

Genome sequencing and neurotoxin diversity of a wandering spider *Pardosa pseudoannulata* (pond wolf spider)

Na Yu^{1,&}, Jingjing Li^{1,&}, Meng Liu^{2,&}, Lixin Huang¹, Haibo Bao¹, Zhiming Yang¹,
Yixi Zhang¹, Haoli Gao¹, Zhaoying Wang¹, Yuanxue Yang¹, Thomas Van Leeuwen³,
Neil S. Millar⁴, Zewen Liu^{1,*}

1, Key laboratory of Integrated Management of Crop Diseases and Pests (Ministry of Education), College of Plant Protection, Nanjing Agricultural University, Weigang 1, Nanjing 210095, China

2, Novogene Bioinformatics Institute, Beijing 100083, China

3, Faculty of Bioscience Engineering, Ghent University, Coupure links 653, Ghent B9000, Belgium

4, Department of Neuroscience, Physiology & Pharmacology, University College London, Gower Street, London WC1E 6BT, United Kingdom

*, Correspondence should be addressed to Zewen Liu (Key laboratory of Integrated Management of Crop Diseases and Pests (Ministry of Education), College of Plant Protection, Nanjing Agricultural University, Weigang 1, Nanjing 210095, P.R. China, Tel: +86 25 84399051; E-mail: liuzewen@njau.edu.cn.

&, these authors contributed equally and should be regarded as co-first authors.

1 **Abstract**

2 Spiders constitute an extensive and diverse branch of the phylum Arthropoda.
3 Whereas the genomes of four web-weaver spider species and a single cave-living
4 spider have been determined, similar studies have not been reported previously for a
5 wandering spider. The pond wolf spider, *Pardosa pseudoannulata*, is a wandering
6 hunter that immobilizes prey using venom rather than a web. It is also an important
7 predator against a range of agriculturally important insect pests. The increasing
8 interest in its wandering lifestyle and in the potential of spider venom as a tool for
9 pest control have prompted a detailed study on this wandering spider species. We
10 have generated a high-quality genome sequence of *P. pseudoannulata* and analysed
11 genes associated with the production of silk and venom toxins. Sequencing reveals
12 that *P. pseudoannulata* has a large genome of 4.26 Gb. The presence of only 16
13 spidroin genes and four types of silk glands is consistent with the moderate use of silk
14 and the lack of a prey-catching web. A large number of genes encode neurotoxins and
15 there is evidence that the majority are highly selective for invertebrates. Comparison
16 between spider species reveals a correlation between spider neurotoxin selectivity for
17 target organisms and spider prosoma size, suggesting a possible coevolution of these
18 two features. The genome data provides valuable insights into the biology of *P.*
19 *pseudoannulata* and its potential role as a natural enemy in pest control.

20 **Keywords:** *Pardosa pseudoannulata*; wandering spider; genome; spidroin; venom;
21 neurotoxin

22 Introduction

23 Spiders are an important group of arthropods with diverse biological and behavioural
24 characteristics, as is illustrated by their use of both silk and venom to incapacitate
25 prey. Some species of spider build webs for a variety of biological functions, but most
26 notably, for the capture of prey ¹. In contrast, other species, including the pond wolf
27 spider, *Pardosa pseudoannulata*, have a wandering lifestyle and employ venom for
28 predation and defence (Fig. 1a). The increasing availability of spider genomic and
29 transcriptomic data is helping to provide a better understanding of their biological and
30 evolutionary importance ². In-depth genome sequencing has been reported previously
31 for three Araneoidea spider species that produce prey-catching webs: the velvet spider
32 *Stegodyphus mimosarum* ³ and *Stegodyphus dumicola* ⁴, the common house spider
33 *Parasteatoda tepidariorum* ^{5,6} and the golden-orb weaver *Nephila clavipes* ⁷. This has
34 helped to provide information on evolutionary relationships and insights into
35 phenomena such as the diversity of genes that are involved in the production of silk
36 proteins and venom. However, despite the importance and diversity of wandering
37 spiders, only a draft genome has been reported (with 40% coverage) for a sit-and-wait
38 spider (*Acanthoscurria geniculata*, a cave-living species) ³.

39 Here, we report the first high-quality genome sequencing of a wandering spider,
40 the pond wolf spider, *Pardosa pseudoannulata*, which belongs to the retrolateral tibial
41 apophysis (RTA) clade of Araneomorphae. An important motivation for undertaking
42 this project was that this would enable a comparison between the genomes of web-
43 building and wandering species, thereby providing insights into their adaptation to
44 differing lifestyles. In addition, *P. pseudoannulata* are predators to a range of insect
45 pests that are of agricultural importance and, as a consequence, *P. pseudoannulata* has
46 been identified as a possible biological control agent. The complete genome

47 sequencing of *P. pseudoannulata* provides a wealth of valuable information,
48 particularly concerning its potential use for insect pest control in integrated pest
49 management.

50 **Results**

51 High quality genomic DNA was extracted from *P.pseudoannulata* adults for
52 sequencing via Illumina and PacBio technologies. Short-insert (250 bp and 350 bp)
53 paired-end libraries, large-insert (2 kb, 5 kb, 10 kb, 15 kb and 20 kb) mate-pair
54 libraries and 10X Genomics linked-read library were sequenced on the Illumina
55 platform and generated 1771.69 Gb raw data (404.03 x coverage). SMRTbell libraries
56 were sequenced on PacBio Sequel platform and generated 87.37 Gb raw data (19.92 x
57 coverage) (Table S1). Raw data and subreads were filtered. The genome size was
58 estimated via k-mer frequency distribution to be ~4.39 Gb (Table S2). Transcriptome
59 sequencing of four pairs of legs, pedipalp, chelicerae, brain, venom gland, fat body,
60 male silk gland and female silk gland was performed and each generated ~7 Gb raw
61 data (Table S3). A draft genome of 4.27 Gb was eventually assembled with a contig
62 N50 of 22.82 kb and scaffold N50 of 699.15 kb (Table 1, Table S4) with GC content
63 counting for 31.36% (Table S5). The reliability and completeness of the genome
64 assembly were evaluated with EST, CEGMA and BUSCO, with 98.79% of the raw
65 sequence reads aligned to the assembly (Table S6), 84.63% mapped to the genome
66 assembly and 76.71% fully covered by one scaffold with more than 90% of the
67 transcripts mapped to one scaffold (Table S7), 91.53% of conserved eukaryotic genes
68 found (Table S8), and 93.7% of the BUSCO dataset identified in the genome
69 assembly (Table S9).

