

1  
2  
3 1 The *Tetragnatha kauaiensis* genome sheds light on the origins of genomic novelty in spiders  
4  
5 2

6 3 José Cerca<sup>1,2,3,\*</sup>, Ellie E. Armstrong<sup>2,4</sup>, Joel Vizueta<sup>5,6</sup>, Rosa Fernández<sup>7</sup>, Dimitar Dimitrov<sup>8</sup>, Bent  
7 4 Petersen<sup>9,10</sup>, Stefan Prost<sup>11,12,13</sup>, Julio Rozas<sup>5</sup>, Dmitri Petrov<sup>4</sup>, Rosemary G. Gillespie<sup>2</sup>  
8  
9 5

10 6 1 - Frontiers in Evolutionary Zoology, Natural History Museum, University of Oslo, Oslo, Norway.  
11  
12 7 2 - Berkeley Evolab, Department of Environmental Science, Policy, and Management, UC Berkeley,  
13  
14 8 CA, USA.

15 9 3 - Holomuseumomics group, Department of Natural History, NTNU University Museum, Norwegian  
16  
17 10 University of Science and Technology, Trondheim, Norway.  
18

19 11 4 - Department of Biology, Stanford University, Stanford, CA, USA.

20 12 5 - Departament de Genètica, Microbiologia i Estadística & Institut de Recerca de la Biodiversitat  
21  
22 13 (IRBio), Universitat de Barcelona, Spain.

23 14 6 - Villum Centre for Biodiversity Genomics, Section for Ecology and Evolution, Department of  
24  
25 15 Biology, University of Copenhagen, Universitetsparken 15, 2100 Copenhagen, Denmark.

26 16 7 - Institute of Evolutionary Biology (CSIC - Universitat Pompeu Fabra). Passeig Marítim de la  
27  
28 17 Barceloneta 37-49. 08003 Barcelona, Spain.

29 18 8 - Department of Natural History, University Museum of Bergen, University of Bergen, P.O. Box  
30  
31 19 7800, 5020 Bergen, Norway.

32 20 9 - Section for Evolutionary Genomics, The GLOBE Institute, Faculty of Health and Medical  
33  
34 21 Sciences, University of Copenhagen, Øster Farimagsgade 5, 1353 Copenhagen, Denmark.

35 22 10 - Centre of Excellence for Omics-Driven Computational Biodiscovery, Faculty of Applied  
36  
37 23 Sciences, AIMST University, Kedah, Malaysia.

38 24 11 - Natural History Museum Vienna, Central Research Laboratories, Burgring 7, 1010 Vienna,  
39  
40 25 Austria.

41 26 12 - University of Veterinary Medicine, Konrad Lorenz Institute of Ethology, Savoyenstraße 1a, 1160  
42  
43 27 Vienna, Austria.

44 28 13 - South African National Biodiversity Institute, National Zoological Garden, Pretoria 0184, South  
45  
46 29 Africa.

47  
48  
49 30  
50 31 \*Corresponding author: jose.cerca@gmail.com  
51  
52  
53

## 32 Abstract

33 Spiders (Araneae) have a diverse spectrum of morphologies, behaviours and physiologies.  
34 Attempts to understand the genomic-basis of this diversity are often hindered by their large,  
35 heterozygous and AT-rich genomes with high repeat content resulting in highly fragmented, poor-  
36 quality assemblies. As a result, the key attributes of spider genomes, including gene family evolution,  
37 repeat content, and gene function, remain poorly understood. Here, we used Illumina and Dovetail  
38 Chicago technologies to sequence the genome of the long jawed spider *Tetragnatha kauaiensis*,  
39 producing an assembly distributed along 3,925 scaffolds with a N50 of ~2 Mb. Using comparative  
40 genomics tools, we explore genome evolution across available spider assemblies. Our findings  
41 suggest that the previously reported and vast genome size variation in spiders is linked to the different  
42 representation and number of transposable elements. Using statistical tools to uncover gene-family  
43 level evolution, we find expansions associated with the sensory perception of taste, immunity and  
44 metabolism. In addition, we report strikingly different histories of chemosensory, venom and silk  
45 gene families, with the first two evolving much earlier, affected by the ancestral whole genome  
46 duplication in Arachnospulmonata (~450 million years ago) and exhibiting higher numbers. Together,  
47 our findings reveal that spider genomes are highly variable and that genomic novelty may have been  
48 driven by the burst of an ancient whole genome duplication, followed by gene family and  
49 transposable element expansion.

## 51 Significance statement

52 Despite being one of the most charismatic animal lineages, progress on spider genome  
53 evolution lags due to the challenges in sequencing and assembling their genomes, which involve  
54 genome size and repeat content. Here, we sequence the genome of *Tetragnatha kauaiensis*, a spider  
55 endemic to Hawai'i, and compare it to other available spider genomes. We find variation in terms of  
56 repeats and transposable elements; expansions in gene-content associated with metabolism, sensory  
57 perception and immunity; and wide variation of chemosensory genes and venom genes.

## 59 Introduction

60 With nearly 50,000 described species ("World Spider Catalog," 2021), and dating back ~350  
61 million years (Fernández et al. 2018), spiders (Chelicerata, Araneae) have conquered most terrestrial  
62 ecosystems, from the cold Arctic to arid deserts (Jackson and Cross 2011; Dimitrov et al. 2012;  
63 Garrison et al. 2016; Fernández et al. 2018). Spiders play a key role in terrestrial ecosystems  
64 regulating community dynamics as major arthropod predators (Herberstein and Wignall 2011; Wilder  
65 2011), having evolved a diverse array of adaptive solutions, which include, a rich cocktail of venoms  
66 to neutralize prey (Binford 2001; King and Hardy 2013), a colour palette essential for camouflaging,  
67 mimicking and signaling (Oxford and Gillespie 1998; Croucher et al. 2013; Cotoras et al. 2016), and

1  
2  
3 68 the ability to produce silk for spinning webs and subduing prey (Vollrath 1999; Garb et al. 2010;  
4 69 Sanggaard et al. 2014).

6 70 Despite the advances in spider ecology, evolution and systematics, knowledge of spider  
7 71 genomes still lags relative to other taxa. Most of the available spider genomes are of poor quality,  
8 72 being highly fragmented (Garb et al. 2018) and lack a substantial part of the genome, with only three  
9 73 recent exceptions involving chromosome-resolved genomes (Escuer et al. 2021; Fan et al. 2021;  
10 74 Sheffer et al. 2021). Several factors contribute to the sparse availability of high-quality spider genome  
11 75 assemblies, including the lack of a model organism among spiders (*sensu Drosophila melanogaster* in  
12 76 flies and *Tribolium castaneum* in beetles) (Brewer et al. 2014), and the challenges associated with  
13 77 sequencing spider genomes, which are characterized by high AT-content, repeats, heterozygosity, and  
14 78 often large genome sizes (Garb et al. 2018). Focus on non-model organism genomes shows that  
15 79 increased taxon-sampling leads to an improved understanding of the diversity and function of  
16 80 molecular mechanisms across the tree of life (McGregor et al. 2008), as it overcomes the biases from  
17 81 the limited number of model taxa, and highlights the idiosyncrasies throughout the tree of life.  
18 82 Consequently, a better representation of spider genomes will certainly help understanding spider  
19 83 diversity and evolution (McGregor et al. 2008).

20 84 A systematic analysis of spider genomes has the potential to unveil the genomic foundation of  
21 85 spider evolution. For example, the detection of duplicate Hox clusters suggested an ancestral whole  
22 86 genome duplication in the common ancestor of modern spiders and scorpions (Arachnoplumonata;  
23 87 Schwager et al. 2007), and this evidence was later on confirmed by the first spider genomes (Clarke et  
24 88 al. 2015; Schwager et al. 2017; Leite et al. 2018). The implications of whole genome duplications  
25 89 may, however, be multifarious and complex (Ohno 1970). On one hand, genome duplication may act  
26 90 as a catalyst for molecular novelty. Under this framework, the retention of duplicated genes and other  
27 91 genetic components may act as ‘reservoirs of genetic variation’, through processes of gene neo- and  
28 92 sub-functionalization (Lynch and Force 2000), and be of use when organisms encounter novel  
29 93 selective pressures (Li et al. 2018; Nieto Feliner et al. 2020; Schmickl and Yant 2021). Considering  
30 94 the evidence for gene duplicates in spider genomes, including spidroins (silk genes) (Sanggaard et al.  
31 95 2014; Clarke et al. 2015; Babb et al. 2017; Garb et al. 2018; Sheffer et al. 2021), venoms (Sanggaard  
32 96 et al. 2014; Gendreau et al. 2017; Haney et al. 2019), chemosensory (Vizueta et al. 2018; Vizueta et  
33 97 al. 2019; Vizueta, Escuer, et al. 2020) gene families may yield insights on phenotypic innovation and  
34 98 the adaptation to novel environments. On the other hand, since genome duplication leads to a  
35 99 significant re-organization of the genome, it may cause deregulation of gene-expression networks or  
36 100 unlock the epigenetic suppression of transposable elements, which may proliferate across the genome  
37 101 and result in decreased fitness for the organism – ‘the genomic shock hypothesis’ (McClintock 1984;  
38 102 Choi et al. 2020). In such a scenario, one expects to find variation in transposable element  
39 103 proliferation across genomes, and ultimately a substantial variation of genome size. The proliferation  
40 104 of transposable elements may thereby underlie genome size variation in spiders, which ranges

1  
2  
3 105 between 0.74 - 5.73 C values (0.7 Gb - 5.6 Gb) (Gregory and Shorthouse  
4 106 2003)(<http://www.genomesize.com/> checked in April 15th 2021; values for: *Habronattus borealis*,  
5 107 *Tetragnatha elongata*, respectively). Comparisons between different genome assemblies may yield  
6 108 important insights on the prevalence of gene duplications, neofunctionalization, and transposable  
7 109 element dynamics across different lineages.

8 110 Here, we report a genome assembly of the Hawaiian spider *Tetragnatha kauaiensis* and place  
9 111 it in the context of currently available spider genomes to assess signatures of genome evolution across  
10 112 spider lineages (Supplementary Table 1). To do so, we first explore the completeness and duplication  
11 113 rates across the spider assemblies. Considering the role of transposable elements in driving genome  
12 114 size variation, we also assess transposable element load in each genome. Third, we quantify the  
13 115 expansion and contraction of gene families (based on gene similarity metrics), and classify the  
14 116 function of these families using Gene Ontology (GO). Finally, we delve deeper into the identification  
15 117 and comparison of chemosensory, venom and spidroin (silk) genes, studying duplicates in a  
16 118 phylogenetic context. Focus on these three categories is grounded on their central role to the survival  
17 119 and fitness of spiders, and benefits from extensive research, including hand curated genes and  
18 120 databases.

## 121 122 **Results**

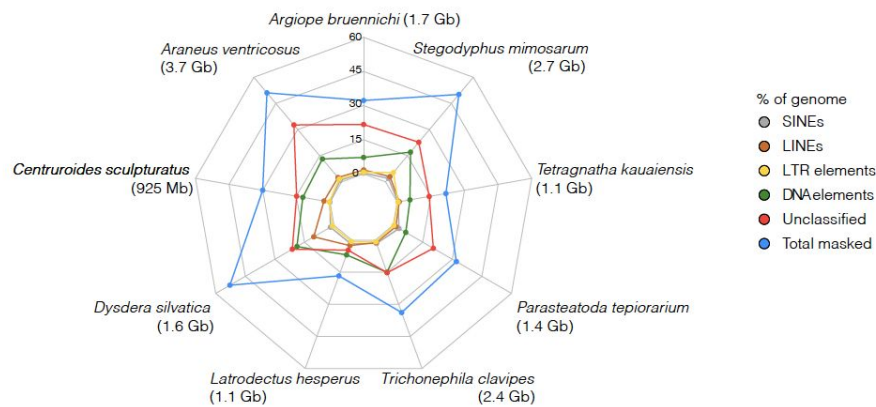
### 123 **The *Tetragnatha kauaiensis* genome**

124 The *T. kauaiensis* genome assembly has a size of ~1.08 gigabases (Gb), distributed along a  
125 total of 132,391 contigs, comprising 3,925 scaffolds. The largest scaffold is ca. 10.5 megabases (Mb),  
126 while the estimated scaffold-N50 for the assembly is ~2 Mb (Supplementary Table 2). The assembly  
127 has a GC content of ~33.3%, in line with the remaining spider genomes (lowest GC content *L.*  
128 *hesperus* with 28.59% and highest content is *S. mimosarum* with a GC content of 33.62;  
129 Supplementary Table 2). The assembly has no obvious contaminants or associated symbionts, as  
130 determined by Blobtools (Supplementary Figure 1). The majority of scaffolds have a similar GC  
131 composition, despite variations in coverage. From all 3,925 scaffolds, 2,774 were labelled as no-hits  
132 (comprising only a total of ~32.46 Mb of the assembly), and 889 labelled as Arthropods (~886 Mb).

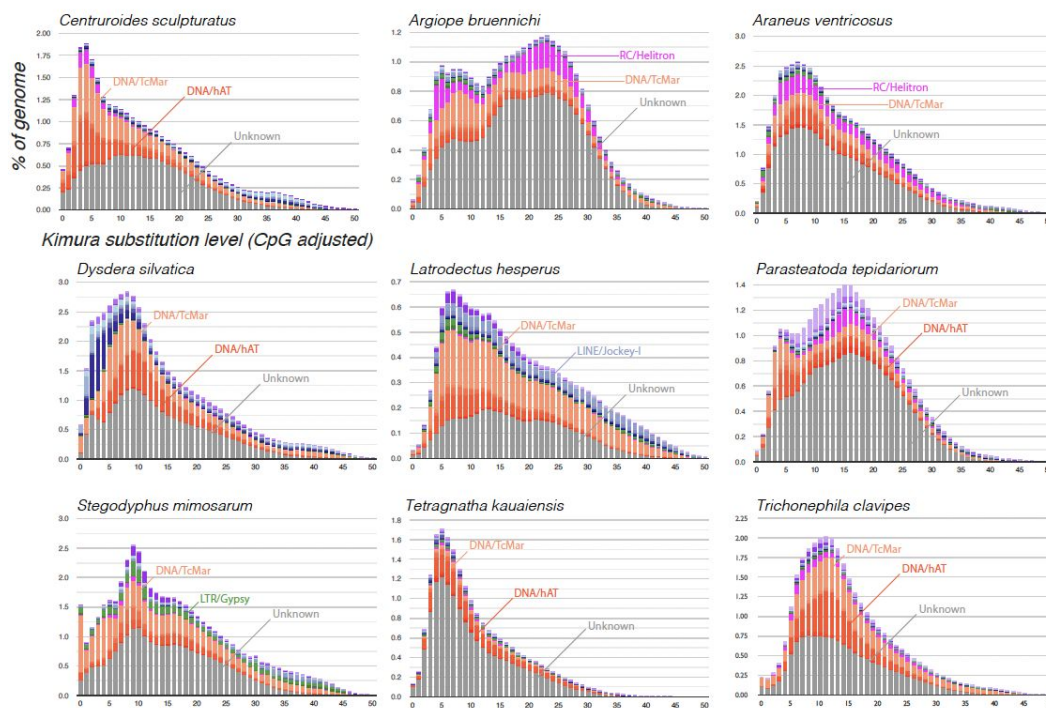
133 Annotation of the *Tetragnatha kauaiensis* genome yielded 38,907 genes, comprising 213,695  
134 exons and 171,423 introns (Supplementary Table 3). Together, all genes cover 290,369,064 bp (290  
135 Mb) representing 26.7% of the genome with 41,209,078 bp (41 Mb, 3.8% of the genome) being  
136 coding sequences (cds). The mean gene length is 7,463 bp (Supplementary Table 3), the longest gene  
137 is 208,580 bp long (208 kb), and 89.7% of BUSCOs are retrieved as complete.

138

A



B



139

140

141

142

143

144

145

146

147

148

149

150

**Figure 1: Transposable element (TE) and repeat characterization** **a)** Web diagram showing the representation of TE and repeats in the assemblies. Assemblies and correspondent assembly sizes are represented on the edges of the web diagram. Different transposable element families or repeats are presented in different colours on the plot, and the total content masked by RepeatMasker is shown in blue. The numbers for each web-line indicate the % of the genome occupied by each transposable element, or the % masked. **b)** Repeat/transposable element landscape plots for the various assemblies. The three most represented transposable element categories are present for every genome (e.g. DNA/TcMar, DNA/hAT, and unknown for *T. kauaiensis*). Each plot shows the Kimura substitution level (x axis) and % of genome covered by repeats (y axis).

**Genome characterization and evolution**

1  
2  
3 151 The analyzed assemblies vary widely in size. *A. ventricosus* has the largest assembly with 3.6  
4 152 Gb (Supplementary Table 2), while *T. kauaiensis* has the smallest assembly with 1,085,571,486 bp  
5 153 (1.1 Gb). In between these extremes, we find the genomes of *S. mimosarum* (2.7 Gb), *T. clavipes* (2.4  
6 154 Gb), *A. bruennichi* (1.7 Gb), *D. silvatica* (1.4 Gb), *P. tepidariorum* (1.5 Gb) and *L. hesperus* (1.1 Gb).

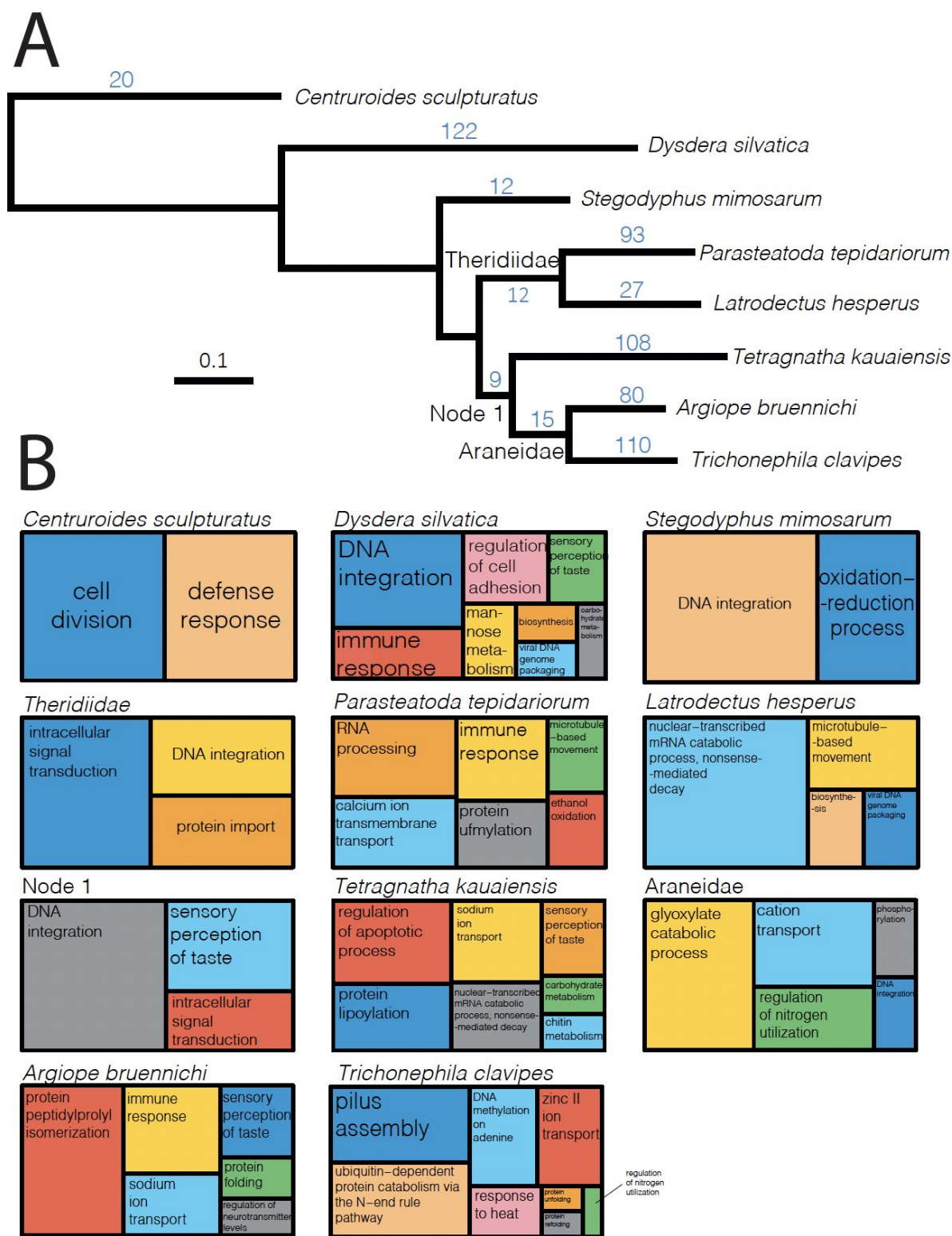
7  
8  
9 155 Considering the 3-fold variation in genome size and the evidence for ancient whole genome  
10 156 duplications in Chelicerata (Shingate et al. 2020) and Arachnida (Schwager et al. 2017; Harper et al.  
11 157 2020), and the suggestion that there has been a large-scale (whole genome or chromosomal)  
12 158 duplication event within spiders (Clarke et al. 2015), we explored the possibility of whole genome  
13 159 duplication private to spider genomes by interrogating the number of homologs in the Hox genes  
14 160 clusters. Using Hox genes 1-5, and based on a threshold of 95% identity, we find no evidence for an  
15 161 additional ancestral whole genome duplication in the studied spider genomes. We found zero, one or  
16 162 two homologs for Hox 1 (Supplementary Table 4). For Hox 2, we found two homologs in all  
17 163 genomes, with the exception of *A. ventricosus*, where we only find a single homolog (Supplementary  
18 164 Table 4). For Hox 3, there was only one homolog in all genomes, with the exception of *P.*  
19 165 *tepidariorum* (2 candidates) and *T. clavipes* (no candidate). For Hox 4, we found two homologous  
20 166 genes in *T. kauaiensis*, *P. tepidariorum*, *L. hesperus* and *S. mimosarum*, one in *T. clavipes* and another  
21 167 in *D. silvatica*. *A. ventricosus*, however, had four homologs for the Hox4 gene. Finally, for Hox 5, we  
22 168 identified one homolog in all genomes, with the exception of *A. ventricosus* and *P. tepidariorum*  
23 169 where we found two homologous genes. This suggests that, with the exception of the outlier with four  
24 170 copies (*Araneus* Hox4), Hox genes are present in 1 or 2 copies.

