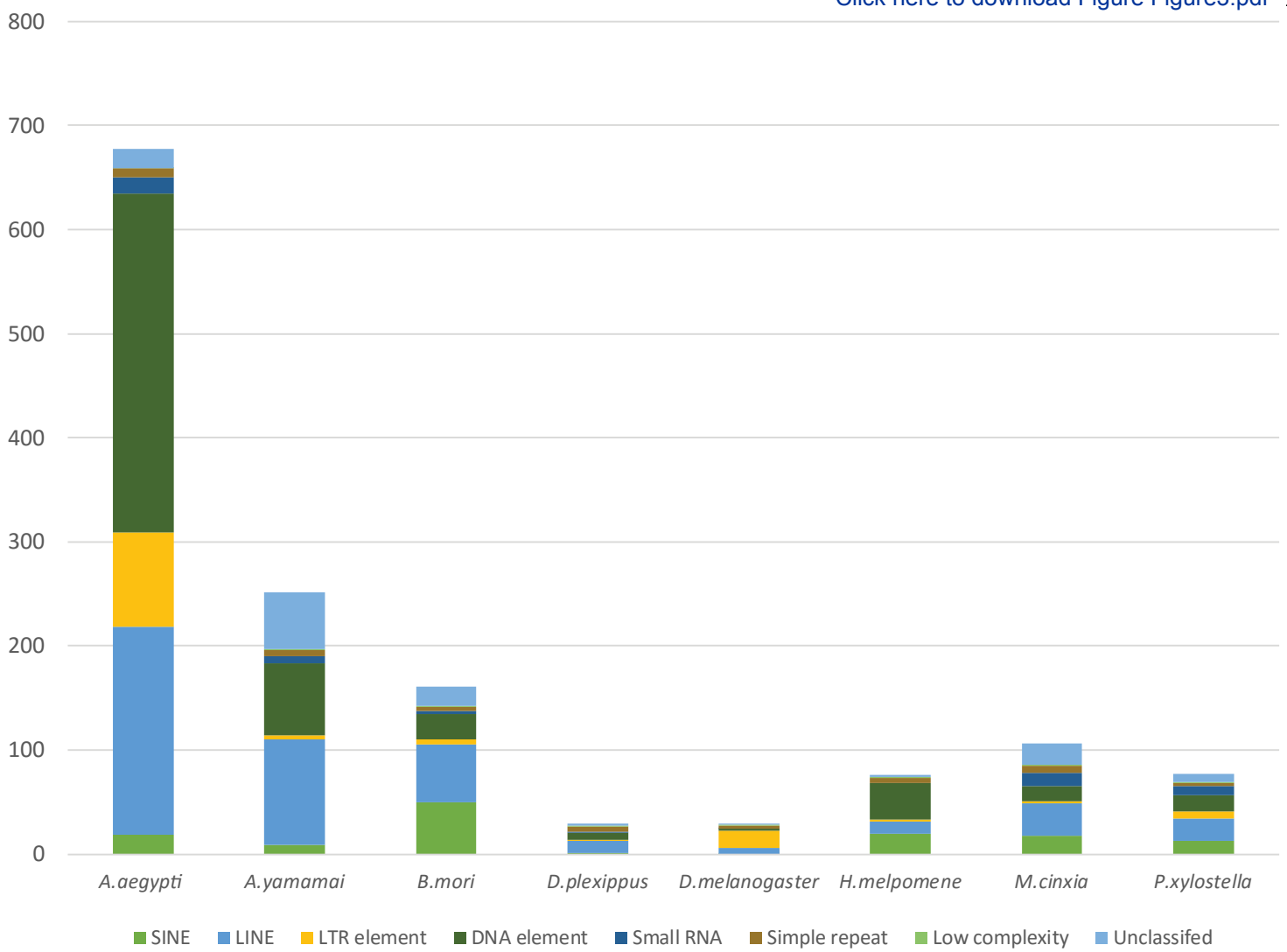