70 **Genome statistics and phylogenomics**

71 The assembled genome was ~4.26 Gb with a contig N50 of 22.82 kb and scaffold
72 N50 of 699.15 kb (Table 1). The repeat sequences accounted for ~51.40% of the *P.*
73 *pseudoannulata* genome and DNA transposons (33.46% of the genome) formed the
74 most abundant category among the TEs, followed by long interspersed elements
75 (LINEs, 3.34%) (Table S10, S11). The gene set contained 23,310 genes, of which
76 98.8% were supported by homologous evidences or transcriptomic data (Fig S1,
77 Table S12, S13). 19,602 protein-coding genes (accounting for 92.0% of the 21,310
78 genes) were annotated with at least one public database (Table S14). It was notable
79 that *P. pseudoannulata* genes were generally composed of short exons and long
80 introns, a typical structural feature for Arachnid genes³ (Table S12, S14).

81 **Table 1. Summary of the *P. pseudoannulata* genome sequence data**

Total sequencing data	1,859.06 Gb
Sequence coverage	423.95 x
Estimated genome size	4.42 Gb
Assembled genome size	4.26 Gb
Repeat content	51.39%
GC content	31.36%
Number of contigs	1,242,313
N50 contig size	22.82 kb
Largest contig	650 kb
Number of scaffolds	1,041,65
N50 scaffold size	699.15 kb
Largest scaffold	8,106.74 kb
Number of protein-coding genes	21,310

Annotated gene number	19,602
Average exon length	179.39 bp
Average intron length	4,132.1 bp

82

83 The phylogenetic relationship of *P. pseudoannulata* with 17 other selected species
84 (Table S15) was analysed using 190 single-copy gene families. According to the
85 phylogenetic analysis, *P. pseudoannulata* diverged from the common ancestral of *S.*
86 *mimosarum* approximately 135.9 million years ago (MYA) and the lineage of *P.*
87 *tepidariorum* and *N. clavipes* diverged from the lineage leading to *P. pseudoannulata*
88 and *S. mimosarum* ~159 MYA (Fig. 1b, Fig. S2). The placing of spiders, ticks,
89 scorpions and mites supported the polyphyletic nature of the Acari ^{3,8}.

90 **Gene family expansion**

91 In *P. pseudoannulata* genome, 11 gene families showed significant expansion and 23
92 gene families showed significant contraction (Fig. 1b). Enrichment of GO terms and
93 KEGG for *P. pseudoannulata* expanded families were performed with the
94 EnrichPipeline ⁹. A false discovery rate (FDR) threshold of 0.05 was used to define
95 GO terms and KEGG that were significantly enriched. Predominantly enriched
96 functional categories for these genes contained several metabolic processes and ion
97 binding that may be related to the environmental adaptation for *P. pseudoannulata*
98 (Table S16, S17).

99 **Spidroin gene and silk gland**

100 *P. pseudoannulata* typically hunt small insects via wandering in the field rather than
101 building prey-catching webs. Therefore, it is of interest to compare the diversity of
102 spidroin genes in a wandering spider compared to that in web waver spiders. In total,
103 sixteen *P. pseudoannulata* putative spidroin genes were obtained and spidroins were

104 classified based on their sequence homology to *N. clavipes* spidroins (Fig. 2a, Table
105 S18). Twelve spidroin genes were designated into five spidroin types as the major
106 ampullate (MaSp, 6), minor ampullate (MiSp, 1), aciniform (AcSp, 3), piriform (PiSp,
107 1) and tubuliform (TuSp, 1) spidroins based on the sequence alignment of the N-
108 terminal domains with those of *N. clavipes*. The other 4 genes were designated as
109 spidron (Sp) due to lack of clear evidence in sequence similarity (Fig. S3, S4). Up to
110 28 spidroins were catalogued into 7 spidroin types in *N. clavipes*⁷ and 19 putative
111 spidroin genes were annotated in *S. mimosarum* genome³. The *P. psuedoannulata*
112 spidroin gene repertoire lacked the flagelliform and aggregate spidroins, but this was
113 not surprising given that they mainly function in prey-catching webs^{1,10-14}. The
114 fifteen complete spidroin proteins ranged from 573 (AcSp_2727) to 1977
115 (MaSp_2831) amino acid residues. Typically, these genes contained only one or two
116 exons but up to 12 exons were also present (e.g. Sp_48488). The spidroins all
117 contained the canonical structure of the repeat region (R) flanked by the relatively
118 conserved N-terminal domain (N) and the C-terminal domain (C) (Fig. 2a, Table S18).

119 Dissection of adult male and female *P. pseudoannulata* identified only four types
120 of silk glands: the major ampullate gland (Ma), minor ampullate gland (Mi),
121 aciniform gland (Ac) and piriform gland (Pi), which were morphologically similar to
122 the reported black widow spider silk glands (Fig. 2b)^{15,16}. Typically seven types of
123 silk glands are present in orb-weaver spiders⁷ and the three absent silk gland types in
124 *P. psuedoannulata* were the tubuliform glands, flagelliform glands and aggregate
125 glands. The observed differences in spidroin genes, spidroin types and silk gland
126 types between *P. pseudoannulata* and web-weaver spiders support the conclusion that
127 silk gland types and silk proteins are specialized in tasks involved in a variety of

128 biological activities^{17,18}. However, whether the spidroin genes were lost in wandering
129 spiders or expanded and differentiated in web-weavers requires further investigation.

130 The relative expression of the 16 spidroin genes in female and male *P.*
131 *pseudoannulata* were quantified. Generally, most spidroin genes were expressed at
132 higher levels in males than in females, however, TuSp_4038 and Sp_48488 were
133 expressed at higher levels in females (Fig. 2c, Table S20). These differences in
134 spidroin expression presumably reflected differences in the requirement for spidroins
135 in each gender. For example, TuSp, which was highly expressed in females, is
136 considered to be the most important component of egg cases¹⁹⁻²². In contrast, PiSp,
137 which was highly expressed in males, may contribute to the silk threads used to attach
138 to substrates^{23,24}. It is also of interest to note that both MaSp and MiSp were
139 abundantly expressed in male *P. pseudoannulata*. Further investigation of the
140 employment of spidroins in male spiders is of interest because, to date, spidroins have
141 been mainly studied in female spiders.