### 171 **Transposable element variation**

172 We find variation in repeat content and tempo of repeat accumulation across the spider  
173 assemblies (Figure 1; Supplementary Table 5). For example, 10.3% of the *D. silvatica* genome is  
174 composed of LINES, whereas all other studied spiders had at most 3% LINES (Figure 1A). *S.*  
175 *mimosarum* had 5.40% of its genome covered by LTR elements, while *A. ventricosus*, which is the  
176 second LTR-element most rich genome, had only 1.60% (Figure 1). Interspersed repeats varied  
177 between 52.84% in *D. silvatica* and 16.53% in *L. hesperus* (Supplementary Table 5). Unclassified  
178 repeats ranged between 32.64% (*A. ventricosus*), and 4.71% *L. hesperus* (Supplementary Table 5).  
179 Overall, Repeatmasker identified between 16.71% -52.84% of total repeat content (Figure 1 A;  
180 Supplementary Table 5). The correlation coefficient (R) between genome size and the % of masked  
181 genome is R=0.65, and the correlation coefficient (R) between total length of the masked genome and  
182 genome size is R=0.962. Finally, we find variability in the accumulation of transposable elements  
183 through time, as represented by the shape of the transposable element/repeat landscape plot curves  
184 (Figure B). For instance, the *A. bruennichi* and *P. tepidariorum* assemblies show two peaks in  
185 transposable element accumulation, whereas all the others display a single peak. *S. mimosarum*,  
186 however, has a recent burst in Tc1/mariner (DNA/TcMar) transposable elements (Figure 1B). Despite  
187 the differences in the accumulation of transposable element/repeats through time, we note that the

1  
2  
3 188 Tc1/mariner group (DNA/TcMar) is present as one of the top three most represented transposable  
4  
5 189 elements in all the assemblies, and and hAT transposons (DNA/hAT) are also among the three-  
6  
7 190 dominant categories in 6 assemblies. There is, however, variation across assemblies, as shown by the  
8  
9 191 high numbers of Helitrons (RC/Helitron) in two of the Araneidae assemblies (*A. bruennichi* and *A.*  
10  
11 192 *ventricosus*), Gypsy (LTRGypsy) in *S. mimosarum*, and Jockey (LINE/Jockey-1) in *L. hesperus*.

12  
13 193 The analysis of genome completeness, as assessed by BUSCO scores, suggests that spider  
14  
15 194 assemblies are considerably fragmented and missing substantial parts of the genome (Supplementary  
16  
17 195 Table 6). For instance, the *D. silvatica*, *L. hesperus* and *T. clavipes* genomes have only, respectively,  
18  
19 196 66%, 38.6% and 52% complete BUSCOs (Arachnid odb10). Completeness in the remaining genomes  
20  
21 197 ranged between 80-99%. Duplicated BUSCOs ranged between 30.5% (*P. tepidariorum*) and 3.2% (*S.*  
22  
23 198 *mimosarum*). Notably, the two biggest genomes, *A. ventricosus* (3.6Gb) and *S. mimosarum* (2.7Gb)  
24  
25 199 have 18.4% and 3.2% duplicated BUSCOs (Supplementary Table 6, Arachnid dataset odb10). The  
26  
27 200 percentage of complete single-copy, duplicated, fragmented, and missing BUSCOs is concordant  
28  
29 201 between the Arthropod and Arachnid sets (Supplementary Table 6).  
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



202  
 203 **Figure 2. Gene family expansion** A) Tree topology obtained for single-copy orthologs. Numbers in  
 204 blue indicate significantly expanded gene families as determined by CAFE. B) Treemap  
 205 representation of Gene Ontology Biological Function Annotation of the significantly expanded gene



206 families as retrieved by REVIGO. Branches/Nodes with significant expansions, including Araneidae,  
207 Theridiidae, and Node 1 are represented together with the different genomes.

208

### 209 **Gene-family evolution**

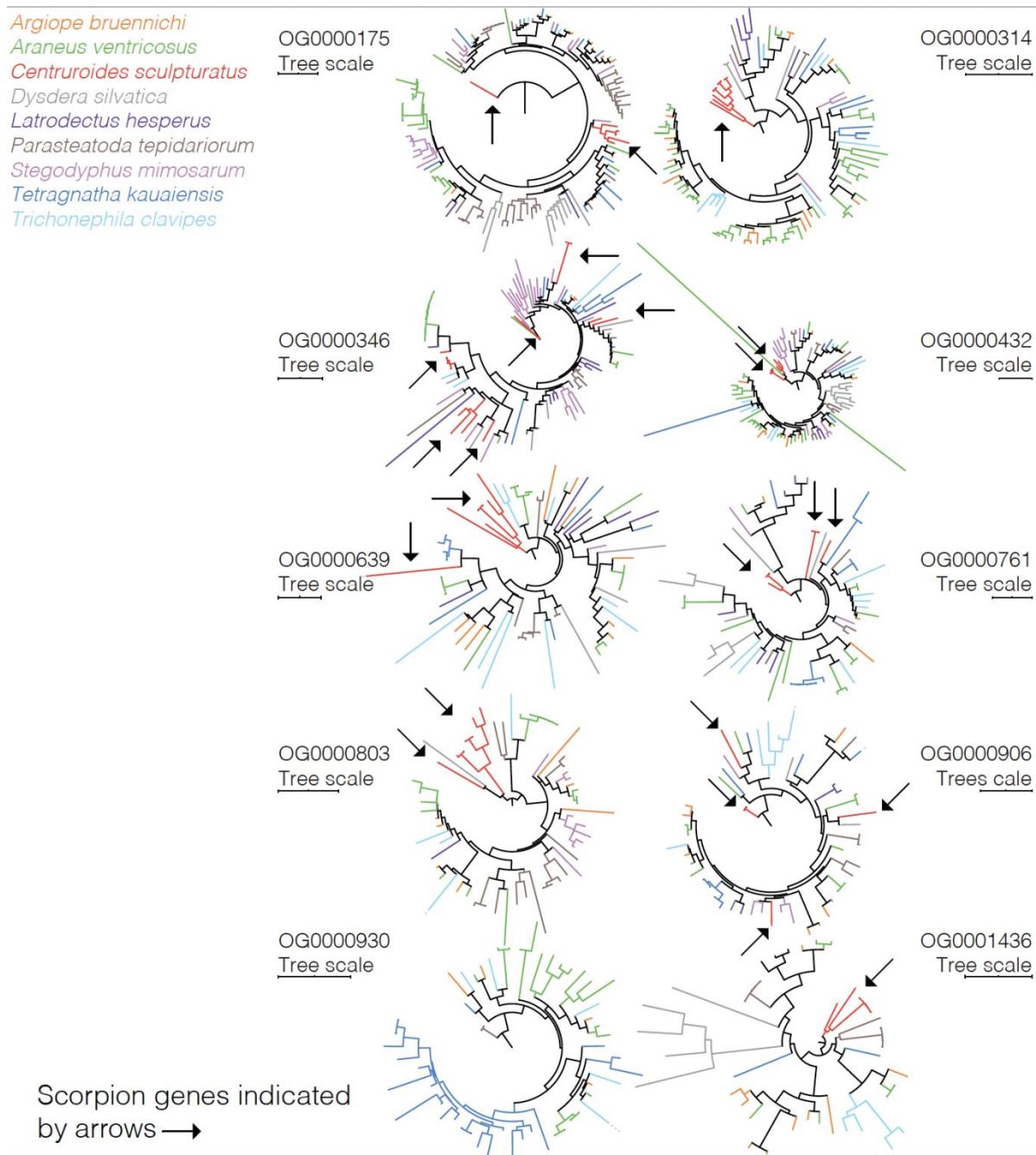
210 Since studying gene family evolution requires a phylogenetic backbone, we used the tree  
211 obtained from OrthoFinder based on 286 single-copy orthologs (orthologs are genes in different  
212 species that evolved from a common ancestral gene; Figure 2A). The tree topology has *T. kauaiensis*  
213 (Tetragnathidae) as sister lineage to the clade comprising the two members of Araneidae (*A.*  
214 *bruennichi* and *T. clavipes*). The clade encompassing all the aforementioned is sister to the  
215 Theridiidae (*L. hesperus* and *P. tepidariorum*). In turn, *S. mimosarium* (Eresidae) is the sister to  
216 Araneoidea (represented here by Tetragnathidae, Araneidae and Theridiidae). *D. silvatica*  
217 (Dysderidae) is the sister to the clade comprising all the aforementioned spiders (Figure 2A). This  
218 topology is in agreement with recent and comprehensive phylogenomic analyses of spiders  
219 (Fernández et al. 2018).

220 From a total of 608 significant gene family expansions in all branches, 572 occurred in  
221 terminal branches (Figure 2B). There were 451 significant expansions, and 157 significant  
222 contractions, of which 124 occurred in terminal branches (Supplementary Figures 1-4).

223 GO annotations of the significantly expanded gene families which were characterized under  
224 ‘biological process’ were organized by REVIGO and are represented in Figure 2B. Broadly, we find  
225 expansions associated with feeding metabolism and sensory perception, mannose metabolism in the  
226 genome of *D. silvatica* and chitin metabolism in *T. kauaiensis* (Figure 2B). Expansions in  
227 carbohydrate metabolism are found in *D. silvatica* and *T. kauaiensis*, while Araneidae has glyoxylate  
228 catabolic process expanded (Figure 2B). Expansions in sensory perception of taste are found in *D.*  
229 *silvatica*, *T. kauaiensis*, *A. bruennichi*, and in Node 1 (Figure 2B). Immune response is found in the  
230 genomes of *D. silvatica*, *P. tepidariorum* and *A. bruennichi*, while sodium ion transport is found in *T.*  
231 *kauaiensis* and *A. bruennichi* (Figure 2B).

232 When considering significant expansions in all GO categories (i.e. biological process,  
233 molecular function and cellular component), we find expansions associated with taste (including  
234 sensory perception of taste in Node 1, *A. bruennichi*, and *D. silvatica*; detection of chemical stimulus  
235 involved in sensory perception of taste in *A. bruennichi* and Node 1; molecular function taste receptor  
236 activity; is found in *A. bruennichi* and *T. kauaiensis*; Supplementary Table 7). We also find evidence  
237 for expansions related to various metabolic processes, including carbohydrate metabolic process, and  
238 mannose metabolic process in *D. silvatica*, while protein catabolic process, 3,4-dihydroxybenzoate  
239 catabolic process, fatty acid catabolic process, pyruvate metabolic process, glucose metabolic process,  
240 protein metabolic process, lipid catabolic process, lipid metabolic process, and fatty acid metabolic  
241 process are found in *T. clavipes*. The *P. tepidariorum* genome includes expansions in peptidoglycan  
242 catabolic process and lipid metabolic process, while that of *T. kauaiensis* includes expansions in chitin

1  
2  
3 243 metabolic process, carbohydrate metabolic process. Theridiidae includes expansions in lipid  
4 244 metabolic process, while Araneidae includes changes in taurine catabolic process. Finally, catalytic  
5 245 activity, is expanded in the genomes of *D. silvatica*, *L. hesperus*, *T. clavipes*, *T. kauaiensis*. Other  
6 246 notable expansions include the regulation of neurotransmitter levels, structural constituent of eye lens  
7 247 in *A. bruennichi*, defence response and toxin activity in *C. sculpturatus*, and response to heat in *T.*  
8 248 *clavipes*. The biological process for ‘sodium channel activity’ is found expanded in *A. bruennichi*, *T.*  
9 249 *clavipes* and *P. tepidariorum*, while the molecular function for ‘sodium channel activity’ is found in  
10 250 *A. bruennichi* and *T. kauaiensis*. Proteolysis (i.e. breakdown of proteins), the break down of process  
11 251 is expanded in *A. bruennichi*, *C. sculpturatus*, *D. silvatica*, *L. hesperus*, *P. tepidariorum*, *S.*  
12 252 *mimosarum*, *T. kauaiensis* and Theridiidae.  
13  
14  
15  
16  
17  
18  
19 253  
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



254

255 **Figure 3. Venom gene phylogenies.** Phylogenies for the 10 largest orthogroups of identified venom

256 genes. For each tree we indicate the Orthogroup ID and tree scale. Different colours correspond to

257 different species, as displayed in the legend. Arrows highlight scorpion toxin genes, and show that

258 most orthogroups in were already present in before the split between scorpions and spiders.

259

260 **Venom gene-family variation**

261 The combination of BLAST and TOXIFY identified a total of 559 toxins in the studied

262 genomes (Supplementary Table 8), included as part of 189 orthogroups. The orthogroups with most

263 genes are displayed on Figure 3 and include OG0000175 (135 genes, Astacin-like metalloproteases as

264 determined by NCBI-blast), OG0000314 (105 genes, Neprilysins or endothelin-converting proteins),

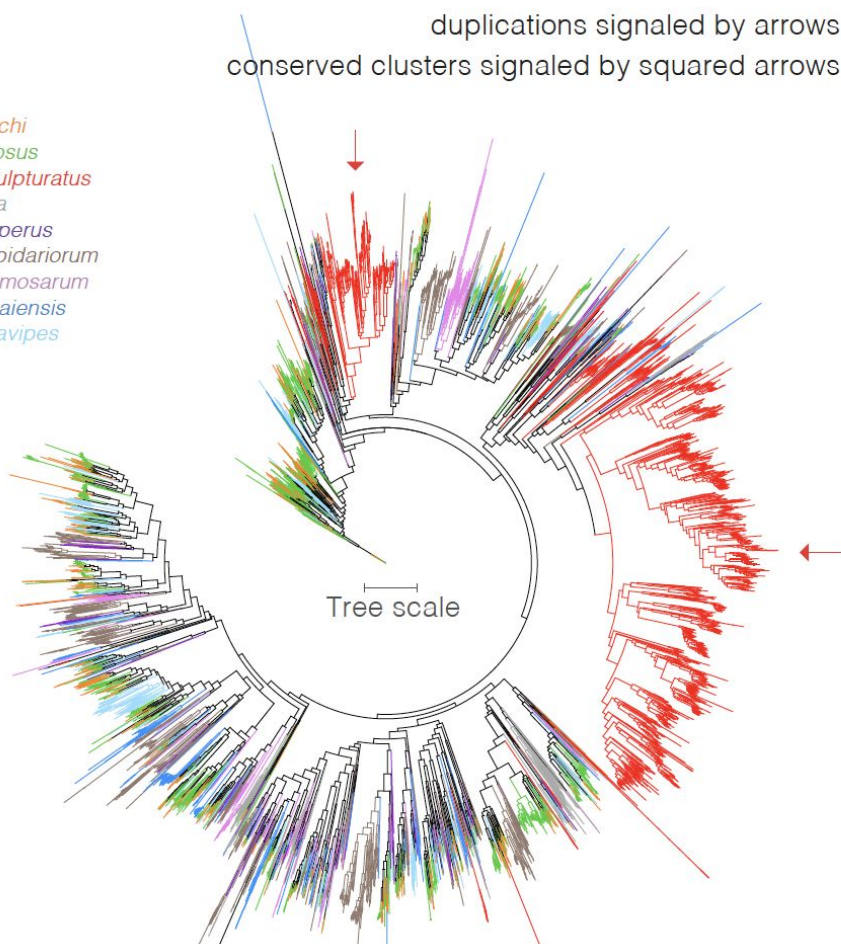
1  
2  
3 265 OG0000346 (99 genes, uncharacterized proteins), OG0000432 (86 genes, Techylectin), OG0000639  
4 266 (68 genes, various toxin-types) , OG0000761 (61 genes, Zonadhesins, various various toxin-types),  
5 267 OG0000803 (59 genes, Astacin-like metalloproteases), OG0000916 (54 genes, Papilins, Kunitz-type  
6 268 serine protease inhibitor) OG0000930 (54 genes, Astacin-like metalloproteases), OG0001436 (41  
7 269 genes, uncharacterized proteins). The two most toxin-rich assemblies were the *A. bruennichi* and *P.*  
8 270 *tepidariorum* where 154 and 200 toxins were identified, respectively. The scorpion genome, *C.*  
9 271 *sculpturatus*, yielded 31 toxins, whereas *D. silvatica* and *L. hesperus* yielded 13 and 16 toxins,  
10 272 respectively (Supplementary Table 8).  
11 273 Phylogenetic analyses of the orthogroups show that most venom families were present before the split  
12 274 between scorpions and spiders (Figure 3). Different spider genomes include species-specific  
13 275 expansions (i.e. groups of 5 or more genes from a single genome that cluster as a monophyletic  
14 276 clade), and many of these have relatively large branch lengths. Specifically, we find evidence for  
15 277 various expansions in *P. tepidariorum* (4 expansions, one with 7 genes, another with 12, one with 7  
16 278 and one with 9 genes), one expansion in *A. ventricosus* (one expansion with 11 closely related genes),  
17 279 one in *D. silvatica* (one expansion in 6 genes) and one in *C. sculpturatus* (5 genes expanded) in  
18 280 OG0000175 (Figure 3). In OG0000314, we found an expansion private to the three Araneidae  
19 281 genomes, including *A. bruennichi*, *A. ventricosus* and *T. clavipes*, various expansions exclusive to the  
20 282 *A. ventricosus* genome, and one expansion specific to the scorpion genome (9 genes). In OG000346,  
21 283 we found various expansions on the *S. mimosarum* (9 genes), *P. tepidariorum* (5 genes), *A.*  
22 284 *ventricosus* (8 genes) genomes. In OG000432 we found genome-specific expansions in *D. silvatica* (8  
23 285 genes; Figure 3). In OG000639, we found an expansion in *C. sculpturatus* (5 genes), and in  
24 286 OG000803 there are two 5-gene expansions, one in *C. sculpturatus*, another in *A. ventricosus*.  
25 287 OG000930 is only present in *T. kauaiensis* (1 expansion with 20 genes), *A. ventricosus*, *A.*  
26 288 *bruennichi*, *T. clavipes* and *P. tepidariorum*. OG0001436 is expanded in *C. sculpturatus* (5 genes).  
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