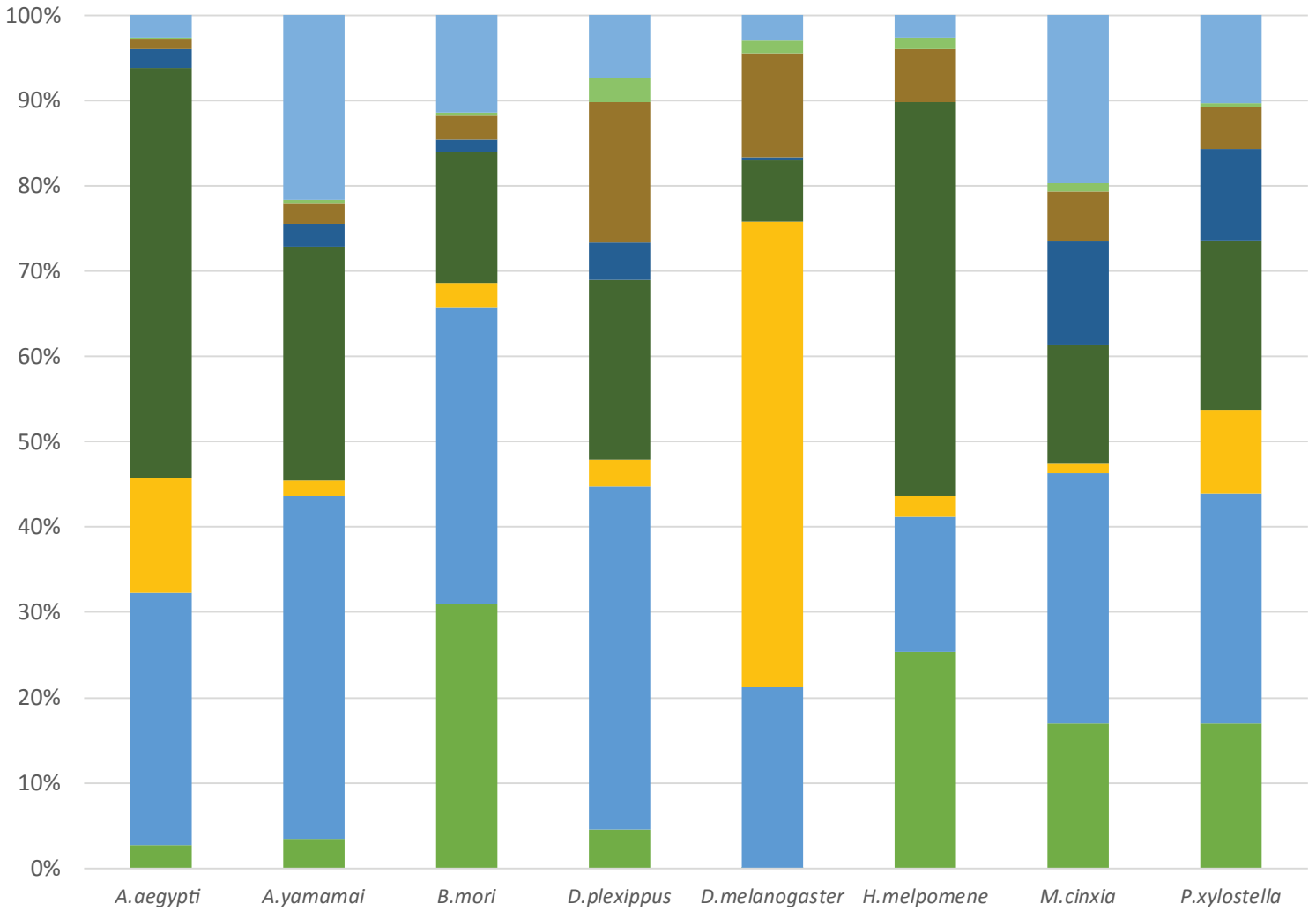