142 **Venom toxin**

143 Spider venom consists of numerous diverse components, however, in the present
144 study we have focused on the best studied neurotoxins. Spider neurotoxins are
145 peptides that contain 6-14 cysteine residues forming disulphide bridges and typically
146 comprise the inhibitor cysteine knot (ICK) motif²⁵. Thirty-two putative neurotoxin
147 precursor genes have been identified from the venom gland transcriptome and
148 genome analysis, forming six distinct families²⁶. The neurotoxin genes of similar
149 sequence often cluster on the same scaffold (Table S21). Other venom components
150 have also been annotated, including venom allergen 5, hyaluronidase, astacin-like
151 metalloprotease toxins, and Kunitz-type protease inhibitors (Table S22).

152 To date, at least 260 spider neurotoxins acting on ion channels have been
153 documented in the spider toxin database ArachnoServer 3.0^{27,28}. With the aim of
154 identifying the possible targets of *P. pseudoannulata* neurotoxins, we conducted a
155 phylogenetic analysis with 29 *P. pseudoannulata* neurotoxins and 48 neurotoxins with
156 known molecular targets from 14 other spider species (Fig. 3a). Further, *P.*
157 *pseudoannulata* neurotoxins are clustered into groups with neurotoxins that target
158 invertebrates only (15/29), vertebrates only (2/29) and both invertebrates and
159 vertebrates (7/29) according to the documented neurotoxins with known selectivity
160 (designated as reference neurotoxins in Fig. 3a). In addition, a single cluster was
161 identified with no similarity to known neurotoxins (5/29) (Fig. 3a). We expressed one
162 of these neurotoxins U1-lycotoxin-Pp1b *in vitro* and studied its toxicity²⁹. The
163 recombinant U1-lycotoxin-Pp1b did not show significant toxicity to mice at a dose of
164 10 mg/kg³⁰ but was found to be toxic to the insect *Nilaparvata lugens*, a typical prey
165 of *P. pseudoannulata*, with an LD_{50} (medium lethal dose) of 1.874 nmol/g insect
166 (13.70 mg/kg insect) (Fig. S5). Therefore, we classified the five unclustered
167 neurotoxins into the group ‘toxic to invertebrates’, based on the observed selectivity
168 of U1-lycotoxin-Pp1b (Fig. 3b, Table S23). While neurotoxins targeting invertebrates
169 account for 69% (20/29) of genes, their proportion in terms of the transcription level
170 is 93% (Fig. 3c, Table S23). Therefore, *P. pseudoannulata* neurotoxins have high
171 selectivity for insects, consistent with its role as a predator of agricultural insect pests.

172 Spiders of different body sizes capture prey ranging from small insects to large
173 rodents³¹. Neurotoxins from *P. pseudoannulata* and 10 other spider species were
174 categorized by their target species as being ‘invertebrate only’, ‘vertebrate only’ or
175 both. Information was retrieved for spider species that have at least 10 neurotoxins
176 with identified targets documented in ArachnoServer 3.0. A clear negative correlation

177 was observed between the percentage of neurotoxins targeting invertebrate only
178 (TX_inv.) and the prosoma length of a spider (liner regression, $F_{1,9} = 10.17$, $P = 0.011$,
179 Percentage of TX_inv. = $-2.96 \times \text{prosoma length} + 89.23$, $R^2 = 0.5304$, Fig. 3d, Fig. S6).
180 A correlation between the size of a spider and its prey is likely to be associated with
181 energy and nutritional requirement³². It seems likely that spider neurotoxins have
182 coevolved with spider body size with a shift in prey from smaller invertebrates to
183 larger vertebrates.

184 Toxins showed considerable diversity in the four species of spiders for which
185 genome sequence data are available (Fig. 3e, Table S24). The annotated toxins can be
186 grouped based on their structural domains. Neurotoxic Knottin toxins comprised
187 nearly half of the toxins in *P. pseudoannulata* and their gene number exceeded that in
188 the three web-weavers. The remarkable abundance of Knottin in *P. pseudoannulata* is
189 consistent with the fact that wandering spiders use venom as their main strategy of
190 predation and defence, whereas orb-weaver spiders rely more on low molecular mass
191 compounds and behavioural adaptations (such as prey-catching webs and sticky glue)
192³¹.

193 **Material and Methods**

194 **Genome sequencing**

195 **Sample preparation and sequencing**

196 For genome sequencing, three batches of *P. pseudoannulata* adults were collected
197 from spiders reared from two individual egg cases (batch #1 of 40 adults from egg
198 case #1; batch #2 of 60 1st-instar spiderlings from egg case #2; and batch #3 of 15 5th-
199 instar spiderlings from egg case #2). The two egg cases were derived from two
200 females collected at different time points from the same field in Jiangsu (λ
201 118.638551, ϕ 32.030345), China.

202 High quality genomic DNA was extracted using the conventional
203 phenol/chloroform extraction protocol³³ and broken into random fragments for
204 whole-genome shotgun sequencing. The genomic DNA was quality-examined with
205 agarose gel electrophoresis and quantified with Qubit™ system. Short-insert (250 bp
206 and 350 bp) paired-end libraries and large-insert (2 kb, 5 kb, 10 kb, 15 kb and 20 kb)
207 mate-pair libraries were prepared using the standard Illumina protocols. All libraries
208 were sequenced on the Illumina HiSeq 2000 platform with paired-end 150 bp and a
209 total of 1,306.48 Gb sequencing data were produced. To promote genome assembly,
210 the technologies of Pacific Bioscience's (PacBio's) single-molecule real-time (SMRT)
211 sequencing and 10x Genomics link-reads were also applied. For PacBio data,
212 SMRTbell libraries were prepared using 20-kb preparation protocols and sequenced
213 on PacBio Sequel platform, which generated 87.37 Gb (19.92x coverage) sequencing
214 data. The 10X Genomics linked-read library was constructed and sequenced on
215 Illumina Hiseq X Ten platform, which generated 465.21 Gb (106.09x coverage) raw
216 reads. For Illumina sequencing, raw data were filtered according to the following
217 criteria: reads containing adapter sequences; reads with $\geq 10\%$ unidentified
218 nucleotides (N); reads with low-quality bases (Q-value <5) more than 20%; duplicated
219 reads generated by PCR amplification during library construction. For PacBio
220 sequencing, subreads were filtered with the default parameters. All sequence data
221 were summarized in Table S1.