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

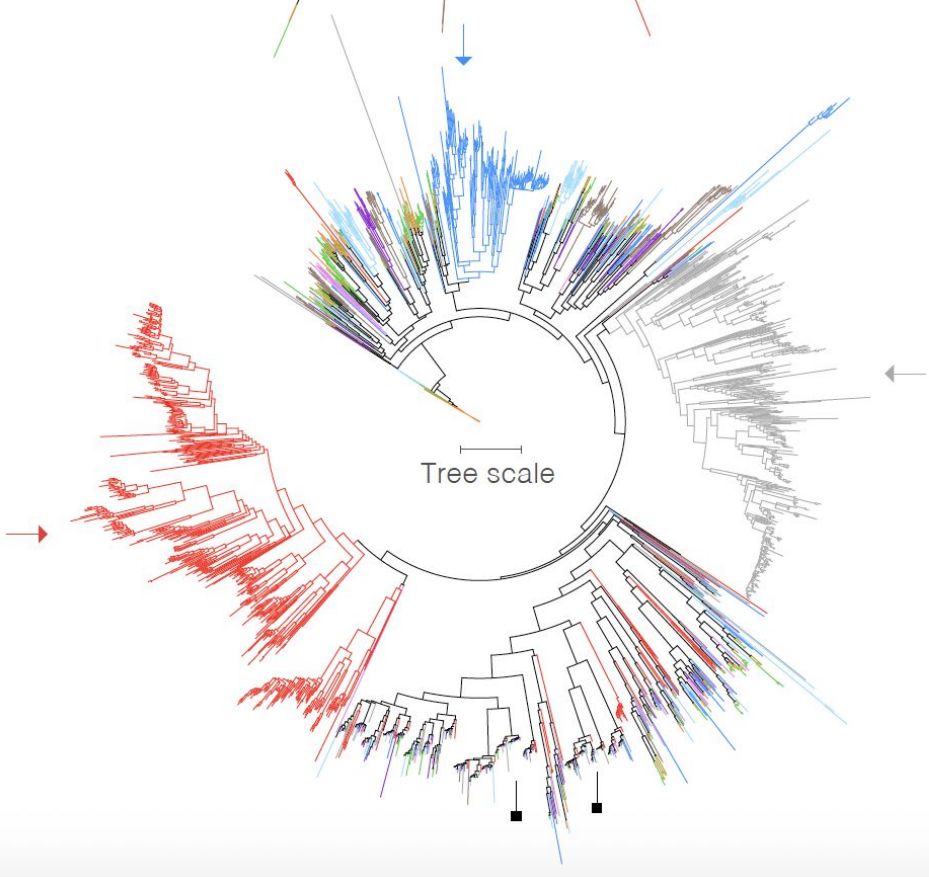
# A

- Argiope bruennichi*
- Araneus ventricosus*
- Centruroides sculpturatus*
- Dysdera silvatica*
- Latrodectus hesperus*
- Parasteatoda tepidariorum*
- Stegodyphus mimosarum*
- Tetragnatha kawaiensis*
- Trichonephila clavipes*

duplications signaled by arrows →  
 conserved clusters signaled by squared arrows ■

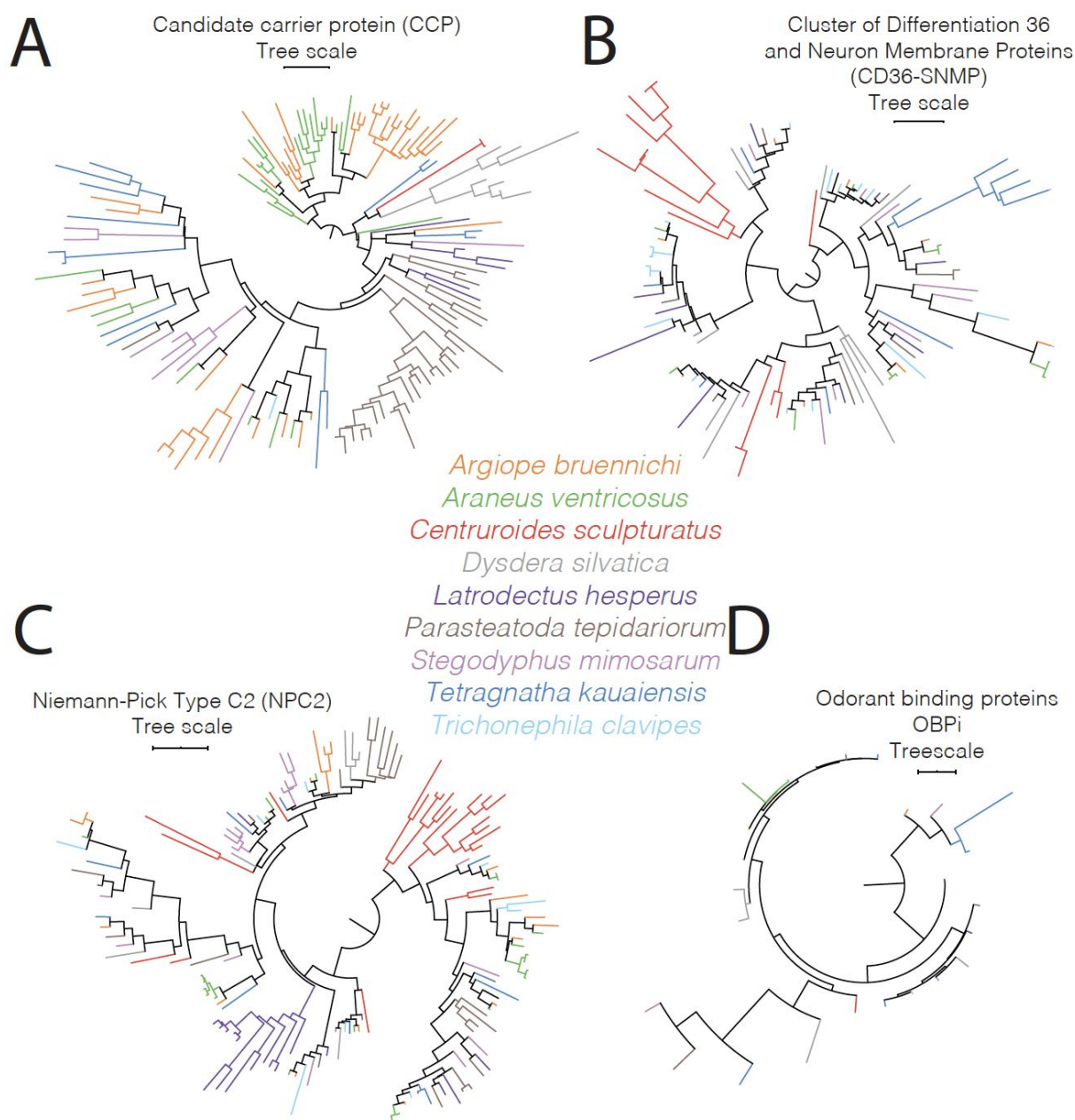


# B



Downloaded from https://academic.oup.com/gbe/advance-article/doi/10.1093/gbe/evab262/6443144 by Stanford Medical Center user on 02 December 2021

293 **Figure 4. Gustatory and Ionotropic reception phylogenies** a) Gustatory receptor phylogeny. The  
 294 phylogeny has 5,595 genes and includes every GR identified in the assemblies herein studied. b)  
 295 Ionotropic receptor phylogeny. The phylogeny has 1,932 genes and includes every IR identified in the  
 296 assemblies herein studied. Arrows indicate major duplications private to specific genomes, whereas  
 297 squared arrows highlight potentially conserved IR genes (small branch length and small duplicates).



303

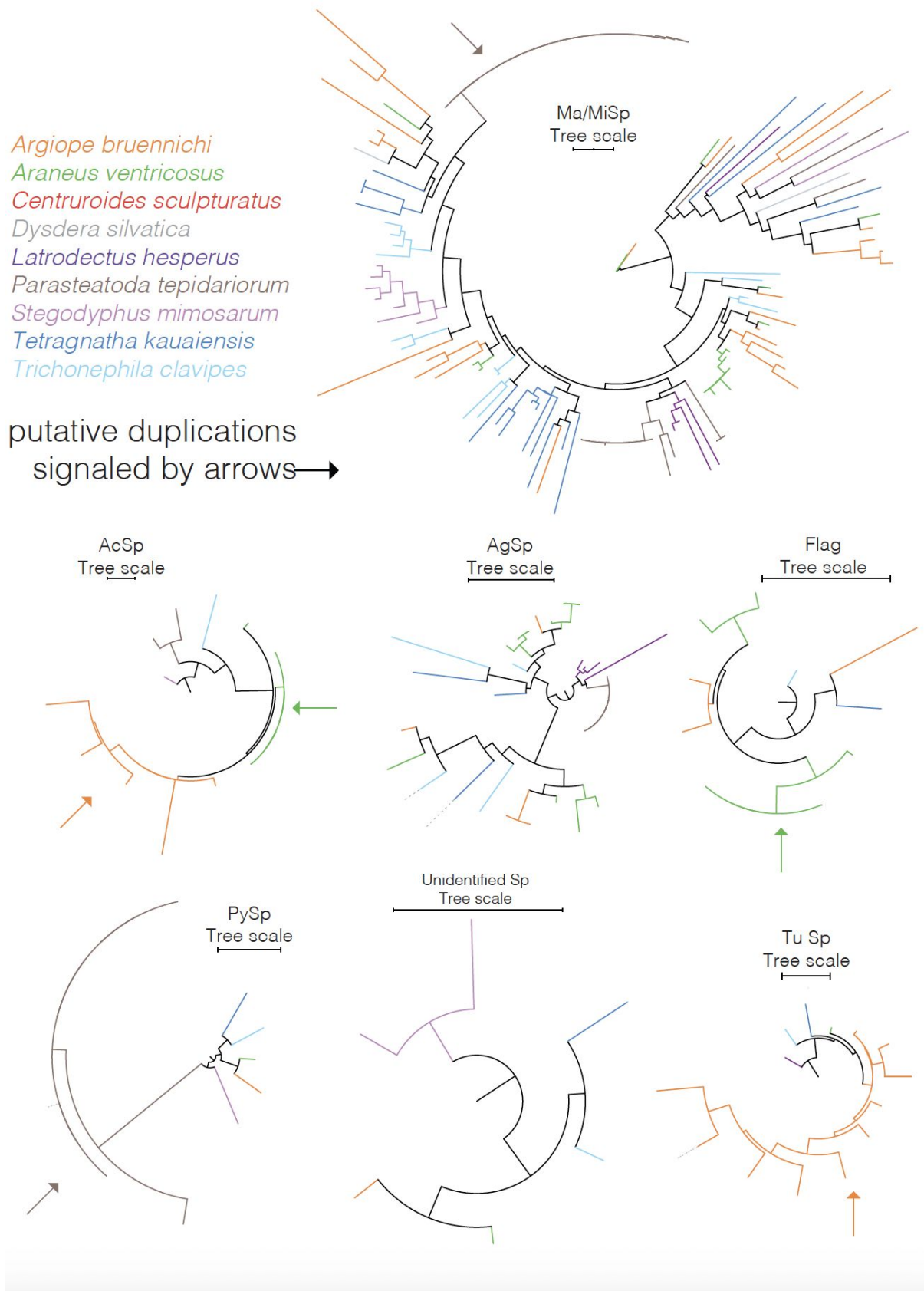
1  
2  
3 304 **Figure 5.** Phylogeny of other chemosensory genes. a) candidate carrier protein (CCP) phylogeny; b)  
4 305 cluster of differentiation 36 and neuron membrane proteins (CD36-SNMP) phylogeny; c) Niemann-  
5 306 Pick type C2 (NPC2), phylogeny; d) Odorant binding proteins (OBP-like) phylogeny.

6  
7  
8 307

9 308 **Chemosensory gene-family variation**

10  
11 309 We identified a total of 5,595 candidate gustatory receptors (GRs), 1,934 candidate ionotropic  
12 310 receptors (IRs), 25 candidate Odorant binding proteins (OBP-like), 147 candidate Niemann-Pick type  
13 311 C2 (NPC2), 137 candidate carrier protein (CCP), and 998 candidate cluster of differentiation 36 and  
14 312 neuron membrane proteins (CD36-SNMP; Supplementary Table 9; Figures 4-5). GRs exhibited a  
15 313 large interspecific variation (Figure 4), ranging between 1,436 GRs in *A. ventricosus* and 84 in *L.*  
16 314 *hesperus*. *C. sculpturatus*, the outgroup, had 1,648 GRs (Supplementary Table 9). The *D. silvatica*  
17 315 genome has the most IR/iGluR genes with 443 genes (Supplementary Table 9; Figure 4). We detected  
18 316 a total of 25 OBP-like genes, with 5 being present in *T. kauaiensis*, 4 in *D. silvatica* and in *S.*  
19 317 *mimosarum*, 3 in *P. tepidariorum* and all remaining genomes having only 1 or 2 OBP-like genes  
20 318 (Supplementary Table 9; Figure 5). From the 147 identified NPC2, *D. silvatica* had the least NPC2-  
21 319 genes (7 genes) and *A. ventricosus* the most (23). *A. bruennichi* had the most CCP, with 41 genes,  
22 320 while *C. sculpturatus* and *T. clavipes* had only 1 CCP (Supplementary Table 9; Figure 5). Finally, we  
23 321 identified at least 8 and at most 16 CD36-SNMP genes. *T. clavipes*, and *C. sculpturatus* had the most  
24 322 CD36-SNMP genes with 16 and 14, respectively, while *P. tepidariorum* and *A. bruennichi* had the  
25 323 least with 8 (Supplementary Table 9).

26  
27 324 Analysis of phylogenetic patterns suggests that the chemosensory portfolio is driven by a  
28 325 highly dynamic diversification process. For instance, within GRs there are two genome-specific  
29 326 expansions of genes in the scorpion, one including 1,237 genes and another 235 genes (Figure 4). A  
30 327 similar pattern is observed in the IRs where we find two genome-specific expansions private to the  
31 328 scorpion genome (88 genes, and 382 genes; Figure 4), a large genome-specific gene group with 392  
32 329 genes in *D. silvatica*, and another in the *Tetragnatha* genome including 139 genes. In CCPs, we  
33 330 found expansions in *A. bruennichi* (5 genes and 13 genes), *P. tepidariorum* (21 genes), *A. ventricosus*  
34 331 (8 genes) and *D. silvatica* (6 genes; Figure 5A). In CD36-SNMP we found expansions in the scorpion,  
35 332 (9 genes), and in *T. kauaiensis* (5 genes; Figure 5B). In NPC2, we found expansions in *L. hesperus*  
36 333 (14 genes), *P. tepidariorum* (6 genes), and *C. sculpturatus* (14 genes; Figure 5C), while in CD36-  
37 334 SNMP (Figure 5D) we found expansions in the *T. kauaiensis* (5 genes), and *C. sculpturatus* (9 genes)  
38 335 genomes.  
39 336



337  
338 **Figure 6.** Silk genes (spidroins) phylogeny These including Major and Minor Ampullate spidroins  
339 (Ma/MiSp), Aciniform spidroins (AcSp), Aggregate spidroins (AgSp), Flagelliform spidroins (Flag),



340 Pyriform spidroins (PySp), an unidentified spidroins group present in the *T. clavipes* genome and  
341 the Tubuliform spidroins (TuSp).

342

343

### 344 **Silk gene-family**

345 We identified a total of 24 putative spidroins in the genome of *T. kauaiensis* (Supplementary Table  
346 9). After querying these to the NCBI protein database, we identified 1 Flagelliform spidroin (Flag), 4  
347 Aggregate spidroins (AgSp), 8 Major Ampullate spidroins (MaSp), 3 Minor Ampullate spidroins  
348 (MiSp), 1 Tubuliform spidroins (TuSp), 1 Pyriform spidroin (PySp) and 1 Aciniform spidroin  
349 (AcSp). There was one spidroin for which NCBI did not yield any results, and 4 where the database  
350 retrieved more than a single gland as a top-hit (Supplementary Table 9). Alignments are provided in  
351 the supplementary.

352 Phylogenetic patterns of spidroin shows several genome-specific expansions of the Ma/Mi spidroins,  
353 including two separate expansions in the *P. tepidariorum* genome (25 genes and 10 genes;  
354 Supplementary Table 10; Figure 6), a single expansion in *S. mimosarum* including 7 genes, another in  
355 *A. ventricosus* including 8 genes, and another in *T. kauaiensis* including 7 genes. In the remaining  
356 spidroins, we find genome-specific expansions in AgSp and PySp in *P. tepidariorum*, with 9 and 6  
357 genes, respectively. In AcSp there are two smaller lineage-specific clades in *A. bruennichi* and *A.*  
358 *ventricosus*. There is a genome-specific expansion in *A. bruennichi* for the TuSp gland, with 7 genes  
359 (Supplementary Table 10; Figure 6).

360

### 361 **Discussion**

362 In this study, we report the sequence assembly of the *Tetragnatha kauaiensis* genome, and  
363 explore genome evolution across the available spider assemblies. To do so, we controlled for the  
364 quality of the assemblies, by focusing on contiguity and completeness (i.e. how complete a genome is  
365 from a gene content perspective based on the presence of universal single copy genes), finding that  
366 many of these assemblies are highly fragmented and incomplete. We find a wide variation in gene  
367 content, repeat content, and genome size in the surveyed spider genomes, which indicates a highly  
368 dynamic pattern of genome evolution. While the low quality of some assemblies did not hamper  
369 comparative analyses of the surveyed spider genomes, results should be interpreted with caution. By  
370 surveying all repeats and transposable elements (hereafter ‘the repeatome’) and studying Hox gene  
371 duplications, we find that the observed genome size differences are likely driven by the expansion of  
372 the repeatome. We also find significant gene-family expansions associated with sensory perception of  
373 taste, immunity and metabolism, which may underlie the diverse biology of spiders. We confirm  
374 previous work showing that venoms and chemosensory genes are present in high numbers across the  
375 assemblies, and discuss the role of a putative ancient whole genome duplication in generating the  
376 diversity we observe in spiders.

1  
2  
3 3774 378 **Repeat content underlie genome size variation in spiders**

5 379 Previous evidence from flow cytometry, Feulgen image analysis densitometry, and genome  
6 380 assembly sizes have found wide variation in genome size in spiders (Gregory and Shorthouse 2003;  
7 381 Sanggaard et al. 2014; Král et al. 2019). For instance, Gregory and Shorthouse (2003) assembled a  
8 382 large dataset comprising 115 species from 19 different families of spiders, finding that spider  
9 383 genomes vary between 5.73 - 0.79 C (~7 Gb for the jumping spider *Habronattus borealis* – ~724 Mb  
10 384 for the long-jawed orbweaver *Tetragnatha elongata*). They also reported a wide variation within  
11 385 relatively closely related species. For instance, genome size in the Salticidae family ranged between  
12 386 1.73 – 5.73 C (between *Habronattus borealis* and the peppered jumping spider *Pelegrina galathea*).  
13 387 Our results are in line with this evidence, since we found variation in genome size among spider  
14 388 assemblies (in our dataset the largest genome was *A. ventricosus* with 3.6 Gb, and the smallest was *T.*  
15 389 *kauaiensis* with 1.08 Gb). We also report variation between relatively closely related species (i.e.  
16 390 within the Araneidae family, where we included three assemblies, genome sizes ranged between 3.6  
17 391 Gb and 1.7 Gb). Similar to previous reports, we do not find a clear phylogenetic pattern of genome  
18 392 size variation across the spider tree of life (Gregory and Shorthouse 2003).

19 393 Genome size may increase through whole genome duplication, where the whole genome  
20 394 doubles itself, or through small scale duplication of genetic elements which may include duplication  
21 395 of genes or transposable elements. Recent evidence, using flow cytometry, has revealed a whole  
22 396 genome duplication in caponiid spiders (Král et al. 2019), which indicates the potential of further  
23 397 whole genome duplications in spiders, other than the duplication ~450 million years ago (Schwager et  
24 398 al. 2007; Schwager et al. 2017). While we have no caponiids in our dataset, we found no evidence of  
25 399 recent whole genome duplication specific to spiders on the analyzed assemblies. This evidence comes  
26 400 from several sources. First, there is a low % of double copy BUSCO genes – a set of highly curated  
27 401 genes, single copy genes. The scorpion assembly has a duplicate BUSCO score of 26 %, whereas  
28 402 spider genomes range between 26% - 0.8 %, in *P. tepidariorum* and *L. hesperus* (note that *L. hesperus*  
29 403 assembly has many missing BUSCOs, which is indicative of a poor assembly quality). Second,  
30 404 analysis of Hox genes shows that these genes are mostly present in 2 copies, with a single exception  
31 405 of four Hox4 in *A. ventricosus*. The four copies of Hox4 in *A. ventricosus* could be an artefact due to  
32 406 the similarity between Hox genes, and we were not able to obtain candidates for Hox1 using the 95%  
33 407 cut-off threshold. The BUSCO-pattern together with that from the Hox genes are in line with the  
34 408 evidence for an ancestral whole genome duplication in Arachnoplumonata. Third, an important finding  
35 409 of our work is that variation in genome size of spiders is largely driven by the duplication of genetic  
36 410 elements, and specifically, the repeatome (transposable elements and repeats). Indeed, we find a  
37 411 R=0.95 correlation between the ‘length of the masked repeats’ and the ‘genome size’ – a strong  
38 412 indication of the role of the repeatome in underlying genome size changes (Figure 1). Expansions of  
39 413 the repeatome are generally constrained in animal lineages since bigger genomes translate to higher

1  
2  
3 414 cell-economy costs through the increase of cell size. In addition to this, proliferation of transposable  
4 415 elements may interfere with gene expression when these selfish elements jump in front of a gene  
5 416 promoter (Choi and Lee 2020). Considering the strikingly different representation of the repeatome  
6 417 that we find here, including the variation in transposable element accumulation through time, we  
7 418 speculate that transposable elements may have had a role in the regulation and variation of gene  
8 419 expression across spiders, likely underlying some of the observed morphological and physiological  
9 420 diversity.