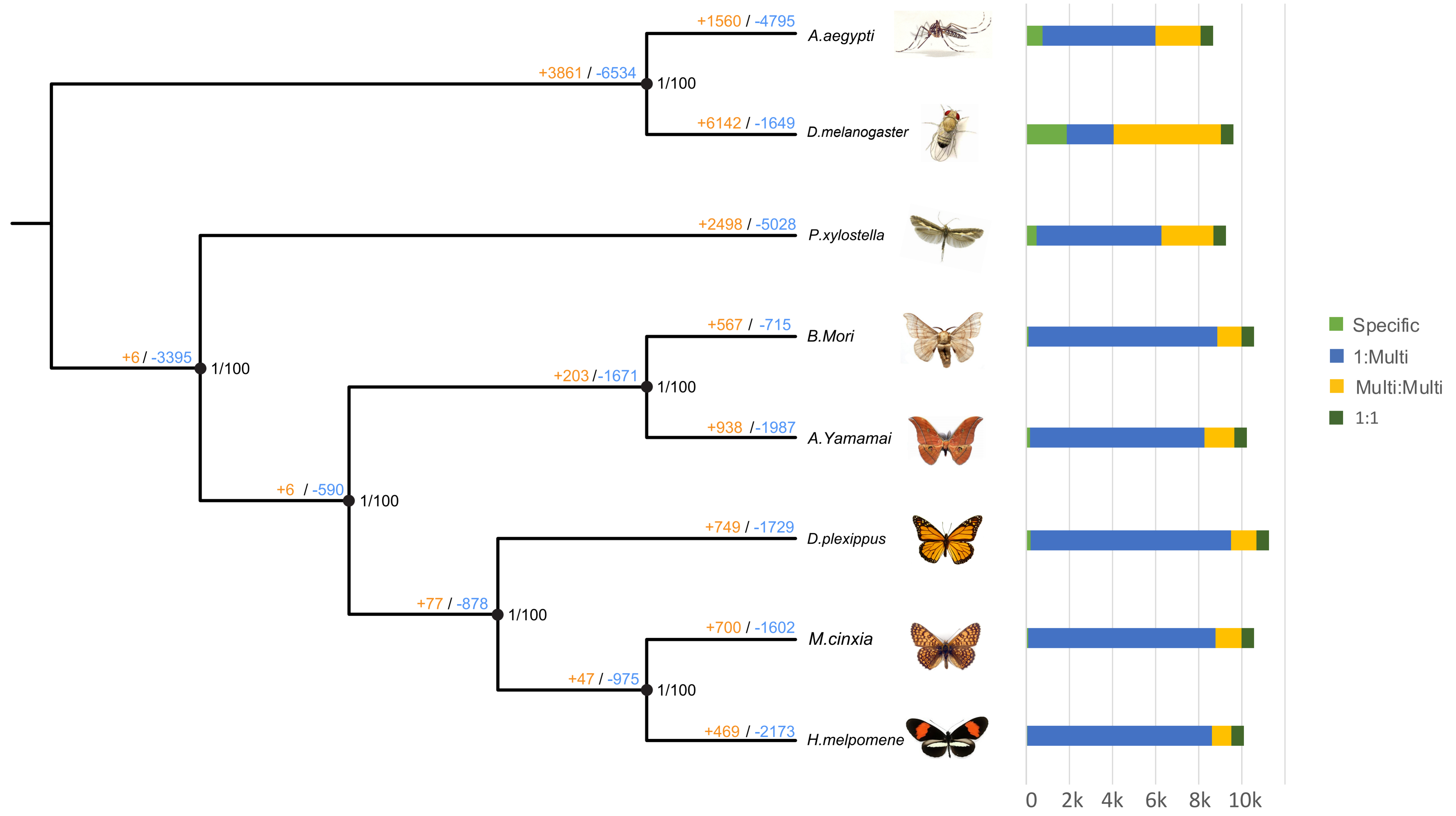


Figure 3  
a

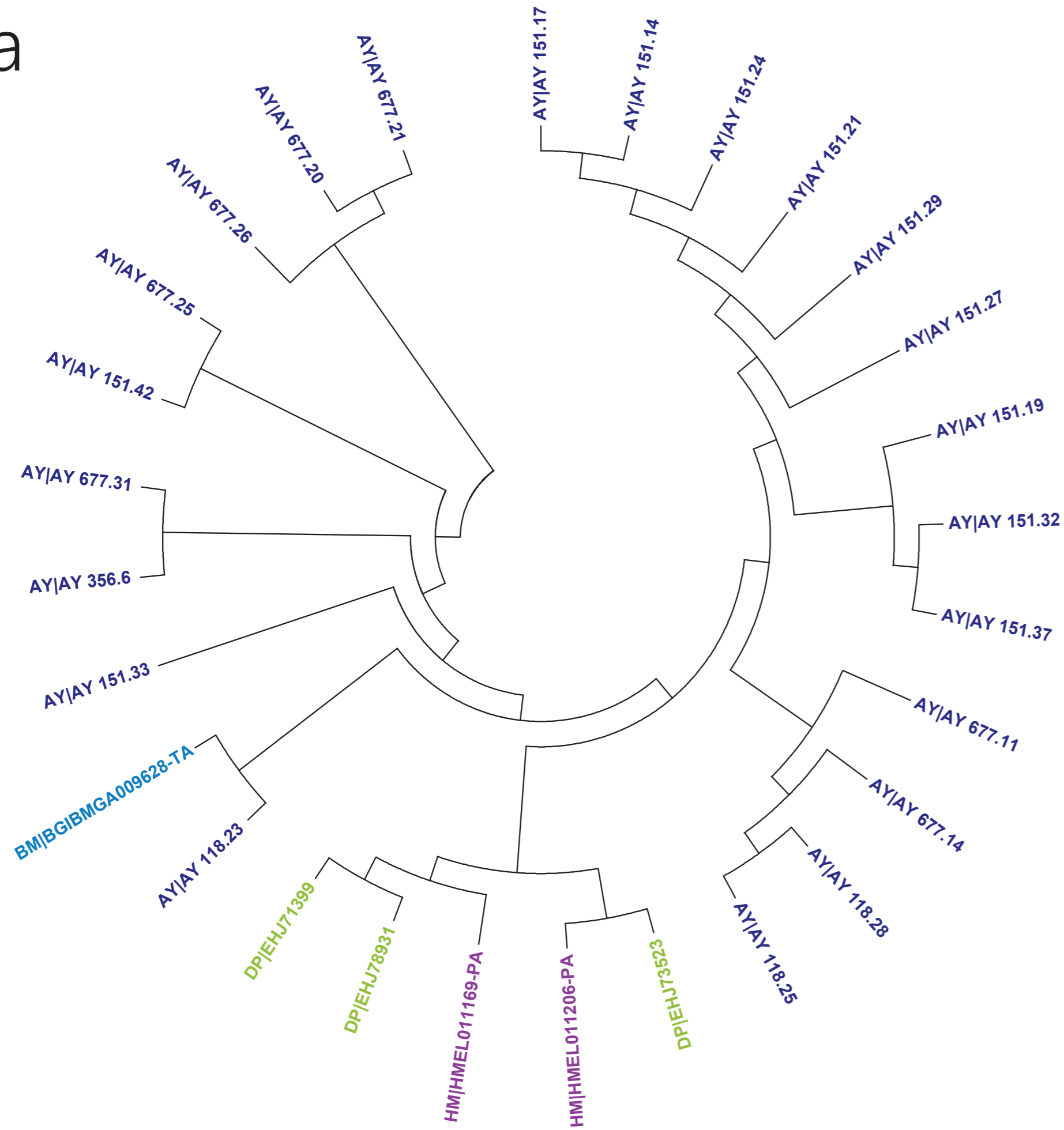


b

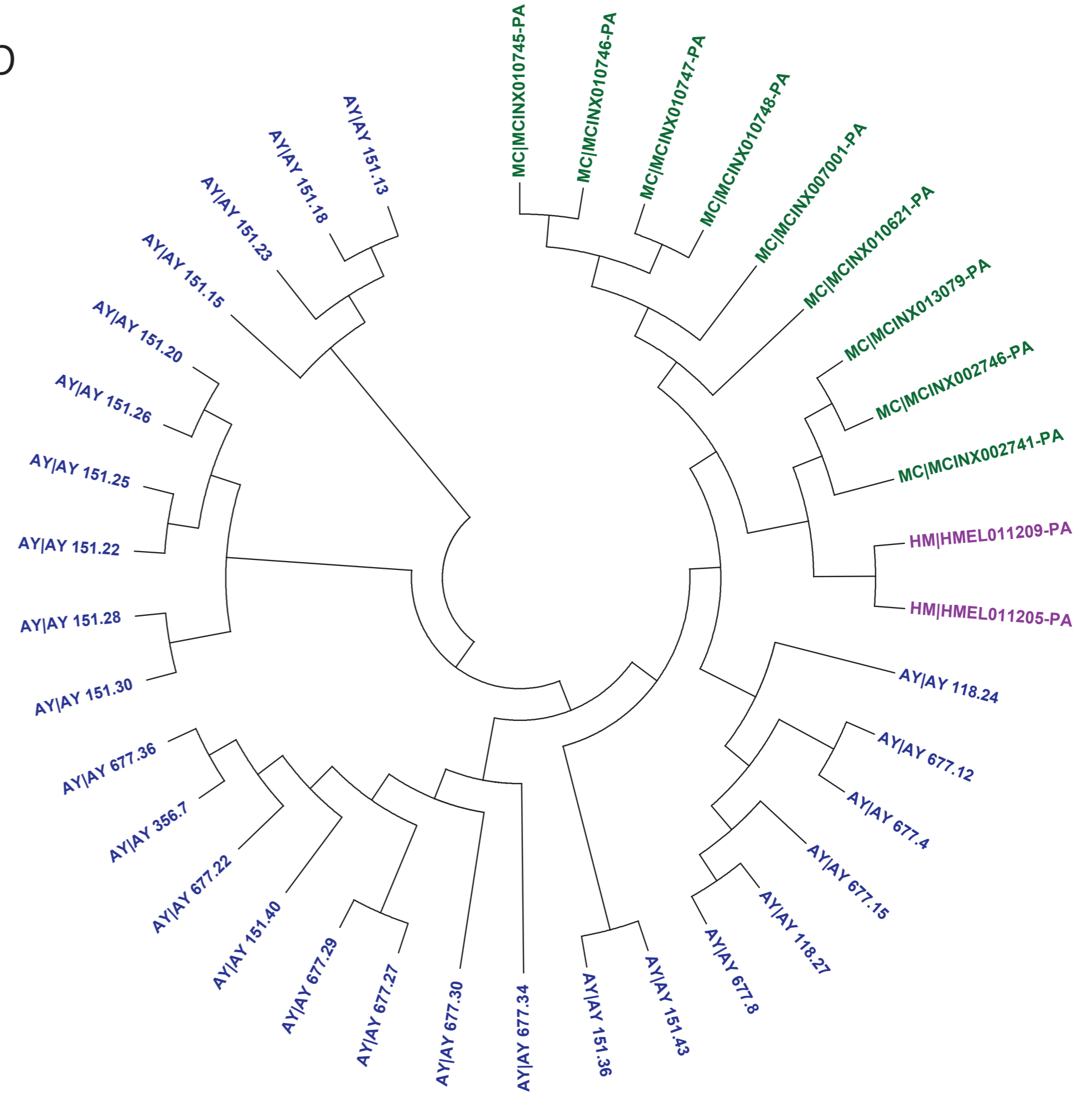


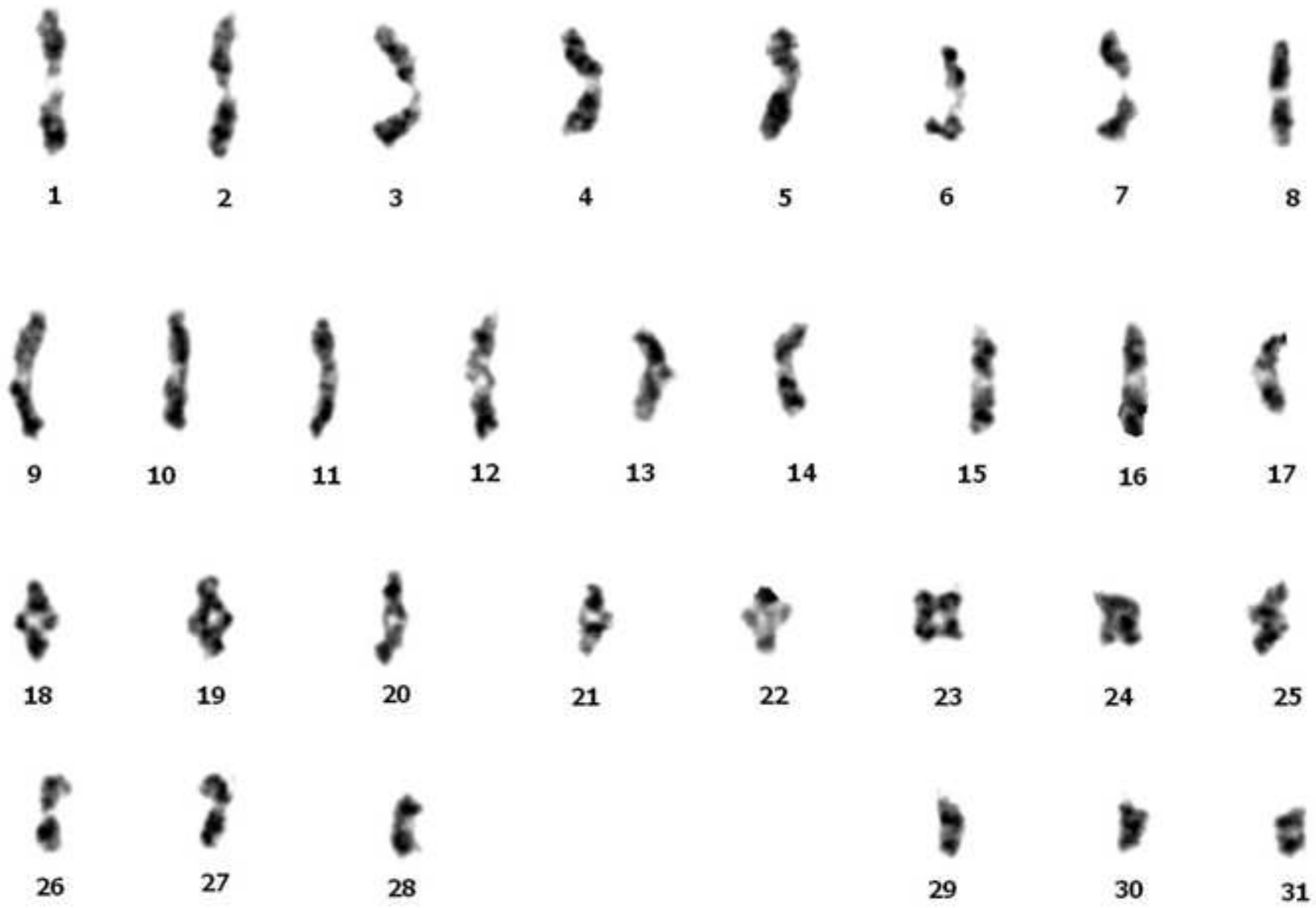


a



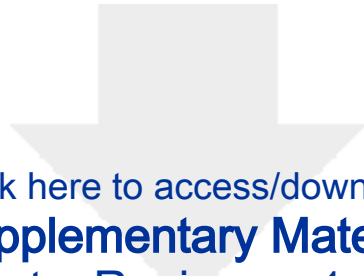
b



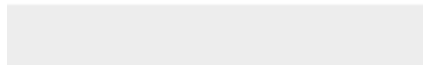


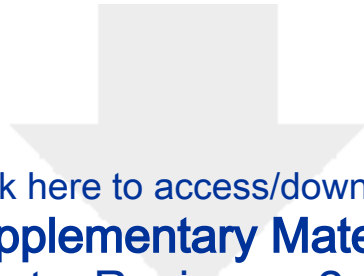






Click here to access/download  
**Supplementary Material**  
Response\_to\_Reviewer\_1\_ENG.docx





Click here to access/download  
**Supplementary Material**  
Response\_to\_Reviewer\_2\_ENG.docx





Sep 17, 2017

Dear Editor of *Gigascience*,

Thank you for sending our manuscript out for review and for obtaining constructive feedback from two expert referees. We appreciate the careful reading of our manuscript by the reviewers. We would also like to thank you for allowing us to submit a fully revised version that addresses all points of the reviewers and the academic editor.

We appreciate suggestions of all the reviewers on our manuscript; all the comments made by the reviewers were quite valid. We have responded to all comments point-by-point in the rebuttal.

Followed the reviewer's suggestion, we tried our best to conduct manual curation for all predicted genes (>20,000). We also changed all related analysis results in the manuscript and added one co-author who helped with this revision.

We hope that our revised manuscript, strengthened by reviewer comments, will meet the high-quality standard of *Gigascience*. We are always ready to strengthen our manuscript once again following the comments of the editors and reviewers at the next revision if necessary.

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](mailto:microsw@cu.ac.kr)