222 **Estimation of genome size**

223 Genome size was estimated by analysing the k-mer frequency. The distribution of k-
224 mer values depends on the genome characteristic and follows a Poisson distribution³⁴.
225 A total of 263 Gb high-quality short-insert reads (350 bp) were used to calculate the
226 17-mer frequency distribution and then estimate the *P. pseudoannulata* genome size

227 using the following formula ³⁵: genome size = (total number of 17-mers)/(position of
228 peak depth).

229 **Transcriptome preparation and sequencing**

230 A batch of adult spiders were collected randomly from a field in Jiangsu (λ
231 118.638551, ϕ 32.030345), China. Tissue samples were dissected from these adults as
232 four pairs of legs, pedipalp, chelicerae, brain, venom gland, fat body, male silk gland
233 and female silk gland. Total RNA for each sample was extracted with TRIZOL
234 Reagent (Thermo Fisher Scientific, Waltham, MA, USA). The concentration and
235 purity of the RNA sample was assessed by Nanodrop spectrophotometer (Thermo
236 Fisher Scientific, Waltham, MA, USA) and the integrity was checked by 2100
237 Bioanalyzer (Agilent, USA). RNA sequencing (RNA-seq) libraries were constructed
238 using the NEBNext® mRNA Library Prep Master Mix Set for Illumina® (New
239 England Biolab, RRID: SCR_013517) according to the manufacturer's instructions.
240 All libraries were sequenced on the Illumina Hiseq X Ten with paired-end 150 bp.
241 Information for all RNA-seq data was summarized in Table S3.

242 **Genome assembly**

243 For genome assembly of *P. pseudoannulata*, Platanus (PLATform for Assembling
244 Nucleotide Sequences, RRID: SCR_015531) (version 1.2.4) ³⁶ was first used to
245 construct the genome assembly backbone with all Illumina reads. Briefly, Platanus
246 carried out following three steps: (1) all short-insert paired-end reads were used to
247 construct *de Bruijn* graphs with automatically optimized k-mer sizes; (2) all short-
248 insert paired-end reads and large-insert mate-pair reads were aligned to the contigs for
249 scaffolding; (3) paired-end reads were aligned to scaffolds to close the gap. Then
250 GapCloser (RRID: SCR_015026) ³⁷ was used to fill the gaps in intra-
251 scaffold. Subsequently, the PacBio data were used to fill additional gaps with the

252 software PBJelly (RRID: SCR_012091) (version 1.3.1)³⁸ with default parameters.
253 After that, the resulting scaffolds were further connected to super-scaffolds using the
254 10X Genomics linked-reads by the software fragScaff (version 140324.1)³⁹
255 The completeness of the genome assembly and the uniformity of the sequencing were
256 evaluated with several approaches. Briefly, BWA (Burrows-Wheeler Aligner, RRID:
257 SCR_010910)⁴⁰ was used to align high-quality short-insert reads onto the *P.*
258 *pseudoannulata* genome with parameters of '-k 32 -w 10 -B 3 -O 11 -E 4'. Gene
259 region completeness was evaluated with the transcripts assembled by Trinity (version
260 2.1.1)⁴¹. All assembled transcript (length >= 200bp) were align onto the genome by
261 the software BLAT (BLAT, RRID: SCR_011919)³⁶ with default parameters.
262 CEGMA (Core Eukaryotic Genes Mapping Approach, RRID: SCR_015055)⁴² was
263 used to identify the exon-intron structures. Further, BUSCO (Benchmarking
264 Universal Single Copy Orthologs, RRID: SCR_015008) (v3.0.2)⁴³ was used to assess
265 the genome completeness with a set of 1066 arthropoda single-copy orthologous
266 genes.

267 **Genome annotation**

268 **Repetitive element identification**

269 Transposable elements (TEs) were identified with homology alignment and *de novo*
270 prediction. A *de novo* repeat library was built using RepeatModeler (RRID:
271 SCR_015027) (version 1.0.4)^{44,45}, RepeatScout (RRID: SCR_014653) (version 1.0.5)
272 ⁴⁶, and LTR_FINDER (RRID: SCR_015247) (version 1.06)⁴⁷ with default parameters.
273 Known TEs were identified via homology-based prediction using the RepeatMasker
274 (RRID: SCR_012954) (version 4.0.5)⁴⁴ with default parameters against the RepBase
275 library^{48,49}. In addition, tandem repeats were identified using Tandem Repeats Finder

276 (RRID: SCR_005659)^{50,51} with parameters "Match=2, Mismatch=7, Delta=7, PM=80,
277 PI=10, Minscore=50, MaxPeriod=2000".