10 421 By conducting a *de novo* annotation of repeats and using the same version and library of repeats for  
11 422 every genome, we guaranteed a standardization of the repeat identification, thereby removing  
12 423 potential biases due to the use of different databases and pipelines. Variation in some elements, both  
13 424 in terms of classes and extent along the genome, was substantial. For instance, Long Interspersed  
14 425 Nuclear Elements (LINEs) occupy less than >2% in every assembly, but occupy 10.3% of the *D.*  
15 426 *silvatica* assembly. This may suggest mechanisms to purge LINEs from some clades, or an expansion  
16 427 specific to *D. silvatica* (and possibly closely related species). Furthermore, DNA elements had a  
17 428 three-fold variation, ranging between 5.59 % (*T. kauaiensis*) and 18.82 % (*D. silvatica*). Despite the  
18 429 overall variation in numbers and accumulation of the repeatome through time, there was a clear  
19 430 dominance of DNA/TcMar and DNA/hAT elements (both DNA elements) across the assembly when  
20 431 considering the top three most represented categories (Figure 1B), suggesting these elements are the  
21 432 most prolific and present across spiders, and potentially scorpions (keep in mind we have single  
22 433 scorpion genome in our analyze using the same version and library of repeats for every genomes).  
23 434 Future studies on spider genome assemblies should put transposable element variation in the context  
24 435 of the spider phylogeny, and should benefit from an increased sampling of spider genomes. The  
25 436 differential presence of repeats and transposable elements may indicate that mechanisms to eliminate  
26 437 these elements such as nonhomologous end joining, or illegitimate recombination may be active in  
27 438 these genomes (Choi et al. 2020). A phylogenetic framework together with ancestral character  
28 439 reconstructions, focusing on transposable element data, will certainly elucidate the patterns of  
29 440 activation and deactivation of certain transposable element classes, and how changes in transposable  
30 441 element proliferation may be linked to particular events in the evolution of spiders. For instance, a  
31 442 caponiid genome, where a more recent genome duplication was detected (Kral et al. 2019), may help  
32 443 understand the impacts of whole genome duplication and transposable element proliferation in  
33 444 spiders. This would allow testing the ‘genomic shock’ hypothesis after genome duplication in spiders.  
34 445 Finally, the variation in the repeatome is in line with those of the remaining arthropods, where  
35 446 variation in transposable elements load was deemed as an important predictor for genome size (Wu  
36 447 and Lu 2019; Gilbert et al. 2021).

448

449 **Gene duplicates**

1  
2  
3 450 Observed patterns in the explored gene families, namely venoms and chemosensory, suggest a  
4 451 central role in the evolution of spiders (Figure 3-5). The presence of most gene families in the  
5 452 scorpion genome and in spider genomes suggests an ancestral status (Vizueta, Escuer, et al. 2020),  
6 453 while variation in gene numbers and their branch lengths along the phylogeny is an indication of  
7 454 divergence, and thereby indirect evidence of the acquisition of novel gene functions (i.e.  
8 455 neofunctionalization). Gene duplicates generally experience relaxation of purifying selection or gene  
9 456 dose compensation and, if one of the copies does not get sub- or neofunctionalized through time, it  
10 457 will be lost. Indeed, we manually curated chemosensory genes, finding a low ratio of pseudogenes  
11 458 (Supplementary Table 9). There are large genome-specific duplications detected in *C. sculpturatus*, *T.*  
12 459 *kauaiensis* and *D. sylvatica* in the two largest chemosensory families (Figure 4 A, B). This is an  
13 460 indicator of the importance of gustatory (GRs) and ionotropic receptors (IRs) in *T. kauaiensis* and *D.*  
14 461 *sylvatica*, and we speculate it may be associated with the colonization of islands (*T. kauaiensis* is part  
15 462 of a Hawaiian radiation of spiders, and *D. sylvatica* is part of a Macaronesian radiation) where  
16 463 environmental conditions can be very different (disharmonic biotas, open ecological niches) (Vizueta  
17 464 et al. 2019). We note that, unfortunately, the taxonomic range (i.e. 1 single genome for Tetragnathidae  
18 465 and 1 single for Dysderidae) does not allow dissecting whether these changes are shared by other  
19 466 members of the families, whether they are private to the species in question (*D. sylvatica*, *T.*  
20 467 *kauaiensis*) or even to the adaptive radiation (in Hawai'i and Macaronesia). Similarly, since we only  
21 468 included a single scorpion assembly, we cannot comment on whether the expansions observed in *C.*  
22 469 *sculpturatus* are specific to all scorpions, or just the *C. sculpturatus* genome.

23 470 Despite the aforementioned evidence, not every gene family is present in very high numbers.  
24 471 For example, we only detected 25 OBP-like genes in all genomes, and the small number of genes  
25 472 together with the short branch lengths confirms that the OBP-like are a relatively conserved family of  
26 473 genes in arachnids (Vizueta et al. 2017). In addition to the OBP-like, we also find few silk genes, with  
27 474 very short branch lengths (notice *P. tepidariorum* in PySp and Ma/MiSp, *A. ventricosus* in Flag and  
28 475 AcSp), which may be indicative of very recent duplications in silk genes (Garb et al. 2007; Clarke et  
29 476 al. 2014; Clarke et al. 2015). These results are in line with those of Clarke et al. (2015) who used  
30 477 transcriptomics to suggest that a large-scale duplication occurred early in the divergence of spiders,  
31 478 and that multiple independent duplication events in silk genes have likely taken place afterwards. Our  
32 479 results, however, have to be interpreted with caution since silk genes are composed of sequences (of  
33 480 often hundreds) of repeated aminoacids (Clarke et al. 2015), being therefore hard to reconstruct in  
34 481 entirety in the gene annotation process, and being typically fragmented onto separate fragments.  
35 482 Considering the fragmentation of most assemblies, it is possible that some duplicates consist of gene  
36 483 fragments.

37 484

38 485 **Significant expansion of metabolism, immunity and sensory perception gene families**

1  
2  
3 486 Using a statistical approach to detect expansion of gene families, we find that most  
4  
5 487 expansions are in terminal branches. As a direct comparison, recent analyses on 76 insect assemblies  
6  
7 488 were able to identify 147 expanded gene families, comprising 9,601 genes, in the branch  
8  
9 489 corresponding to insects ('the Last-Insect-Common-Ancestor'; Thomas et al. 2020), thereby  
10  
11 490 providing evidence for 'ancient expansions' particular to insects. Thomas et al (2020), however,  
12  
13 491 included 10 times more genomes than we did, and some of the spider genomes in our dataset lack  
14  
15 492 substantial data, as indicated by the BUSCO scores (Supplementary Table 6). Thus, it is possible that  
16  
17 493 spiders have their own set of 'ancient expansions', which we were not able to detect due to the  
18  
19 494 limitations of our dataset. It is also possible that the inclusion of fragmented assemblies (*D. silvatica*  
20  
21 495 and *L. hesperus*) leads to an inflation of expanded gene families on closely related assemblies (e.g. *P.*  
22  
23 496 *tepidariorum*). We expect that the addition of more highly completed spider genomes will help to  
24  
25 497 further our understanding of the evolutionary history of gene families in spiders.

26  
27 498 Despite the challenges in the dataset, we find notable evidence for various gene families  
28  
29 499 expansions in spiders. Specifically, using gene ontology annotations (GO) we find that gene families  
30  
31 500 associated with various metabolic functions, sensory perception of taste, and immune functions are  
32  
33 501 expanded. This pattern is similar to the pattern found in arthropods which includes expansions of  
34  
35 502 metabolic genes (Thomas et al. 2020). These independent pieces of evidence suggest that gene  
36  
37 503 duplications associated with metabolism, immunity and sensory functions may have been  
38  
39 504 instrumental to the evolution of arthropods in general, but also spiders specifically. We speculate that  
40  
41 505 these expansions may contribute to the success, in terms of number of species and adaptation to  
42  
43 506 different environments in spiders. As chromosome resolved assemblies become cheaper and  
44  
45 507 technically less challenging, revising the role of gene expansions and gene contractions will certainly  
46  
47 508 yield important insights towards the understanding of genome evolution of spiders.

## 48 509 **Conclusion**

49 510 We have sequenced the *Tetragnatha kauaiensis* genome, and explored patterns of genome  
50  
51 511 evolution across various genome assemblies. Comparative genomics analyses including *T. kauaiensis*,  
52  
53 512 1 scorpion (outgroup), and 7 additional spiders assemblies suggest that variation of transposable  
54  
55 513 elements and repeat content are associated with the wide variation of spider genome sizes. We also  
56  
57 514 found many duplications in chemosensory and venom genes, consistent with the evidence that the  
58  
59 515 evolution of toxins and the ability to perceive the environment are ancestral attributes of spider  
60  
61 516 evolution. Our results suggest that the evolutionary history of spiders is characterized by gene-family  
62  
63 517 expansions associated with sensory perception of taste, metabolism and immune responses, and by  
64  
65 518 multiple gene duplication events. While we uncovered interesting patterns of genome evolution, we  
66  
67 519 acknowledge the limitations of this work due to the lack of high-quality genomes. We hope that,  
68  
69 520 however, this work catalyzes enthusiasm in the spider research community to produce and analyse  
70  
71 521 more high-quality genomes.  
72  
73 522

## 523 **Methods**

### 524 ***Tetragnatha kauaiensis* - Genome sequencing, assembly, annotation and quality verification**

525 We sequenced the genome of a single individual of *T. kauaiensis* using a paired-end and a  
526 non-size selected mate-pair library on a lane of Illumina HiSeq4000 (individual ID AJR402, collected  
527 31/May/2013 by AJ Rominger in Kaua'i, at 22.1412, -159.6206). Using these libraries we built a base  
528 assembly using ALLPATHS-LG with default parameters in addition to 'HALOIDIFY = True'  
529 (Gnerre et al. 2011). We then sequenced an additional individual using the Dovetail Chicago method  
530 (AJR443, collected 03/June/2013 by AJ Rominger in Kaua'i, at 22.1469, -159.6638), which was used  
531 to scaffold the initial assembly using the HiRise software (Koch 2016; Putnam et al. 2016).

532 The quality of the assembly was first assessed using BUSCO v3.0.2 arthropoda db v9 (Simão  
533 et al. 2015), which searches for highly conserved genes in the assembly. Then we used the  
534 Assemblathon 2 script (<https://github.com/ucdavis-bioinformatics/assemblathon2-analysis>) (Bradnam  
535 et al. 2013), which assesses scaffold and contig statistics, to evaluate the quality of the assembly.  
536 Annotation of repeats was carried out by identifying and building a database of repeats along the  
537 genome using RepeatModeler followed by masking them using RepeatMasker (Tarailo-Graovac and  
538 Chen 2009). We explored the draft assembly for contaminants, including gut-microbiota and wet-lab  
539 contaminants using Blobtools (Koutsovoulos et al., 2016; Laetsch et al., 2017; Dominik R Laetsch  
540 and Blaxter, 2017)(Supplementary Figure 1).

541 To determine protein-coding genes and their locations along the genome, we used  
542 BRAKERv1 (Hoff et al. 2019). We used whole-body *T. kauaiensis* transcriptome reads previously  
543 generated by (Yim et al. 2014) (SRR1313313, SRR1427109). Raw transcriptomic reads were cleaned  
544 using Trimmomatic (Bolger et al. 2014) and aligned to the generated genome using STAR (Dobin et  
545 al. 2013). The resulting binary alignment map (BAM) file was provided to BRAKERv1 as RNA-  
546 based evidence. The final annotation was assessed by BUSCOv4.0.1 (Seppey et al. 2019), using the  
547 Arthropoda10 (1,013 genes) and Arachnida10 (2,943 genes) gene sets.

548

### 549 **Genomes used for comparative genomics**

550 We searched the I5K and NCBI databases and the literature for published and available spider  
551 genomes (data consulted on 23rd October 2019). In total, we downloaded nine spider genomes  
552 (Supplementary Table 1), their general feature format (gff3), and predicted protein files (faa;  
553 Supplementary Table 1). From the available genomes, we selected those with a contig-N50 above  
554 8,000 bp in order to avoid genomes that were highly fragmented. This included the genomes of  
555 *Stegodyphus mimosarum* (Sanggaard et al. 2014), *Latrodectus hesperus* (BCM-HGSC website),  
556 *Parasteatoda tepidariorum* (Gendreau et al. 2017), *Trichonephila clavipes* (Babb et al. 2017),  
557 *Dysdera silvatica* (Sánchez-Herrero et al. 2019), *Araneus ventricosus* (Kono et al. 2019) and *Argiope*  
558 *bruennichi* (Sheffer et al. 2021). Additionally, we downloaded the genome of the bark scorpion  
559 *Centruroides sculpturatus* (Schwager et al. 2017) as an outgroup.

560

**561 Characterization of spider genomes**

562 We characterized spider genomes based on the (i) continuity and completeness of the  
563 assemblies, (ii) assembly size, (iii) repeat-content, and (iv) broad genomic features. Specifically, (i)  
564 the continuity of each genome serves as a proxy of the overall quality of an assembly, and it affects  
565 the detection of genes, repeat sequences and transposable elements (Peona et al. 2018). We  
566 characterized the contiguity of the assemblies using the Assemblathon 2 script, as described above for  
567 *T. kauaiensis*, retrieving contig-N50, scaffold-N50, total number of contigs, total number of scaffolds,  
568 maximum scaffold size, assembly size and GC content. (ii) The ‘completeness’ of the assemblies, is  
569 generally defined as an overview of the genes which may be missing, fragmented, duplicated or  
570 present in a single copy in an assembly. To assess the completeness of the genomes, we used BUSCO  
571 v4.0.1 as outlined above for *Tetragnatha kauaiensis* (the Arthropoda10 set including 1013 genes; and  
572 the Arachnida10 set including 2,943 genes). (iii) To assess repeat content, we used Repeat-Modeler  
573 v2.0.1 and Repeat-Masker-v4.1.0. Repeat content in the genome includes simple repeats (typically 1-5  
574 base pairs, e.g. AAA, TTTT), tandem repeats (100-200 base pairs), segmental duplications (10,000 -  
575 300,000 base pairs), and interspersed repeats (SINES, which are non-functional copies of RNA genes  
576 that were reintegrated into the genome; DNA transposons; LINES, which are non-retrovirus  
577 retrotransposons). We ran RepeatModeler and RepeatMasker for each genome to screen and annotate  
578 DNA sequences *de novo*, thereby annotating and masking repeats. We retrieved repeat-statistics  
579 including % of the genome covered by different repeats and transposable element landscape plots.  
580 Finally, (iv) we assessed broad genomic features including, among others, the number of genes,  
581 coding sequences, introns, gene length using Another Gff Analysis Toolkit v0.4.0 (AGAT available at  
582 <https://github.com/NBISweden/AGAT/>; `agat_sp_functional_statistics.pl`, and `agat_sp_statistics.pl`).  
583 The association between total genome size, and % of masked sequences and total length of masked  
584 genome was assessed with a correlation using the `cor()` function in R.

585

**586 Spider genome evolution**

587 Previous work suggests that the whole genome duplication in the common ancestor of  
588 scorpions and spiders can be linked to the diversification of spiders (Schwager et al. 2007; Schwager  
589 et al. 2017). To better understand the presence of whole genome duplication in the studied lineages,  
590 we used two complementary approaches. We first analyzed repeat content variation in the available  
591 spider genomes (as described above), since differences in repeat content may translate to differences  
592 in genome size. Second, we downloaded the Hox genes 1-5 from the *P. tepidariorum* genome, and  
593 searched for these in the remaining spider genomes using BLAST (Altschul et al. 1990). Hox gene-  
594 copies are prime candidates for detecting whole genome duplications since they are functionally  
595 constrained (Leite et al. 2018). For example, a 1:4 ortholog ratio is maintained between the  
596 *Drosophila melanogaster* genome and vertebrate genomes, indicating the two whole genome

1  
2  
3 597 duplications which occurred in the lineage of modern vertebrates (Hakes et al. 2007; Schwager et al.  
4 598 2017).

599

### 600 **Spider gene-family evolution**

601 Another component of genome evolution is gene-family expansion and reduction, or the gain  
602 and loss of gene-copies. Focusing on the predicted-proteins resulting from the annotations of the  
603 spider genomes, we first cleaned and filtered sequences using Kinfin's  
604 filter\_fastas\_before\_clustering.py (Laetsch and Blaxter 2017) removing sequences shorter than 30  
605 amino acids. We then removed all isoforms of a given gene, keeping only the longest isoform using  
606 in-house scripts. For this analysis, we removed the genome of *A. ventricosus* since it has twice the  
607 number of genes compared to the other spider genomes, and this biases the analysis. Cleaned and  
608 isoform-free prediction-proteins were then analyzed using Computational Analysis of Family  
609 Evolution (CAFE v 4.2.1) (De Bie et al. 2006). Briefly, we first determined gene-similarity (based on  
610 BLAST- e-values) in the dataset using an all-by-all blast approach. We then applied a Markov Cluster  
611 algorithm (MCL; mexload, mcl mcxdump) (Enright et al. 2002), and parsed the output using the  
612 mcl2rawcafe.py script. These clusters (gene-families) are then integrated in a phylogenetic-backbone,  
613 which was retrieved from OrthoFinder's single-copy orthologs (Emms and Kelly 2015). This tree was  
614 then converted to an ultrametric format with r8s (Sanderson 2003), using the divergence time of 175  
615 million years between Tetragnathidae (*T. kauaiensis*) and Araneidae (*A. bruennichi*) as a calibration  
616 point (Fernández et al. 2018). We used Dendroscope's Graphical User Interface (GUI) to visualize  
617 trees and remove bootstrap support (Huson and Scornavacca 2012). Using the main pipeline of  
618 CAFE, we estimated the birth-death parameter lambda ( $\lambda = 0.0021$ ) for the dataset and obtained  
619 information on gene-family under significant evolution.

620 Genes belonging to gene-families that have undergone significant changes, that is, fast  
621 evolving families, were annotated using Gene Ontology terms (GO:terms) using the command-line  
622 version of Interproscan v5.34-73.0 (Ashburner et al. 2000). GO term annotations for genes belonging  
623 to expanded or reduced gene families were summarized and plotted as a treemap using R (Team and  
624 Others 2013) with REVIGO's treemap script (Supek et al. 2011) .

625

### 626 **Silk, chemosensory and venom gene variation**

627 To investigate venom gene evolution, we downloaded all toxin sequences available in the  
628 Arachnoserver v3.0 (Pineda et al. 2018), and used these as a database to query proteins from the  
629 spider and scorpion genomes with BLAST. Hits with e-values below  $1e-10$  were considered as  
630 candidate venom-genes. However, since venom proteins are potentially highly divergent and typically  
631 short, BLAST searches may result in a high proportion of false positives. To address this issue, we ran  
632 TOXIFY on the candidates, a pipeline specifically designed to identify toxins using deep learning  
633 algorithms (Cole and Brewer 2019). TOXIFY generates a prediction score between 0 and 1 where the



1  
2  
3 634 higher the score, the more likely a molecule is to be a venom, and we selected values above 0.75 as a  
4 635 criterion here. After TOXIFY, we kept a list of 589 putative venom genes across the assemblies. We  
5 636 then used OrthoFinder, obtaining an orthogroup-assignment for each of these 589 venom genes,  
6 637 finding that they group in 189 orthogroups. From these 189 groups, we selected the 10 biggest (in  
7 638 terms of gene number), identified the toxin-group using NCBI nr protein database, and aligned the  
8 639 genes within orthogroups using mafft v7.455 (Kato and Standley 2013). These alignments were then  
9 640 used to obtain a maximum likelihood (ML) phylogenetic tree with bootstrap estimate (automatic  
10 641 determination of the substitution model) using IQ-Tree v1.6.12 (Nguyen et al. 2015; Chernomor et al.  
11 642 2016; Kalyaanamoorthy et al. 2017; Hoang et al. 2018). The resulting phylogeny was plotted,  
12 643 formatted, coloured and labelled using the iTOL web server (Letunic and Bork 2019).