278 **Protein-coding gene prediction**

279 Protein-coding genes were predicted with a combination of homology-based
280 prediction, *de novo* prediction, and transcriptome sequencing-based prediction
281 methods. For the homology-based gene prediction, protein sequences from four
282 species including *Stegodyphus mimosarum*, *Parasteatoda tepidariorum*, *Tetranychus*
283 *urticae* and *Drosophila serrata*⁵² were aligned to our assembled genome using
284 TBLASTN (RRID: SCR_011822)^{53,54} with e-value $\leq 1e-5$. The BLAST hits were
285 conjoined with the software Solar⁵⁵. Then GeneWise (RRID: SCR_015054) (version
286 2.2.0)^{56,57} was applied to predict gene models based on the alignment sequences. The
287 *de novo* prediction was performed using Augustus (RRID: SCR_008417) (version
288 3.0.2)^{58,59}, GeneScan (RRID: SCR_012902) (version 1.0)^{60,61}, GeneID (version 1.4)
289 ^{62,63}, GlimmerHMM (RRID: SCR_002654) (version 3.0.4)^{64,65} and SNAP (Semi-
290 HMM-based Nucleic acid Parser, RRID: SCR_007936)^{66,67} on the repeat-masked
291 genome. For the transcriptome-based prediction, RNA-Seq data from different tissues
292 including four pairs of legs, pedipalp, chelicerae, brain, venom gland and fat body
293 were aligned to the *P. pseudoannulata* genome using TopHat (RRID: SCR_013035)
294 (version 2.0.13)^{68,69} and gene structures were predicted with Cufflinks (RRID:
295 SCR_014597) (version 2.1.1)^{70,71}. In addition, the RNA-Seq data was assembled by
296 Trinity (RRID: SCR_013048) (version 2.1.1)^{41,72}. These assembled sequences were
297 aligned against our assembled genome by PASA (Program to Assemble Spliced
298 Alignment, RRID: SCR_014656)^{73,74} and generated gene models were used as the
299 training set for the softwares Augustus, GlimmerHMM and SNAP (Semi-HMM-
300 based Nucleic Acid Parser, RRID: SCR_002127)⁶⁷. Eventually, gene models

301 obtained from all the methods were integrated into a comprehensive and non-
302 redundant gene set with the software EVIDENCEModeler (EVM, RRID: SCR_014659)
303 ^{75,76}.

304 **Functional annotation**

305 To obtain functional annotation, all predicted protein-coding sequences in *P.*
306 *pseudoannulata* genome were aligned to public databases including National Center
307 for Biotechnology Information nonredundant protein (NR) ⁷⁷ and SwissProt ^{78,79}. The
308 known motifs and domains were annotated by searching InterPro databases ⁸⁰
309 including Pfam (RRID: SCR_004726) (version 27.0) ^{81,82}, PRINTS (RRID:
310 SCR_003412) (version 42.0) ^{83,84}, PROSITE (RRID: SCR_003457) (version 20.89)
311 ^{85,86}, ProDom (RRID: SCR_006969) (version 2006.1) ^{87,88}, SMART (RRID:
312 SCR_005026) (version 6.2) ^{89,90} and PANTHER (RRID: SCR_004869) (version 7.2)
313 ^{91,92} with the software InterProScan (RRID: SCR_005829) (version 4.7) ^{80,93}. Gene
314 Ontology (GO, RRID: SCR_002811) ^{94,95} terms for each gene were obtained from the
315 corresponding InterPro entry. Kyoto Encyclopedia of Genes and Genomes (KEGG,
316 RRID: SCR_012773) databases ^{96,97} were searched to identify the pathways in which
317 the genes might be involved.

318 **Phylogeny and divergence time estimation**

319 Gene family analysis was performed with 18 species including *Caenorhabditis*
320 *elegans*, *Drosophila melanogaster*, *Apis mellifera*, *Tribolium castaneum*, *Nilaparvata*
321 *lugens*, *Acyrtosiphon pisum*, *Bombyx mori*, *Hyalella Azteca*, *Daphnia magna*,
322 *Eurytemora affinis*, *Tetranychus urticae*, *Ixodes scapularis*, *Metaseiulus occidentalis*,
323 *Centruroides sculpturatus*, *Stegodyphus mimosarum*, *Parasteatoda tepidariorum*,
324 *Nephila clavipes*, and *P. pseudoannulata* (Table S15). Only the longest transcript of a
325 gene was retained as the representative if the gene had alternative splicing isoforms

326 identified. Genes with protein sequences shorter than 30 amino acids were removed.
327 Then, the similarities between genes in all selected genomes were identified using all-
328 versus-all BLASP with an E-value threshold of $1e-7$ and all the blast hits were
329 concatenated by the software Solar⁵⁵. Finally, gene families were constructed using
330 OrthoMCL^{98,99} with the setting of “-inflation 1.5”. In total, the protein-coding genes
331 were clustered into 29,995 gene families and 190 single-copy orthologs.

332 The phylogenetic relationship of *P. pseudoannulata* with the other 17 selected
333 species was analysed using the 190 single-copy gene families. Protein sequences of
334 the ortholog genes were aligned using the multiple alignment software MUSCLE with
335 default parameters¹⁰⁰. Then the alignments of each family were concatenated into a
336 super alignment matrix and RAxML (version 8.0.19)^{101,102} was used to reconstruct
337 the phylogenetic tree through maximum likelihood methods with default substitution
338 model-PROTGAMMAAUTO. Divergence times of these species were estimated
339 using the MCMCtree program in PAML^{103,104} with the parameters of ‘burn-in=10000,
340 sample-number=100,000 and sample-frequency=2’. Calibration points applied in
341 present study were obtained from the TimeTree database (*Drosophila melanogaster*,
342 *Bombyx mori*, *Tribolium castaneum*, *Apis mellifera*, 238~ 377 MYA; *Apis mellifera*,
343 *Tribolium castaneum*, *Bombyx mori*, *Drosophila melanogaster*, *Acyrtosiphon pisum*,
344 *Nilaparvata lugens*, 295~305 MYA; *Caenorhabditis elegans* and other species,
345 521~581 MYA).^{105,106}.

346 **Gene family contraction and expansion**

347 Expansion and contraction analysis of orthologous gene families was performed using
348 CAFÉ program (version 2.1) (Computational Analysis of gene Family Evolution,
349 RRID: SCR_005983)^{107,108}. The program uses a random birth and death model to
350 infer changes of gene families along each lineage of phylogenetic tree. Based on a

351 probabilistic graphical model, this method calculates a p-value for transitions between
352 parent and child nodes gene family size over a phylogeny. The gene families were
353 significantly expanded or contracted in the *P. pseudoannulata* genome with a p-value
354 of 0.05.