13 644         Considering the recent evidence on the wide variation in chemosensory gene-family size in  
14 645 Chelicerates (Vizueta et al. 2017; Vizueta et al. 2018), we searched the available genomes for  
15 646 Gustatory Receptors (GRs), Ionotropic Receptors (IRs), Niemann-Pick Type C2 (NPC2), Odorant  
16 647 binding proteins (OBP-like), Candidate carrier protein (CCP), Cluster of Differentiation 36 and  
17 648 Neuron Membrane Proteins (CD36-SNMP). To do so, we used BITACORA v1.2 (Vizueta, Escuer, et  
18 649 al. 2020; Vizueta, Sánchez-Gracia, et al. 2020), using its GeMoMa algorithm (Keilwagen et al. 2019),  
19 650 benefiting from a curated chemosensory database used in Vizueta et al (2018). To ensure the quality  
20 651 of the annotations, we performed a round of manual curation of the results, guaranteeing that (i) only  
21 652 a single isoform was selected and (ii) that putative annotation artefacts including small fragments,  
22 653 chimeric annotations or identical proteins by misassembly of duplicated contigs were removed.  
23 654 Finally, curated gene members were classified as pseudogenes (i.e. sequences with in-frame stop  
24 655 codons), partial or putatively complete functional proteins. The identified GRs, IRs, NPC2, OBP-like,  
25 656 CCP and CD36-SNMP were aligned using mafft, and a tree was generated and plotted using IQ-Tree  
26 657 and iTOL as described above.

27 658         We next identified spidroins (silk genes). To do so, we used a combination of BLAST  
28 659 searches using N-domains published with the *T. clavipes* genome, and the NCBI accession numbers  
29 660 for N-terminals and C-terminals from Vienneau-Hathaway et al. (2017). We extracted hits with an e-  
30 661 value below 1e-10 and candidate silk genes were then queried in NCBI nr database search (blastp) to  
31 662 classify the gland to which they belong based on NCBI's top hit. After labelling the gland, we did an  
32 663 orthogroup assignment using OrthoFinder as described above, and built a phylogeny for the silks in  
33 664 each gland, using the same method as described above for venom genes.

34 665

#### 35 666 **Data availability statement**

36 667 The raw data is available through ENA ( <https://www.ebi.ac.uk/ena/browser/home> ), ID:  
37 668 PRJEB48087. The assembly and annotation is available through DRYAD  
38 669 (<https://doi.org/10.5061/dryad.b2rbnzsg> ).  
39 670

40 670

**671 Acknowledgements**

672 JC is immensely grateful to Torsten H. Struck, for the freedom to pursue his interests, and for  
673 mentorship and stewardship (this is NHM genomics contribution X). JC is extremely grateful to Mark  
674 Blaxter for receiving him in his laboratory and opening the assembly, annotation and comparative  
675 genomics world – ‘lang may yer lum reek’. We thank Andy J. Rominger for the field collection of  
676 specimens. JC thanks Lewis Steven, Andrea Martínez Martínez, and Dom Laetsch for their time,  
677 patience and expertise on gene and repeat annotation. A note of appreciation to Samuel Abalde, for  
678 advice on venom identification and classification. JC is thankful to Nina Sokolov, Wagner Menezes,  
679 Lisa Carroll, Katherine Magoulick, Aahan Agrawaland, Leke Hutchins, and Nik Susič, Leif Egil Loe,  
680 for support, friendship and advice (‘em cada esquina, um amigo’). This paper was possible due to a  
681 Peder Sæther grant which funded JC to visit and stay with RG. A NORBIS travel internationalization  
682 grant guaranteed funding for JC to be trained by Mark Blaxter in Edinburgh. RF acknowledges  
683 support from the Ministerio de Economía y Competitividad and the Ministerio de Ciencia of Spain  
684 (RyC2017-22492 and PID2019-108824GA-I00). JR and JV are supported by the Ministerio de  
685 Economía y Competitividad and the Ministerio de Ciencia of Spain (CGL2016-75255 and PID2019-  
686 103947GB). We are indebted to three anonymous reviewers for their insightful and detailed  
687 comments which have made this work considerably more coherent and stronger.

688

**689 References**

690

691

692 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J.*  
693 *Mol. Biol.* 215:403–410.

694 Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Michael Cherry J, Davis AP, Dolinski K,  
695 Dwight SS, Eppig JT, et al. 2000. Gene Ontology: tool for the unification of biology. *Nat. Genet.*  
696 25:25–29.

697 Babb PL, Lahens NF, Correa-Garhwal SM, Nicholson DN, Kim EJ, Hogenesch JB, Kuntner M,  
698 Higgins L, Hayashi CY, Agnarsson I, et al. 2017. The *Nephila clavipes* genome highlights the  
699 diversity of spider silk genes and their complex expression. *Nat. Genet.* 49:895–903.

700 Binford GJ. 2001. Differences in venom composition between orb-weaving and wandering Hawaiian  
701 Tetragnatha (Araneae). *Biol. J. Linn. Soc. Lond.* [Internet]. Available from:  
702 <https://academic.oup.com/biolinnean/article-abstract/74/4/581/2639728>

703 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data.  
704 *Bioinformatics* 30:2114–2120.

705 Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman JA,  
706 Chapuis G, Chikhi R, et al. 2013. Assemblathon 2: evaluating de novo methods of genome  
707 assembly in three vertebrate species. *Gigascience* 2:10.

708 Brewer MS, Cotoras DD, Croucher PJP, Gillespie RG. 2014. New sequencing technologies, the

- 1  
2  
3 709 development of genomics tools, and their applications in evolutionary arachnology. *Arachnol.*  
4 710 *Mitt.* 42:1–15.  
5
- 6 711 Chernomor O, von Haeseler A, Minh BQ. 2016. Terrace Aware Data Structure for Phylogenomic  
7 712 Inference from Supermatrices. *Syst. Biol.* 65:997–1008.  
8
- 9 713 Choi I-Y, Kwon E-C, Kim N-S. 2020. The C- and G-value paradox with polyploidy, repeatomes,  
10 714 introns, phenomes and cell economy. *Genes & Genomics* [Internet] 42:699–714. Available from:  
11 715 <http://dx.doi.org/10.1007/s13258-020-00941-9>  
12
- 13 716 Choi JY, Lee YCG. 2020. Double-edged sword: The evolutionary consequences of the epigenetic  
14 717 silencing of transposable elements. *PLoS Genet.* 16:e1008872.  
15
- 16 718 Clarke TH, Garb JE, Hayashi CY, Arensburger P, Ayoub NA. 2015. Spider Transcriptomes Identify  
17 719 Ancient Large-Scale Gene Duplication Event Potentially Important in Silk Gland Evolution.  
18 720 *Genome Biol. Evol.* 7:1856–1870.  
19
- 20 721 Clarke TH, Garb JE, Hayashi CY, Haney RA, Lancaster AK, Corbett S, Ayoub NA. 2014. Multi-  
21 722 tissue transcriptomics of the black widow spider reveals expansions, co-options, and functional  
22 723 processes of the silk gland gene toolkit. *BMC Genomics* 15:365.  
23
- 24 724 Cole TJ, Brewer MS. 2019. TOXIFY: a deep learning approach to classify animal venom proteins.  
25 725 *PeerJ* 7:e7200.  
26
- 27 726 Cotoras DD, Brewer MS, Croucher PJP, Oxford GS, Lindberg DR, Gillespie RG. 2016. Convergent  
28 727 evolution in the colour polymorphism of Selkirkiella spiders (Theridiidae) from the South  
29 728 American temperate rainforest. *Biol. J. Linn. Soc. Lond.* 120:649–663.  
30
- 31 729 Croucher PJP, Brewer MS, Winchell CJ, Oxford GS, Gillespie RG. 2013. De novo characterization of  
32 730 the gene-rich transcriptomes of two color-polymorphic spiders, *Theridion grallator* and *T.*  
33 731 *californicum* (Araneae: Theridiidae), with special reference to pigment genes. *BMC Genomics*  
34 732 14:862.  
35
- 36 733 De Bie T, Cristianini N, Demuth JP, Hahn MW. 2006. CAFE: a computational tool for the study of  
37 734 gene family evolution. *Bioinformatics* 22:1269–1271.  
38
- 39 735 Dimitrov D, Lopardo L, Giribet G, Arnedo MA, Alvarez-Padilla F, Hormiga G. 2012. Tangled in a  
40 736 sparse spider web: single origin of orb weavers and their spinning work unravelled by denser  
41 737 taxonomic sampling. *Proc. Biol. Sci.* 279:1341–1350.  
42
- 43 738 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR.  
44 739 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29:15–21.  
45
- 46 740 Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome comparisons  
47 741 dramatically improves orthogroup inference accuracy. *Genome Biol.* 16:157.  
48
- 49 742 Enright AJ, Van Dongen S, Ouzounis CA. 2002. An efficient algorithm for large-scale detection of  
50 743 protein families. *Nucleic Acids Res.* 30:1575–1584.  
51
- 52 744 Escuer P, Pisarenco VA, Fernández-Ruiz AA, Vizueta J, Sánchez-Herrero JF, Arnedo MA,  
53 745 Sánchez-Gracia A, Rozas J. 2021. The chromosome-scale assembly of the Canary Islands  
54 746 endemic spider *Dysdera silvatica* (Arachnida, Araneae) sheds light on the origin and genome  
55 747 structure of chemoreceptor gene families in chelicerates. *Molecular Ecology Resources*  
56 748 [Internet]. Available from: <http://dx.doi.org/10.1111/1755-0998.13471>  
57
- 58 749 Fan Z, Yuan T, Liu P, Wang L-Y, Jin J-F, Zhang F, Zhang Z-S. 2021. A chromosome-level genome  
59 750 of the spider *Trichonephila antipodiana* reveals the genetic basis of its polyphagy and evidence  
60

- 1  
2  
3 751 of an ancient whole-genome duplication event. *Gigascience* [Internet] 10. Available from:  
4 752 <http://dx.doi.org/10.1093/gigascience/giab016>  
5
- 6 753 Fernández R, Gabaldon T, Dessimoz C. 2020. Orthology: Definitions, prediction, and impact on  
7 754 species phylogeny inference. *Phylogenetics in the Genomic Era*:2–4.
- 9 755 Fernández R, Kallal RJ, Dimitrov D, Ballesteros JA, Arnedo MA, Giribet G, Hormiga G. 2018.  
10 756 Phylogenomics, Diversification Dynamics, and Comparative Transcriptomics across the Spider  
11 757 Tree of Life. *Curr. Biol.* 28:2190–2193.
- 13 758 Garb JE, Ayoub NA, Hayashi CY. 2010. Untangling spider silk evolution with spidroin terminal  
14 759 domains. *BMC Evol. Biol.* 10:243.
- 16 760 Garb JE, DiMauro T, Lewis RV, Hayashi CY. 2007. Expansion and intragenic homogenization of  
17 761 spider silk genes since the Triassic: evidence from Mygalomorphae (tarantulas and their kin)  
18 762 spidroins. *Mol. Biol. Evol.* 24:2454–2464.
- 20 763 Garb JE, Sharma PP, Ayoub NA. 2018. Recent progress and prospects for advancing arachnid  
21 764 genomics. *Curr Opin Insect Sci* 25:51–57.
- 23 765 Garrison NL, Rodriguez J, Agnarsson I, Coddington JA, Griswold CE, Hamilton CA, Hedin M, Kocot  
24 766 KM, Ledford JM, Bond JE. 2016. Spider phylogenomics: untangling the Spider Tree of Life.  
25 767 *PeerJ* 4:e1719.
- 27 768 Gendreau KL, Haney RA, Schwager EE, Wierschin T, Stanke M, Richards S, Garb JE. 2017. House  
28 769 spider genome uncovers evolutionary shifts in the diversity and expression of black widow  
29 770 venom proteins associated with extreme toxicity. *BMC Genomics* 18:178.
- 31 771 Gilbert C, Peccoud J, Cordaux R. 2021. Transposable Elements and the Evolution of Insects. *Annu.*  
32 772 *Rev. Entomol.* 66:355–372.
- 34 773 Gnerre S, Maccallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, Sharpe T, Hall G, Shea TP,  
35 774 Sykes S, et al. 2011. High-quality draft assemblies of mammalian genomes from massively  
36 775 parallel sequence data. *Proc. Natl. Acad. Sci. U. S. A.* 108:1513–1518.
- 38 776 Gregory TR, Shorthouse DP. 2003. Genome sizes of spiders. *J. Hered.* 94:285–290.
- 40 777 Hakes L, Pinney JW, Lovell SC, Oliver SG, Robertson DL. 2007. All duplicates are not equal: the  
41 778 difference between small-scale and genome duplication. *Genome Biology* [Internet] 8:R209.  
42 779 Available from: <http://dx.doi.org/10.1186/gb-2007-8-10-r209>
- 44 780 Haney RA, Matte T, Forsyth FS, Garb JE. 2019. Alternative Transcription at Venom Genes and Its  
45 781 Role as a Complementary Mechanism for the Generation of Venom Complexity in the Common  
46 782 House Spider. *Front Ecol Evol* [Internet] 7. Available from:  
47 783 <http://dx.doi.org/10.3389/fevo.2019.00085>
- 49 784 Harper A, Gonzalez LB, Schönauer A, Seiter M, Holzem M, Arif S, McGregor AP, Sumner-Rooney  
50 785 L. 2020. Widespread retention of ohnologs in key developmental gene families following whole  
51 786 genome duplication in arachnoplumonates. *bioRxiv* [Internet]:2020.07.10.177725. Available  
52 787 from: <https://www.biorxiv.org/content/10.1101/2020.07.10.177725v1.abstract>
- 54 788 Herberstein ME, Wignall A. 2011. Introduction: spider biology. In: Spider behaviour: flexibility and  
55 789 versatility. Cambridge University Press. p. 1–30.
- 57 790 Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS. 2018. UFBoot2: Improving the  
58 791 Ultrafast Bootstrap Approximation. *Mol. Biol. Evol.* 35:518–522.
- 59  
60