355 **Spidroin gene classification and quantification**

356 Multiple rounds of BLAST (RRID: SCR_004870) were run in the genome database to
357 identify putative spidroin genes. The spidroin genes in *N. clavipes*⁷ and *S.*
358 *mimosarum*³ were first used as queries for BLAST and the resultant *P.*
359 *pseudoannulata* spidroin genes were then added into the query repertoire to run
360 further BLAST searches⁷. The scaffolds containing the putative spidroin genes were
361 then subjected to Augustus (RRID: SCR_008417)¹⁰⁹ for gene prediction and the
362 predicted spidroin genes were manually checked for the presence of structural feature,
363 namely, the N-terminal domain, the repeat region and the C-terminal domain. The
364 manually checked spidroin genes were then examined by BLAST in the silk gland
365 transcriptomes of both males and females for their corresponding transcripts. Spidroin
366 genes were confirmed if at least one transcript aligned with 95% identity. The N-
367 terminal domain (130 amino acids) of the spidroins in *P. pseudoannulata* and *N.*
368 *clavipes* were aligned with ClustalW function and a phylogenetic tree was constructed
369 with maximum-likelihood method (1000 replicates) in MEGA (RRID: SCR_000667,
370 version 7)¹¹⁰ (Fig. S3, S4). Gene structures were drawn to scale in IBS (version 1.0.3)
371¹¹¹.

372 Silk glands were dissected and identified following protocols relating to the
373 western black widow spider^{15,16}. Images were obtained with a portable video
374 microscope (3R-MSA600, Anyty, 3R Eddyteck Corp., China) and contrasted with
375 Photoshop CS6.

376 Spider adults were collected from the paddy fields in Nanjing (Jiangsu, China), and
377 reared in laboratory conditions for at least two weeks. Spiders were anesthetized with
378 CO₂ and silk glands were carefully collected. The entire silk glands from 10 females
379 and 15 males were pooled as one sample, respectively. Three samples were prepared
380 for each gender. Each sample was kept in 200 µL RNAlater (Thermo Fisher
381 Scientific, Waltham, MA, USA) at -80 °C until total RNA extraction.

382 Total RNA was extracted from silk gland samples with GeneJET RNA Purification
383 Kit (Thermo Fisher Scientific, Waltham, MA, USA) after removing the RNAlater and
384 eluted with 44 µL nuclease-free water. Genomic DNA was removed with TURBO
385 DNA-free Kit (Thermo Fisher Scientific, RRID:SCR_008452) following the
386 manufacturer's instructions. The quality and quantity of the total RNAs were
387 monitored with NanoDrop spectrophotometer (Thermo Fisher Scientific) and 2%
388 agarose gel electrophoresis. RNA samples were stored at -80 °C. cDNA was
389 synthesized with 2 µg RNA using PrimeScript RT Reagent Kit (TaKaRa, Kyoto,
390 Japan) and then stored at -20 °C. Primers for quantitative real-time PCR (qPCR) were
391 designed using Beacon Designer (version 7.92, PREMIER Biosoft International, CA,
392 USA) (Table S19). Two pairs of universal primers were designed for MaSp
393 (MaSp_691565, MaSp_3359, MaSp_4789, MaSp_258724, MaSp_2831) and AcSp
394 (AcSp_1925.1, AcSp_1925.2), respectively, due to the high similarity of their
395 sequences. Glyceraldehyde-3-phosphatedehydrogenase (GAPDH) and elongation
396 factor 1-alpha (EF1 α) genes were selected as the reference genes for spidroin gene
397 quantification. The specificity and efficiency of the primers were validated via
398 standard curves with five serial cDNA dilutions and the melt curve with a temperature
399 range 60-95 °C. qPCR was performed using SYBR Premix Ex Taq Kit (TaKaRa,
400 Kyoto, Japan) following the manufacturer's instructions on a 7500 Real-Time PCR

401 System (Applied Biosystems, RRID:SCR_005039). Reagents were assembled in a 20
402 μL reaction containing 10 μL SYBR Premix Ex Taq, 6.8 μL sterile water, 0.4 μL
403 forward primer, 0.4 μL reverse primer, 0.4 μL ROX Reference Dye II and 2 μL
404 cDNA. The reaction program was 95 °C for 30 sec, 40 cycles of 95 °C for 5 sec and
405 60 °C for 34 sec. No template control (NTC) and no reverse transcriptase control
406 (NRT) were included as negative controls to eliminate the possibilities of reagent
407 contamination and genomic DNA contamination. Each reaction was performed in two
408 technical replicates and three biological samples were tested. Ct values of qPCR were
409 exported from 7500 Real-Time PCR Software (RRID:SCR_014596) (version 2.0.6).
410 The expression levels of target genes were relative to the geometric mean of two
411 reference genes¹¹² following the $2^{-\Delta\text{CT}}$ method¹¹³. Statistical analyses were performed
412 with GraphPad Prism (RRID: SCR_002798) (version 7).

413 **Neurotoxin identification and bioassay**

414 **Identification of neurotoxin genes and comparative analysis among spiders**

415 Neurotoxin candidates were retrieved via BLAST in the genome database with
416 neurotoxins from ArachnoServer 3.0 as queries^{27,28}. They were identified as
417 neurotoxins when the proteins met the criteria such as containing 6-14 cysteine
418 residues and the canonical neurotoxin domains. The neurotoxins were characterized
419 with their signal peptide via SignalP 4.1 Server^{114,115}, propeptide, and the Cys-Cys
420 disulfide bridge pattern. Inhibitor Cystine Knots (ICK) were predicted on the
421 KNOTTIN database¹¹⁶. Toxin genes subjected to interspecific comparison were
422 retrieved from the genome annotation with astacin-like metalloprotease toxin
423 excluded. Non-Knottin toxins were subjected to NCBI domain analysis and grouped
424 accordingly, putative neurotoxins with Spider_toxin or toxin_35 domain, cysteine
425 protease inhibitors containing TY (thyroglobulin type I repeats, accession no.

426 cd00191), serine protease inhibitors containing KU (BPTI/Kunitz family of serine
427 protease inhibitors, accession no. cd00109), Trypsin-like serine proteases with
428 Tryp_SPc (Trypsin-like serine protease, accession no. cd00190), SVWC family
429 proteins containing SVWC (single-domain von Willebrand factor type C proteins,
430 accession no. pfam15430), colipases with COLIPASE (Colipases, accession no.
431 smart00023) and the rest designated as other toxins.