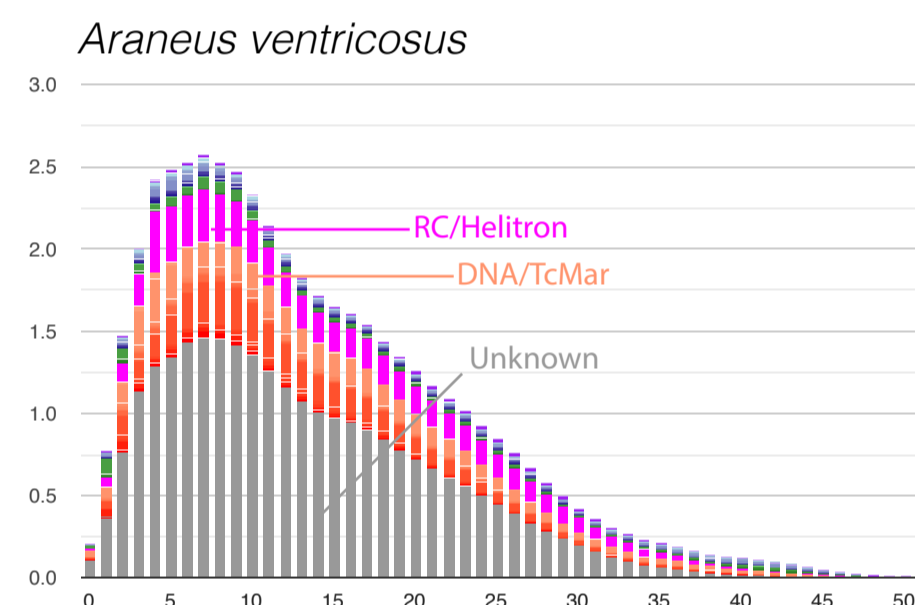
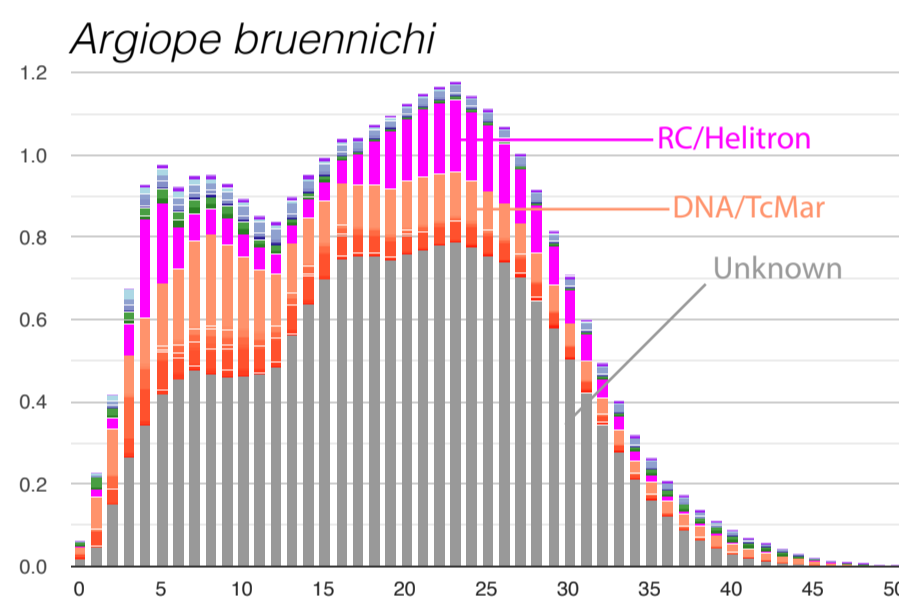
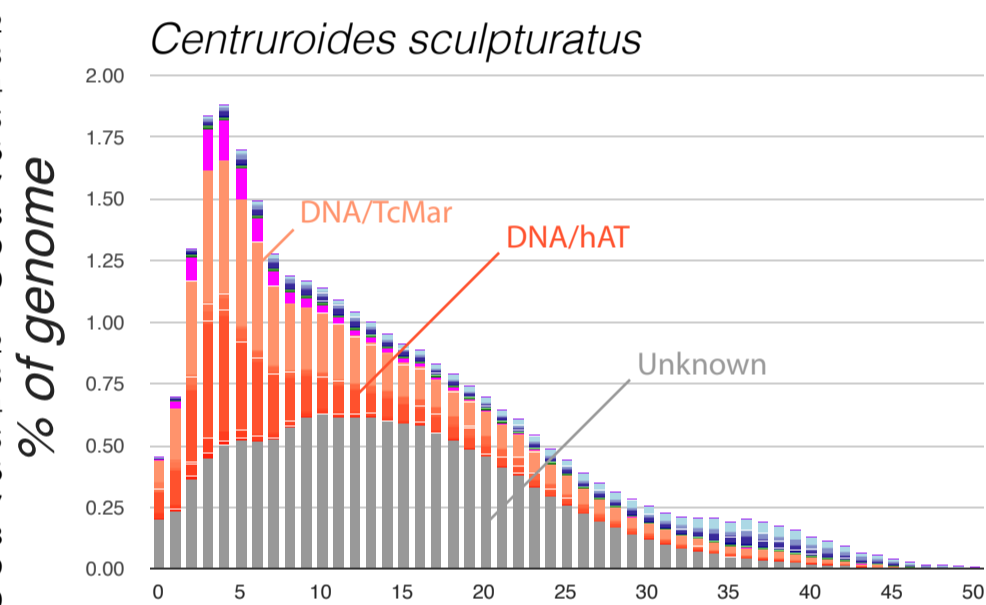
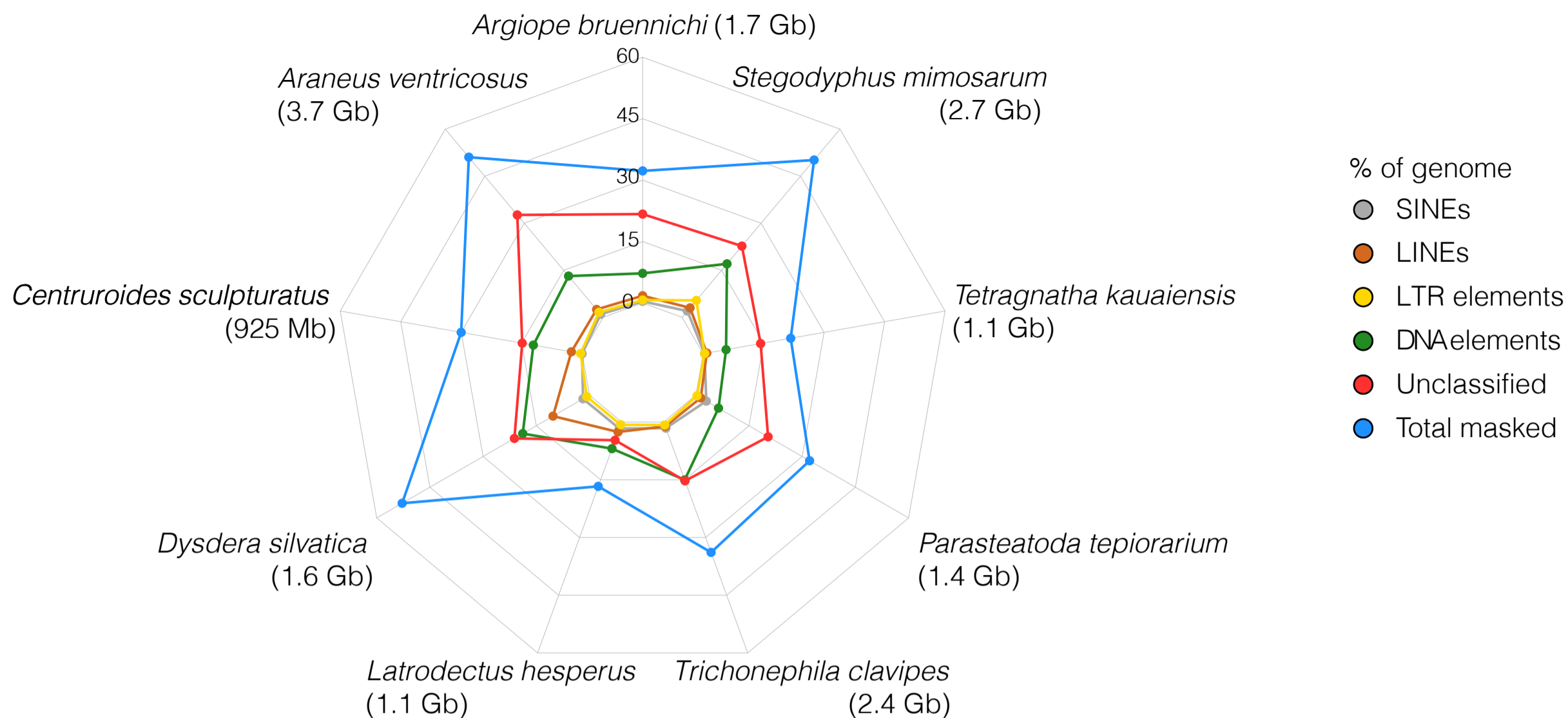
- 1  
2  
3 792 Hoff KJ, Lomsadze A, Borodovsky M, Stanke M. 2019. Whole-Genome Annotation with BRAKER.  
4 793 In: Kollmar M, editor. Gene Prediction: Methods and Protocols. New York, NY: Springer New  
5 794 York. p. 65–95.
- 6  
7 795 Huson DH, Scornavacca C. 2012. Dendroscope 3: An Interactive Tool for Rooted Phylogenetic Trees  
8 796 and Networks. *Systematic Biology* [Internet] 61:1061–1067. Available from:  
9 797 <http://dx.doi.org/10.1093/sysbio/sys062>
- 10  
11 798 Jackson RR, Cross FR. 2011. Spider Cognition. In: Casas J, editor. Advances in Insect Physiology.  
12 799 Vol. 41. Academic Press. p. 115–174.
- 13  
14 800 Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS. 2017. ModelFinder: fast  
15 801 model selection for accurate phylogenetic estimates. *Nature Methods* [Internet] 14:587–589.  
16 802 Available from: <http://dx.doi.org/10.1038/nmeth.4285>
- 17  
18 803 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:  
19 804 improvements in performance and usability. *Mol. Biol. Evol.* 30:772–780.
- 20  
21 805 Keilwagen J, Hartung F, Grau J. 2019. GeMoMa: Homology-Based Gene Prediction Utilizing Intron  
22 806 Position Conservation and RNA-seq Data. *Methods Mol. Biol.* 1962:161–177.
- 23  
24 807 King GF, Hardy MC. 2013. Spider-venom peptides: structure, pharmacology, and potential for control  
25 808 of insect pests. *Annu. Rev. Entomol.* 58:475–496.
- 26  
27 809 Koch L. 2016. Chicago HighRise for genome scaffolding. *Nat. Rev. Genet.* 17:194–194.
- 28  
29 810 Kono N, Nakamura H, Ohtoshi R, Moran DAP, Shinohara A, Yoshida Y, Fujiwara M, Mori M,  
30 811 Tomita M, Arakawa K. 2019. Orb-weaving spider *Araneus ventricosus* genome elucidates the  
31 812 spidroin gene catalogue. *Sci. Rep.* 9:8380.
- 32  
33 813 Král J, Forman M, Kořínková T, Lerma ACR, Haddad CR, Musilová J, Řezáč M, Herrera IMÁ,  
34 814 Thakur S, Dippenaar-Schoeman AS, et al. 2019. Insights into the karyotype and genome  
35 815 evolution of haplogyne spiders indicate a polyploid origin of lineage with holokinetic  
36 816 chromosomes. *Sci. Rep.* 9:3001.
- 37  
38 817 Laetsch DR, Blaxter ML. 2017. KinFin: Software for Taxon-Aware Analysis of Clustered Protein  
39 818 Sequences. *G3* 7:3349–3357.
- 40  
41 819 Leite DJ, Baudouin-Gonzalez L, Iwasaki-Yokozawa S, Lozano-Fernandez J, Turetzek N, Akiyama-  
42 820 Oda Y, Prpic N-M, Pisani D, Oda H, Sharma PP, et al. 2018. Homeobox Gene Duplication and  
43 821 Divergence in Arachnids. *Mol. Biol. Evol.* 35:2240–2253.
- 44  
45 822 Letunic I, Bork P. 2019. Interactive Tree Of Life (iTOL) v4: recent updates and new developments.  
46 823 *Nucleic Acids Res.* 47:W256–W259.
- 47  
48 824 Li Z, Tiley GP, Galuska SR, Reardon CR, Kidder TI, Rundell RJ, Barker MS. 2018. Multiple large-  
49 825 scale gene and genome duplications during the evolution of hexapods. *Proc. Natl. Acad. Sci. U.*  
50 826 *S. A.* 115:4713–4718.
- 51  
52 827 Lynch M, Conery JS. 2000. The evolutionary fate and consequences of duplicate genes. *Science*  
53 828 290:1151–1155.
- 54  
55 829 Lynch M, Force A. 2000. The probability of duplicate gene preservation by subfunctionalization.  
56 830 *Genetics* 154:459–473.
- 57  
58 831 McClintock B. 1984. The significance of responses of the genome to challenge. *Science* 226:792–801.
- 59  
60

- 1  
2  
3 832 McGregor AP, Hilbrant M, Pechmann M, Schwager EE, Prpic N-M, Damen WGM. 2008. *Cupiennius*  
4 833 *salei* and *Achaearanea tepidariorum*: Spider models for investigating evolution and development.  
5 834 *BioEssays* [Internet] 30:487–498. Available from: <http://dx.doi.org/10.1002/bies.20744>  
6  
7 835 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective stochastic  
8 836 algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32:268–274.  
9  
10 837 Nieto Feliner G, Casacuberta J, Wendel JF. 2020. Genomics of Evolutionary Novelty in Hybrids and  
11 838 Polyploids. *Front. Genet.* 11:792.  
12  
13 839 Ohno S. 1970. The enormous diversity in genome sizes of fish as a reflection of nature’s extensive  
14 840 experiments with gene duplication. *Trans. Am. Fish. Soc.* 99:120–130.  
15  
16 841 Oxford GS, Gillespie RG. 1998. Evolution and ecology of spider coloration. *Annu. Rev. Entomol.*  
17 842 43:619–643.  
18  
19 843 Peona V, Weissensteiner MH, Suh A. 2018. How complete are “complete” genome assemblies?-An  
20 844 avian perspective. *Mol. Ecol. Resour.* 18:1188–1195.  
21  
22 845 Pineda SS, Chaumeil P-A, Kunert A, Kaas Q, Thang MWC, Le L, Nuhn M, Herzig V, Saez NJ,  
23 846 Cristofori-Armstrong B, et al. 2018. ArachnoServer 3.0: an online resource for automated  
24 847 discovery, analysis and annotation of spider toxins. *Bioinformatics* 34:1074–1076.  
25  
26 848 Putnam NH, O’Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, Troll CJ, Fields A, Hartley  
27 849 PD, Sugnet CW, et al. 2016. Chromosome-scale shotgun assembly using an in vitro method for  
28 850 long-range linkage. *Genome Res.* 26:342–350.  
29  
30 851 Sánchez-Gracia A, Vieira FG, Rozas J. 2009. Molecular evolution of the major chemosensory gene  
31 852 families in insects. *Heredity* 103:208–216.  
32  
33 853 Sánchez-Herrero JF, Frías-López C, Escuer P, Hinojosa-Alvarez S, Arnedo MA, Sánchez-Gracia A,  
34 854 Rozas J. 2019. The draft genome sequence of the spider *Dysdera silvatica* (Araneae,  
35 855 *Dysderidae*): A valuable resource for functional and evolutionary genomic studies in  
36 856 chelicerates. *Gigascience* [Internet] 8. Available from:  
37 857 <http://dx.doi.org/10.1093/gigascience/giz099>  
38  
39 858 Sanderson MJ. 2003. r8s: inferring absolute rates of molecular evolution and divergence times in the  
40 859 absence of a molecular clock. *Bioinformatics* 19:301–302.  
41  
42 860 Sanggaard KW, Bechsgaard JS, Fang X, Duan J, Dyrland TF, Gupta V, Jiang X, Cheng L, Fan D,  
43 861 Feng Y, et al. 2014. Spider genomes provide insight into composition and evolution of venom  
44 862 and silk. *Nat. Commun.* 5:3765.  
45  
46 863 Schmickl R, Yant L. 2021. Adaptive introgression: how polyploidy reshapes gene flow landscapes.  
47 864 *New Phytol.* 230:457–461.  
48  
49 865 Schwager EE, Schoppmeier M, Pechmann M, Damen WGM. 2007. Duplicated Hox genes in the  
50 866 spider *Cupiennius salei*. *Front. Zool.* 4:10.  
51  
52 867 Schwager EE, Sharma PP, Clarke T, Leite DJ, Wierschin T, Pechmann M, Akiyama-Oda Y, Esposito  
53 868 L, Bechsgaard J, Bilde T, et al. 2017. The house spider genome reveals an ancient whole-genome  
54 869 duplication during arachnid evolution. *BMC Biol.* 15:62.  
55  
56 870 Seppely M, Manni M, Zdobnov EM. 2019. BUSCO: Assessing Genome Assembly and Annotation  
57 871 Completeness. *Methods Mol. Biol.* 1962:227–245.  
58  
59 872 Sheffer MM, Hoppe A, Krehenwinkel H, Uhl G, Kuss AW, Jensen L, Jensen C, Gillespie RG, Hoff

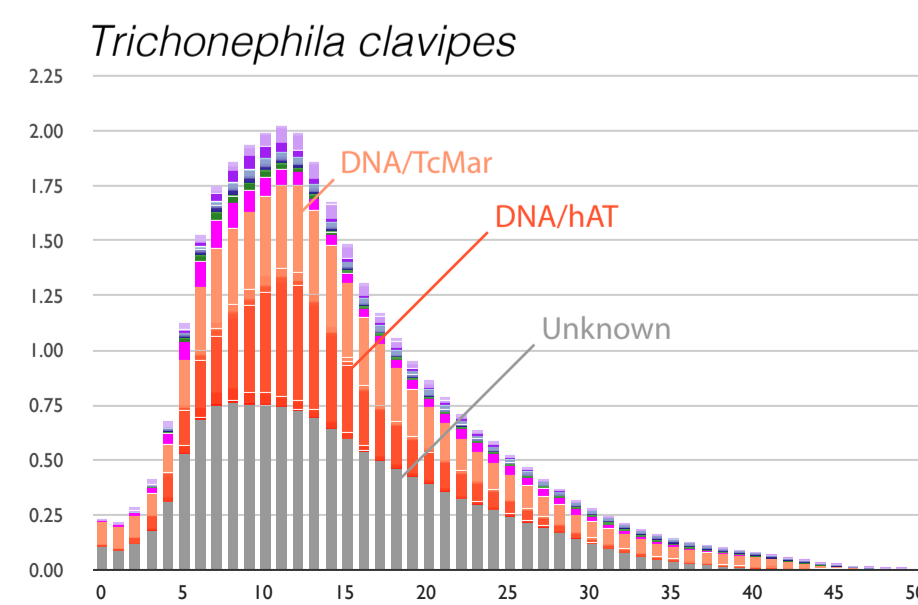
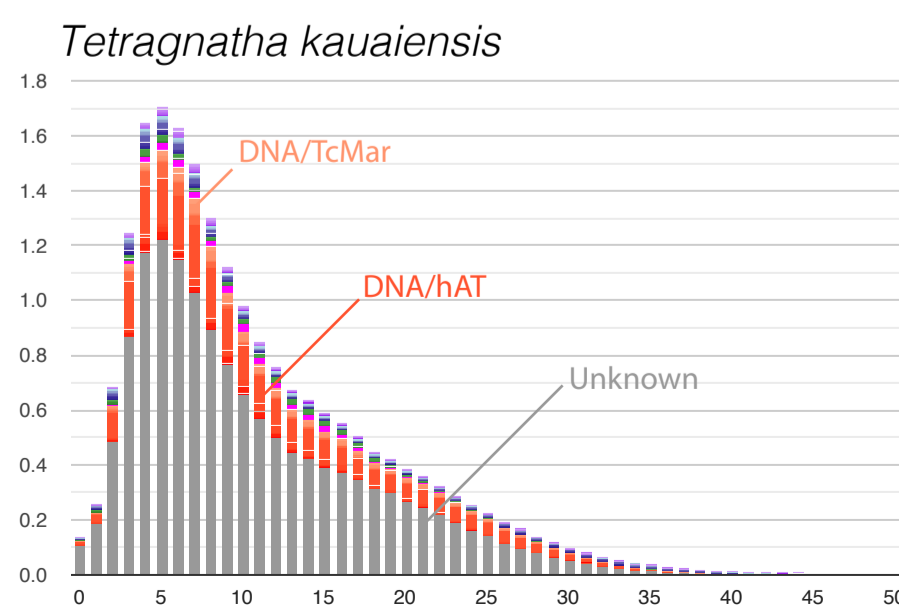
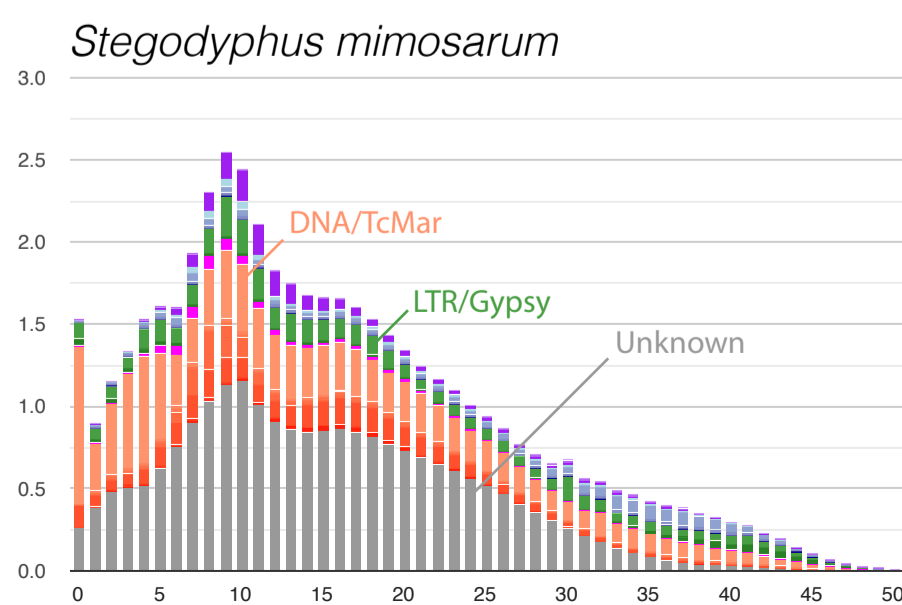
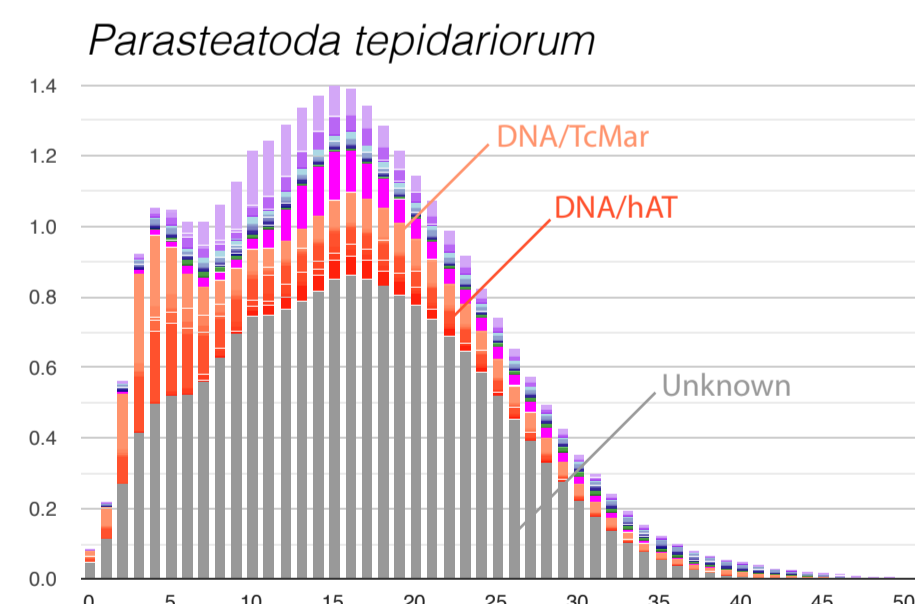
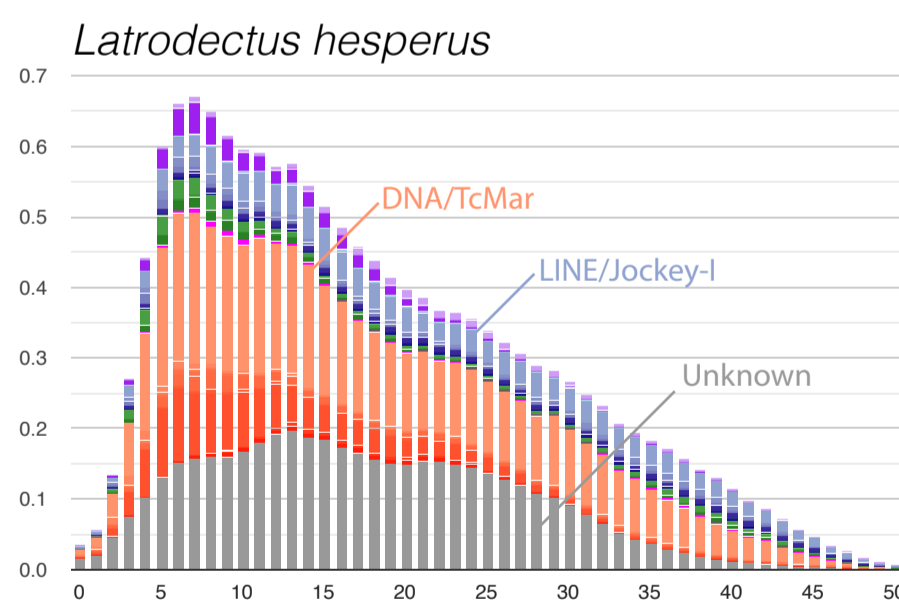
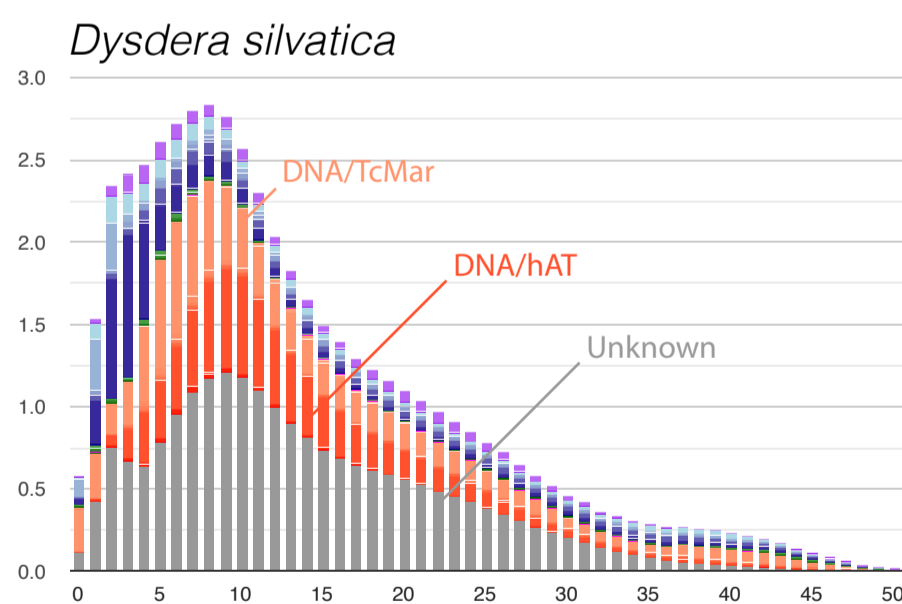
- 1  
2  
3 873 KJ, Prost S. 2021. Chromosome-level reference genome of the European wasp spider *Argiope*  
4 874 *bruennichi*: a resource for studies on range expansion and evolutionary adaptation. *Gigascience*  
5 875 [Internet] 10. Available from: <http://dx.doi.org/10.1093/gigascience/giaa148>  
6
- 7 876 Shingate P, Ravi V, Prasad A, Tay B-H, Garg KM, Chattopadhyay B, Yap L-M, Rheindt FE,  
8 877 Venkatesh B. 2020. Chromosome-level assembly of the horseshoe crab genome provides insights  
9 878 into its genome evolution. *Nat. Commun.* 11:2322.
- 10  
11 879 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing  
12 880 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*  
13 881 31:3210–3212.
- 14  
15 882 Supek F, Bošnjak M, Škunca N, Šmuc T. 2011. REVIGO summarizes and visualizes long lists of  
16 883 gene ontology terms. *PLoS One* 6:e21800.
- 17  
18 884 Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to identify repetitive elements in genomic  
19 885 sequences. *Curr. Protoc. Bioinformatics* Chapter 4:Unit 4.10.
- 20  
21 886 Team RC, Others. 2013. R: A language and environment for statistical computing. Available from:  
22 887 <https://cran.microsoft.com/snapshot/2014-09-08/web/packages/dpLR/vignettes/xdate-dpLR.pdf>
- 23  
24 888 Thomas GWC, Dohmen E, Hughes DST, Murali SC, Poelchau M, Glastad K, Anstead CA, Ayoub  
25 889 NA, Batterham P, Bellair M, et al. 2020. Gene content evolution in the arthropods. *Genome Biol.*  
26 890 21:15.
- 27  
28 891 Vieira FG, Sánchez-Gracia A, Rozas J. 2007. Comparative genomic analysis of the odorant-binding  
29 892 protein family in 12 *Drosophila* genomes: purifying selection and birth-and-death evolution.  
30 893 *Genome Biol.* 8:R235.
- 31  
32 894 Vienneau-Hathaway JM, Brassfield ER, Lane AK, Collin MA, Correa-Garhwal SM, Clarke TH,  
33 895 Schwager EE, Garb JE, Hayashi CY, Ayoub NA. 2017. Duplication and concerted evolution of  
34 896 MiSp-encoding genes underlie the material properties of minor ampullate silks of cobweb  
35 897 weaving spiders. *BMC Evol. Biol.* 17:78.
- 36  
37 898 Vizueta J, Escuer P, Sánchez-Gracia A, Rozas J. 2020. Genome mining and sequence analysis of  
38 899 chemosensory soluble proteins in arthropods. *Methods Enzymol.* 642:1–20.
- 39  
40 900 Vizueta J, Frías-López C, Macías-Hernández N, Arnedo MA, Sánchez-Gracia A, Rozas J. 2017.  
41 901 Evolution of Chemosensory Gene Families in Arthropods: Insight from the First Inclusive  
42 902 Comparative Transcriptome Analysis across Spider Appendages. *Genome Biol. Evol.* 9:178–196.
- 43  
44 903 Vizueta J, Macías-Hernández N, Arnedo MA, Rozas J, Sánchez-Gracia A. 2019. Chance and  
45 904 predictability in evolution: The genomic basis of convergent dietary specializations in an  
46 905 adaptive radiation. *Mol. Ecol.* 28:4028–4045.
- 47  
48 906 Vizueta J, Rozas J, Sánchez-Gracia A. 2018. Comparative Genomics Reveals Thousands of Novel  
49 907 Chemosensory Genes and Massive Changes in Chemoreceptor Repertoires across Chelicerates.  
50 908 *Genome Biology and Evolution* [Internet] 10:1221–1236. Available from:  
51 909 <http://dx.doi.org/10.1093/gbe/evy081>  
52 910
- 53  
54 910 Vizueta J, Sánchez-Gracia A, Rozas J. 2020. bitacora: A comprehensive tool for the identification and  
55 911 annotation of gene families in genome assemblies. *Mol. Ecol. Resour.* 20:1445–1452.
- 56  
57 912 Vollrath F. 1999. Biology of spider silk. *Int. J. Biol. Macromol.* 24:81–88.
- 58  
59 913 Wilder SM. 2011. Spider Nutrition: An Integrative Perspective. In: Casas J, editor. *Advances in Insect*  
60 914 *Physiology*. Vol. 40. Academic Press. p. 87–136.

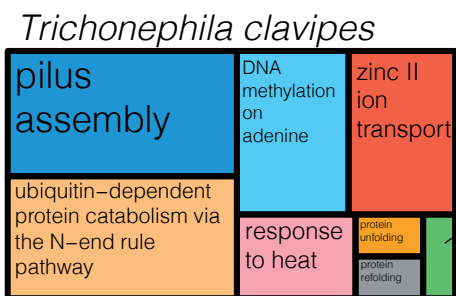
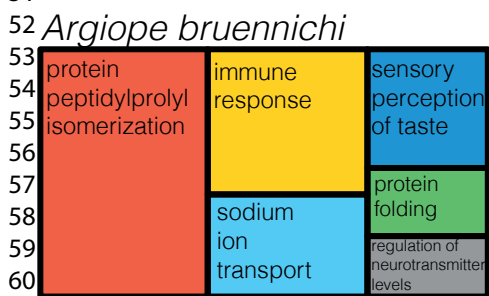
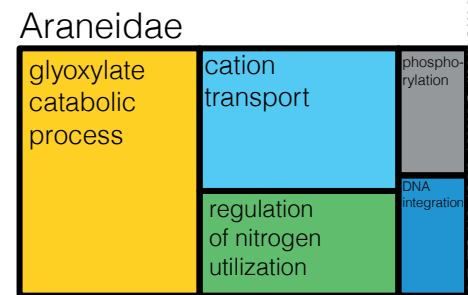
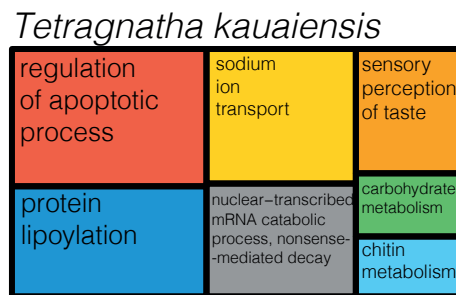
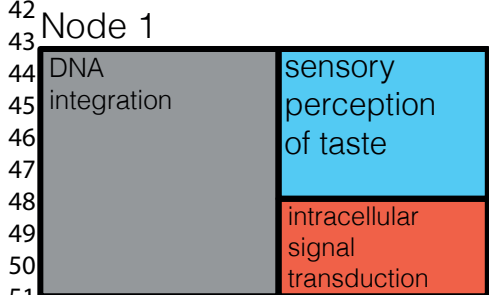
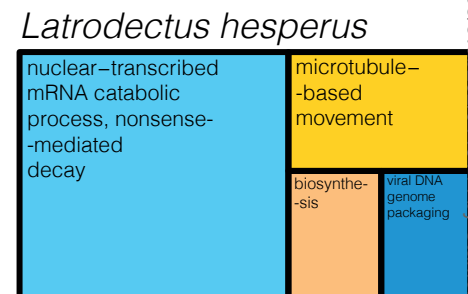
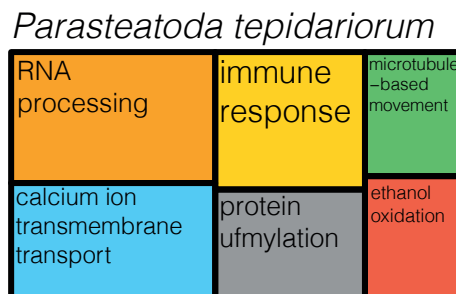
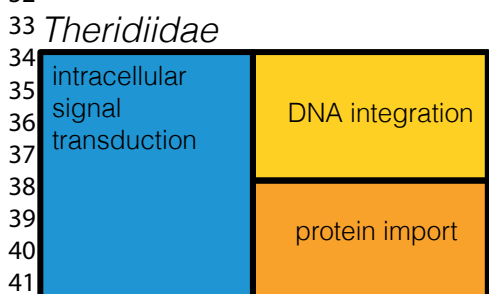
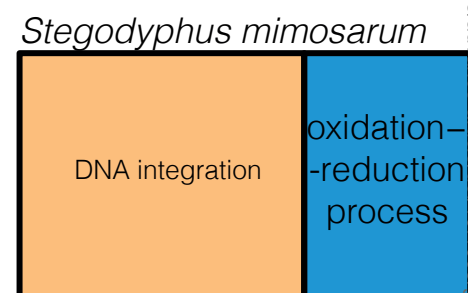
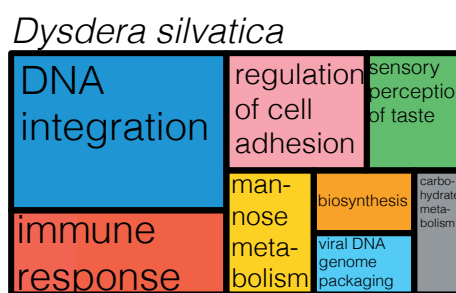
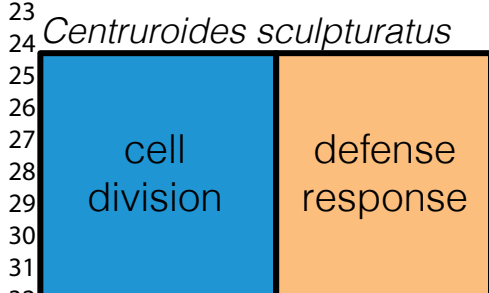
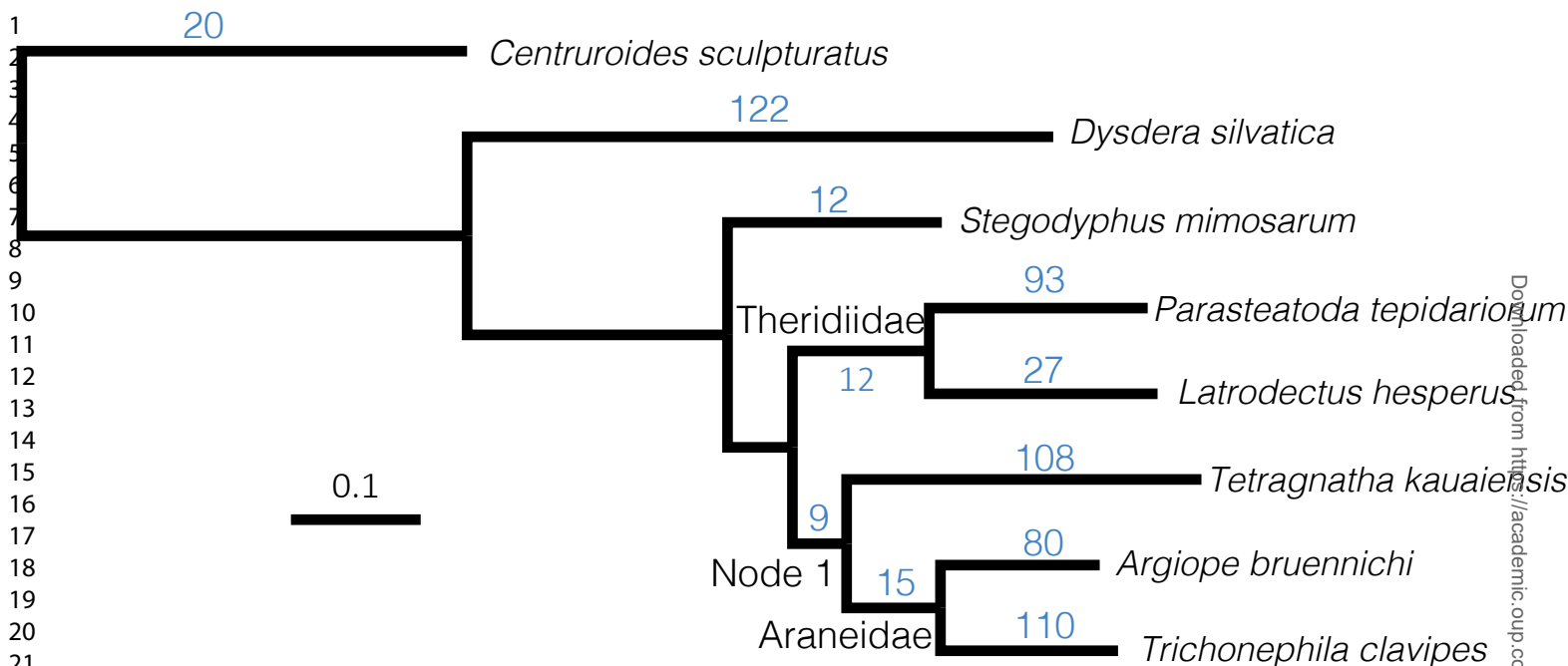
- 1  
2  
3 915 Wu C, Lu J. 2019. Diversification of Transposable Elements in Arthropods and Its Impact on Genome  
4 916 Evolution. *Genes* [10](http://dx.doi.org/10.3390/genes10050338). Available from: <http://dx.doi.org/10.3390/genes10050338>  
5  
6 917 Yim KM, Brewer MS, Miller CT, Gillespie RG. 2014. Comparative transcriptomics of maturity-  
7 918 associated color change in Hawaiian spiders. *J. Hered.* 105 Suppl 1:771–781.  
8  
9 919  
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





Kimura substitution level (CpG adjusted)



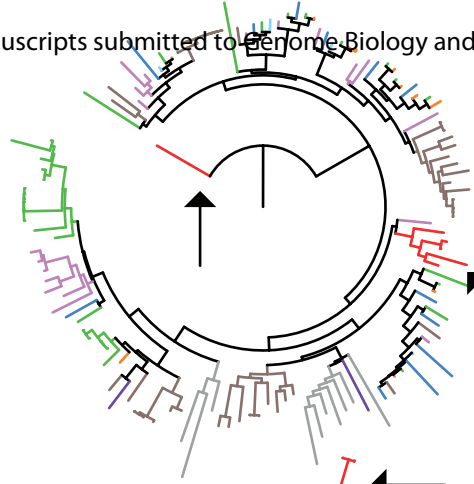


regulation of nitrogen utilization

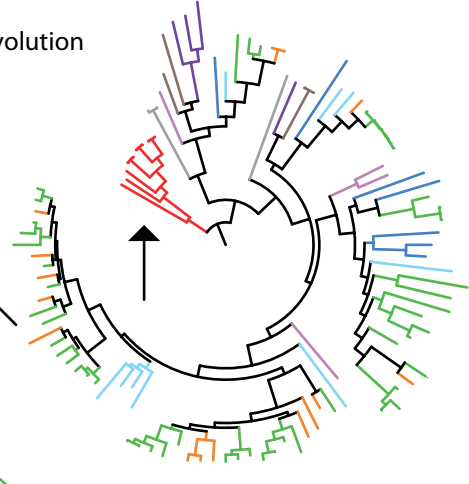
Downloaded from https://academic.oup.com/gbe/advance-article-abstract/doi/10.1093/gbe/evab025/6449441 by Stanford Medical Center user on 02 December 2021

- 1 *Argiope bruennichi*
- 2 *Araneus ventricosus*
- 3 *Gentruoides sculpturatus*
- 4 *Lysdera silvatica*
- 5 *Latrodectus hesperus*
- 6 *Parasteatoda tepidariorum*
- 7 *Stegodyphus mimosarum*
- 8 *Tetragnatha kauaiensis*
- 9 *Trichonephila clavipes*