432 **Construction of neurotoxin expression vector**

433 The open reading frame of neurotoxin U1-lycotoxin-Pp1b was cloned into the
434 prokaryotic expression vector pLicC-MBP-APETx2^{29,117}. *Kpn* I and *Ava* I were used
435 for double digestion, and primer sequences were:
436 cggggtaccccgaaatctgtatttcagggaaggcatgcacccaaggtttac (forward) and
437 ccctcagggttaaccgaatagagtcttaattctgcc (reverse). For the PCR amplification, the high-
438 fidelity PrimerSTAR (TaKaRa, Tokyo, Japan) was used, and the amplification
439 program was 94°C for 3 minutes, followed by 30 cycles of 94°C for 30 seconds, 55°C
440 for 30 seconds, 72°C for 30 seconds, and finally 72°C for 10 minutes. The PCR
441 products were gel-purified using a Gel Extraction Kit (CW BIO, Nanjing, China),
442 ligated into sequencing vector, and sequenced at Genscript Biotechnology Co. Ltd.
443 (Nanjing, China). Plasmid extraction was performed using Mini Plasmid Extraction
444 Kit I (OMEGA, Guangzhou, China), and double digestion using *Kpn* I (TaKaRa,
445 Tokyo, Japan) and *Ava* I (TaKaRa, Tokyo, Japan). After the target gene and vector
446 were recovered, they were ligated with T4 DNA ligase (TaKaRa, Tokyo, Japan) at
447 4°C overnight. Subsequently, the ligation product was transformed into *Escherichia*
448 *coli* BL21 strain, and sequenced at Genscript Biotechnology Co. Ltd. (Nanjing,
449 China). A positive clone was selected and the final concentration of 40% glycerol was
450 added for preservation at -80°C.

451 **Expression and purification of the recombinant neurotoxin**

452 Positive clones were cultivated in LB liquid medium containing 100 mg/L
453 carbenicillin at 37°C in a constant temperature shaker with 250 r/min for 14h. The
454 culture broth was inoculated in an LB liquid medium containing 100 mg/L
455 carbenicillin at a ratio of 1:100 and cultured with 250 r/min at 37°C for 4h. IPTG was
456 added into the culture at a final concentration of 0.5 mmol/L, 160r/min, and 25°C to
457 induce the expression for 4h. Bacteria were collected after centrifuging at 4000rpm
458 for 10min at 4°C and resuspended with 20mM Tris-HCl (pH 7.4). Then, 100µg/ml
459 lysozyme, 0.1% Triton X-100, 0.5mM PMSF were added to the suspension and the
460 digestion was performed at 37°C for 1 hour. Bacteria were then ultrasonicated and
461 debris were removed by centrifugation with 16,000 rpm for 20min at 4°C, and the
462 supernatant was filtered through a 0.22 µm filter.

463 The recombinant toxin was purified using AKTA Avant automated protein
464 purification system (GE, Uppsala, Sweden). Initially, the fusion protein was collected
465 by affinity chromatography using nickel column HisTrap HP (GE, Uppsala, Sweden),
466 and then transferred to TEV protease buffer using desalting column HiPrep 26/10
467 Desalting (GE, Uppsala, Sweden) and digested with TEV protease (Solarbio, Beijing,
468 China) at 16°C for 10h. Next, TEV protease buffer was replaced by PBS using the
469 desalting column. Finally, the protein label was removed by affinity chromatography
470 using nickel column, and the purity of recombinant toxin was detected by SDS-PAGE.

471 **Biological activity assay of recombinant neurotoxin**

472 The insecticidal activity of the recombinant neurotoxin against 5th-instar *Nilaparvata*
473 *lugens* nymphs was determined by microinjection¹¹⁸. *N. lugens* nymphs were injected
474 with neurotoxin solutions at a series of concentrations with 3 replicates of 20
475 individuals per replicate. PBS was injected as control solution. Before injection, the

476 test insects were anesthetized with CO₂ and each insect was injected with 30nl of
477 recombinant neurotoxin. After injection, the test insects were checked 1h or 12h later.
478 The toxicity of the recombinant neurotoxin against mice was measured by lateral
479 ventricle injection ¹¹⁹.

480 **Discussions**

481 The complete genome sequencing of the wandering spider *P. pseudoannulata*
482 provides valuable insights into the aspects of the biology of wandering spiders,
483 including diversity of spidroin genes and invertebrate-specific neurotoxins which may
484 have potential importance in developing novel pest management strategies.

485 The evolutionary diversification of spiders has long been discussed. The focus of
486 the debate has been whether the orb web origin is monophyletic or polyphyletic ^{1,2,120}.
487 Notably, the cursorial, non-web building spider taxa from RTA clade has proven to be
488 more important than previously thought in the phylogenetic analysis with
489 morphological, behavioural and molecular evidences ^{2,120}. In addition to the
490 distinction in silk use, wandering spiders differs from the web weaver spiders in
491 habitat and biotic interaction, which can also promote the diversification ². Therefore,
492 the massive genomic information of *P. pseudoannulata* offers added value to the
493 diversification analysis. Spiders of different ecological niches have evolved their
494 corresponding behaviours and lifestyle. Genome comparison analysis of spiders of
495 different lifestyles or habitats will reveal primary hints for molecular mechanisms of
496 spiders' evolutionary adaptation.

497 The present studies were prompted, in part, an interest in the genetic basis for
498 differences in methods of predation by spiders and, in particular, in the use of
499 neurotoxins to incapacitate insect prey. Although further work will be required to
500 understand the molecular targets and mode of actions of *P. psuedoannulata*

501 neurotoxins, it is possible that a better understanding of spider toxins may lead to the
502 development of novel pest control agents applicable to integrated pest management
503 and other potential biomedical applications.

504 **Availability of supporting data and materials**

505 The *P. pseudoannulata* genome has been deposited in GenBank under accession No.
506 SBLA00000000 and transcriptomes have been deposited to NCBI Sequence Read
507 Archive under accession No. SRR8083387-SRR8083398.

508 **Additional files**

509 **Table S1-S3.** Statistics of genome and transcriptome sequencing

510 **Table S4-S6.** Genome assembly

511 **Table S7-S9.** Genome evaluation

512 **Table S10-S14.** Genome annotation

513 **Table S15-S17.** Comparative genomes, gene expansion and contraction

514 **Table S18-S20.** Spidroin gene classification and quantification

515 **Table S21-S24.** Venom components and neurotoxins

516 **Figure S1.** Venndiagram of gene sets obtained using three prediction methods (*de*
517 *novo*, homology-based, RNAseq-based).