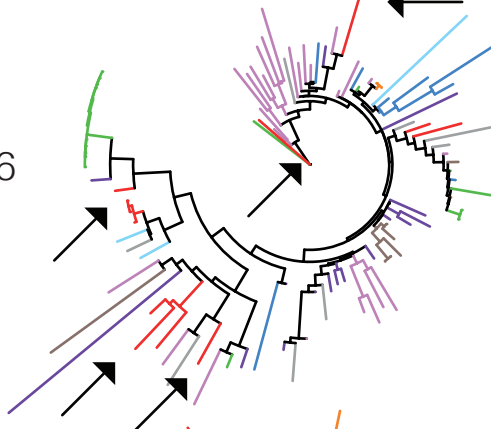
OG0000175  
Tree scale



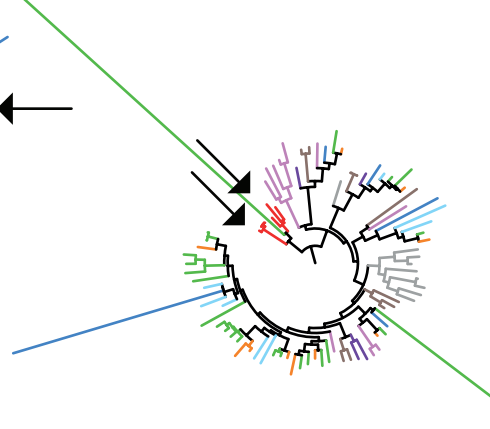
OG0000314  
Tree scale



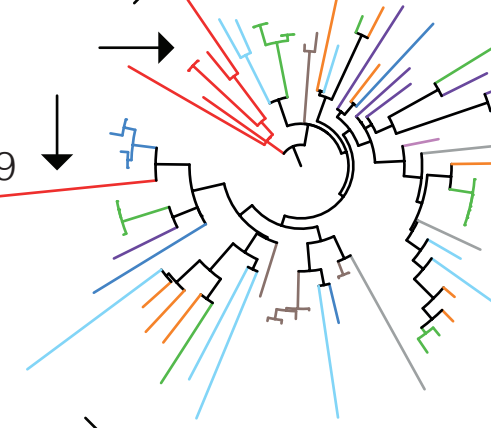
OG0000346  
Tree scale



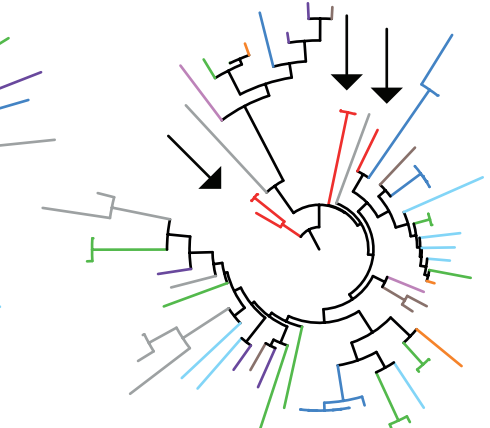
OG0000432  
Tree scale



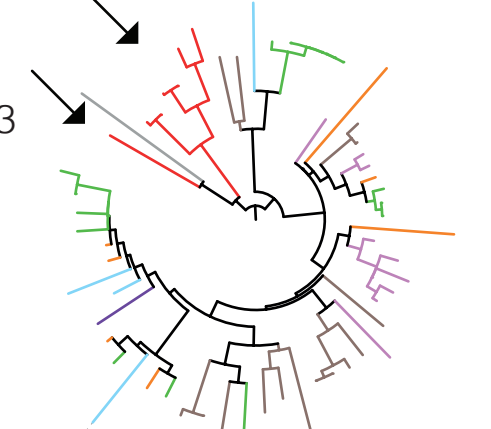
OG0000639  
Tree scale



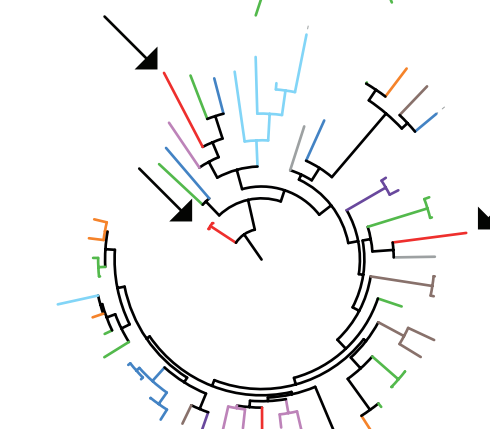
OG0000761  
Tree scale



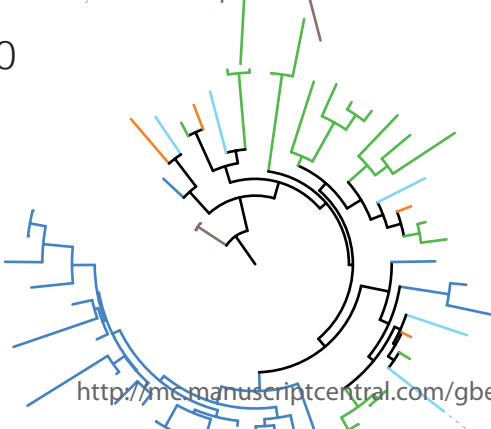
OG0000803  
Tree scale



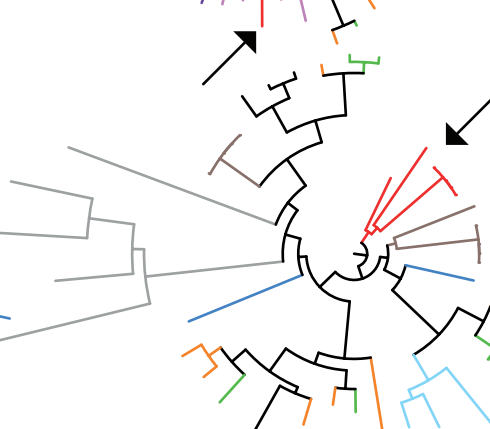
OG0000906  
Tree scale



OG0000930  
Tree scale



OG0001436  
Tree scale



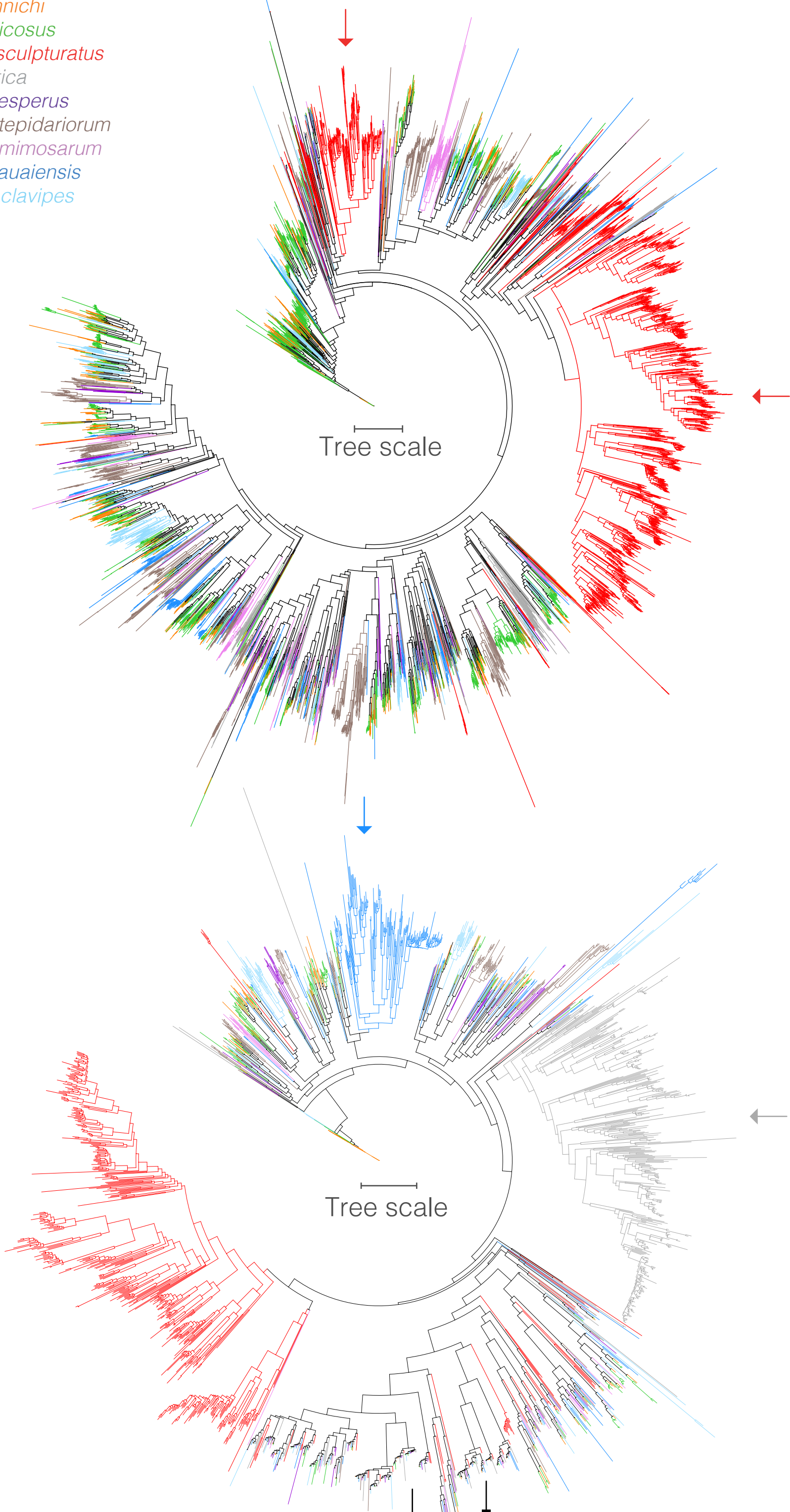
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

Scorpion genes indicated by arrows →

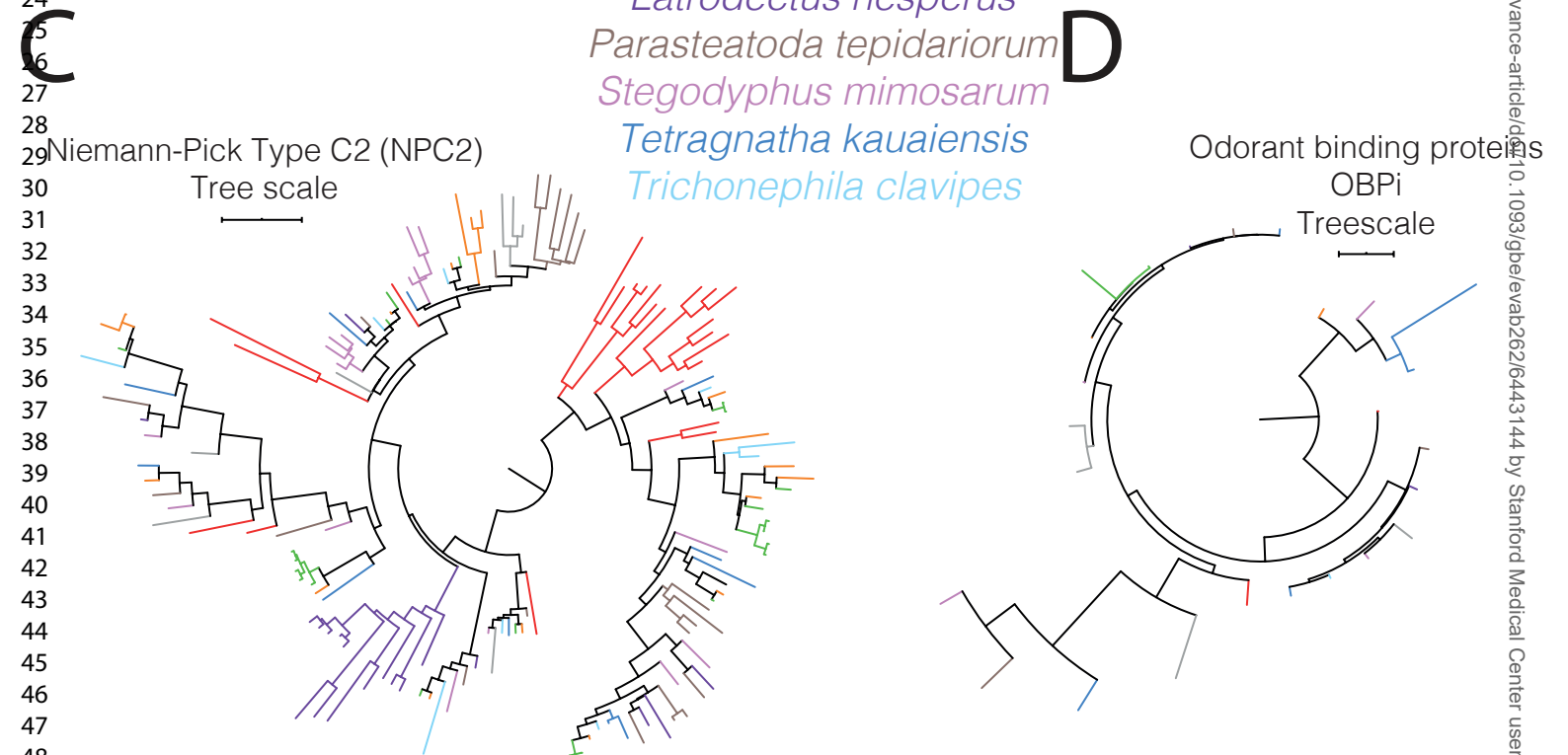
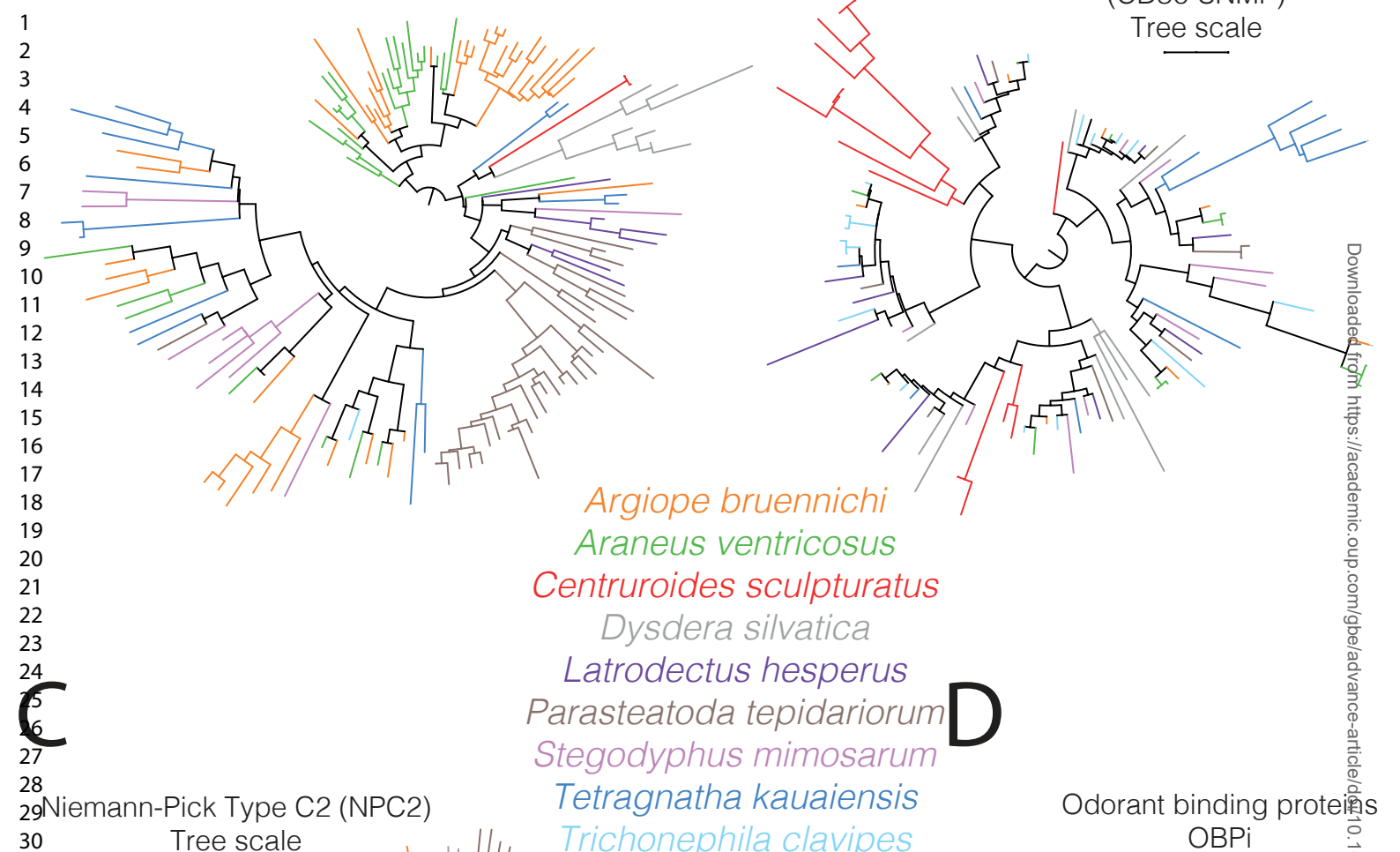
duplications signaled by arrows →

conserved clusters signaled by squared arrows ◼

- 1
- 2
- 3
- 4
- 5 *Argiope bruennichi*
- 6 *Araneus ventricosus*
- 7 *Centruroides sculpturatus*
- 8 *Dysdera silvatica*
- 9 *Latrodectus hesperus*
- 10 *Parasteatoda tepidariorum*
- 11 *Stegodyphus mimosarum*
- 12 *Tetragnatha kauaiensis*
- 13 *Trichonephila clavipes*
- 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



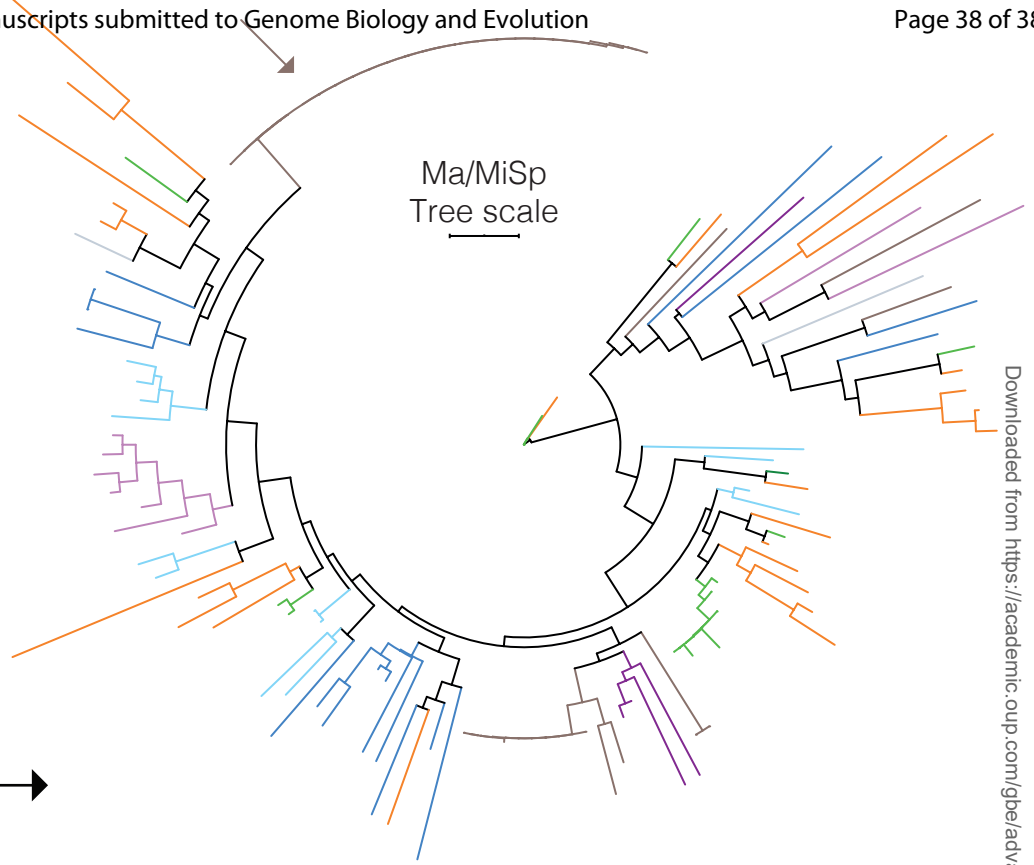
Downloaded from https://academic.oup.com/gbe/advance-article/doi/10.1093/gbe/evab262/6443144 by Stanford Medical Center user on 02 December 2021



Downloaded from <https://academic.oup.com/gbe/advance-article/doi/10.1093/gbe/evab262/6443144> by Stanford Medical Center user on 02 December 2021

- 1
- 2
- 3
- 4 *Argiope bruennichi*
- 5 *Araneus ventricosus*
- 6 *Centruroides sculpturatus*
- 7
- 8
- 9 *Dysdera silvatica*
- 10 *Latrodectus hesperus*
- 11
- 12 *Parasteatoda tepidariorum*
- 13 *Stegodyphus mimosarum*
- 14
- 15 *Tetragnatha kauaiensis*
- 16 *Trichonephila clavipes*
- 17
- 18
- 19

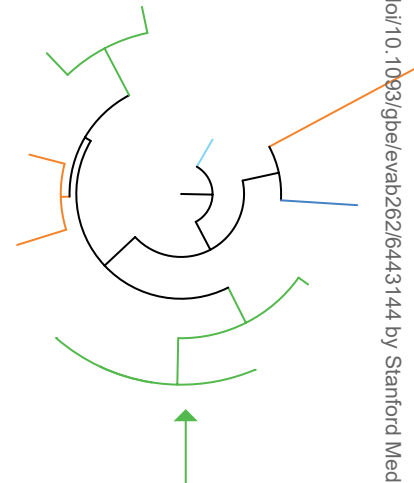
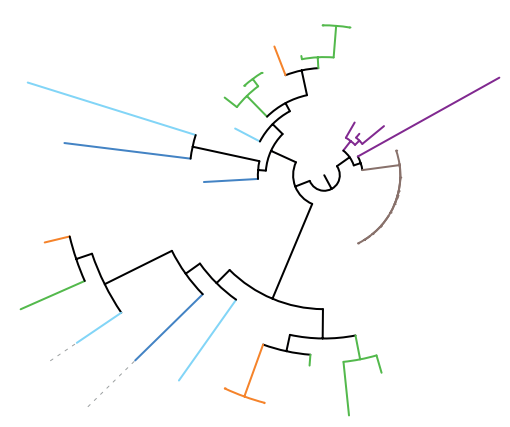
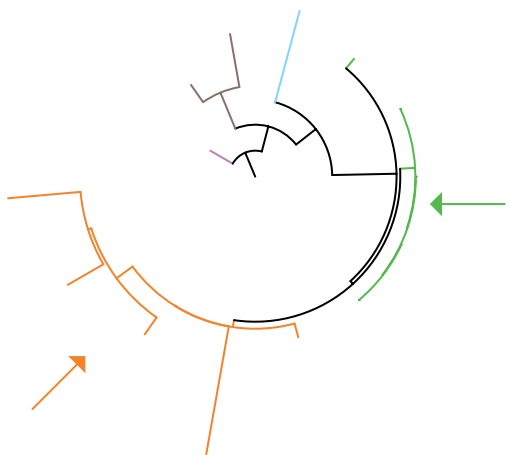
20 putative duplications  
 21 signaled by arrows →



26 AcSp  
Tree scale

27 AgSp  
Tree scale

28 Flag  
Tree scale



44 Unidentified Sp  
Tree scale

46 PySp  
Tree scale

48 Tu Sp  
Tree scale