518 **Figure S2.** Estimation of divergence time.

519 **Figure S3.** Phylogenetic tree of 16 *P. pseudoannulata* spidroins.

520 **Figure S4.** Phylogenetic tree of 16 *P. pseudoannulata* spidroins and 26 complete *N.*
521 *clavipes* spidrons.

522 **Figure S5.** Toxicity assay of the recombinant neurotoxin U1-lycotoxin-Pp1.

523 **Figure S6.** Linear regression analysis of spider prosoma length and the number of
524 neurotoxins targeting invertebrate prey species.

525 **Abbreviations**

526 Ac: aciniform; bp: base pair; BUSCO: benchmarking universal single-copy orthologs;
527 BWA: Burrows-Wheeler Alignment tool; CDS: coding sequence; CEGMA: core
528 eukaryotic genes mapping approach; EF1 α : elongation factor 1-alpha; EST: expressed
529 sequence tag; FDR: false discovery rate; FPKM: Fragments Per Kilobase of exon
530 model per Million mapped fragments; GAPDH: glyceraldehyde-3-
531 phosphatedehydrogenase ; Gb: gigabases; GO: Gene Ontology ; KEGG: Kyoto
532 Encyclopedia of Genes and Genomes ; Ma: major ampullate ; Mi: minor ampullate ;
533 MYA: million years ago; NJ: neighbour joining; NRT: no reverse transcriptase
534 control; NTC: no template control; Pi: piriform; RTA: retrolateral tibial apophysis; Sp:
535 spidroin; TE: transposable element; Tu: tubuliform.

536 **Animal care**

537 Animal experimental procedures were approved by the Laboratory Animal Ethical
538 Committee of Nanjing Agricultural University (No. PZ2019021) and performed
539 accordingly.

540 **Competing interests**

541 The authors declare that they have no competing interests.

542 **Authors' contributions**

543 Z. Liu initiated and supervised the project. H. Bao and Y. Yang contributed to the
544 sample collection and handling. M. Liu performed genome sequencing, assembly and
545 primary annotation. M. Liu, J. Li and H. Gao performed the comparative genomic
546 analysis. N. Yu, Y. Zhang, H. Bao, T. Van Leeuwen, N. S. Millar and Z. Liu
547 contributed to the data mining and analysis. N. Yu and Z. Yang conducted the
548 analysis and experimental validation of spidroin genes. L. Huang and Z. Wang

549 performed the analysis and experimental validation of neurotoxins. N. Yu and J. Li
550 submitted data to NCBI. N. Yu and Z. Liu wrote the initial draft of the manuscript. N.
551 Yu, T. Van Leeuwen, N. Millar and Z. Liu revised the manuscript. All authors have
552 read and approved the final manuscript.

553 These authors contributed equally: Na Yu, Jingjing Li and Meng Liu.

554 **Current address:**

555 Yuanxue Yang, Cotton Research Center, Shandong Academy of Agricultural
556 Sciences, Jinan 250100, China

557 **Acknowledgments**

558 We thank Dr. Huixing Lin (MOE joint international research laboratory of animal
559 health and food safety, Nanjing Agricultural University) for the neurotoxin assay with
560 mice. The work was supported by National Natural Science Foundation of China
561 (grant number 31772185, 31601656, 31701823).

562 **References**

- 563 1. Blackledge, T.A. *et al.* Reconstructing web evolution and spider diversification in
564 the molecular era. *Proceedings of the National Academy of Sciences of the United*
565 *States of America* **106**, 5229-5234 (2009).
- 566 2. Fernandez, R. *et al.* Phylogenomics, diversification dynamics, and comparative
567 transcriptomics across the spider tree of life. *Current Biology* **28**, 2190-2193
568 (2018).
- 569 3. Sanggaard, K.W. *et al.* Spider genomes provide insight into composition and
570 evolution of venom and silk. *Nature Communications* **5** (2014).
- 571 4. Liu, S., Aagaard, A., Bechsgaard, J. & Bilde, T. DNA methylation patterns in the
572 social spider, *Stegodyphus dumicola*. *Genes* **10** (2019).
- 573 5. Gendreau, K.L. *et al.* House spider genome uncovers evolutionary shifts in the
574 diversity and expression of black widow venom proteins associated with extreme
575 toxicity. *BMC Genomics* **18**, 178 (2017).

- 576 6. Schwager, E.E. *et al.* The house spider genome reveals an ancient whole-genome
577 duplication during arachnid evolution. *BMC Biology* **15** (2017).
- 578 7. Babb, P.L. *et al.* The *Nephila clavipes* genome highlights the diversity of spider
579 silk genes and their complex expression. *Nature Genetics* **49**, 895-903 (2017).
- 580 8. Regier, J.C. *et al.* Arthropod relationships revealed by phylogenomic analysis of
581 nuclear protein-coding sequences. *Nature* **463**, 1079-1083 (2010).
- 582 9. Chen, S. *et al.* *De novo* analysis of transcriptome dynamics in the migratory
583 locust during the development of phase traits. *Plos One* **5** (2010).
- 584 10. Adrianos, S.L. *et al.* *Nephila clavipes* flagelliform silk-like GGX motifs
585 contribute to extensibility and spacer motifs contribute to strength in synthetic
586 spider silk fibers. *Biomacromolecules* **14**, 1751-1760 (2013).
- 587 11. Hayashi, C.Y. & Lewis, R.V. Molecular architecture and evolution of a modular
588 spider silk protein gene. *Science* **287**, 1477-1479 (2000).
- 589 12. Chores, O., Bayarmagnai, B. & Lewis, R.V. Spider web glue: two proteins
590 expressed from opposite strands of the same DNA sequence. *Biomacromolecules*
591 **10**, 2852-2856 (2009).
- 592 13. Jain, D., Amarpuri, G., Fitch, J., Blackledge, T.A. & Dhinojwala, A. Role of
593 hygroscopic low molecular mass compounds in humidity responsive adhesion of
594 spider's capture silk. *Biomacromolecules* **19**, 3048-3057 (2018).
- 595 14. Singla, S., Amarpuri, G., Dhopatkar, N., Blackledge, T.A. & Dhinojwala, A.
596 hygroscopic compounds in spider aggregate glue remove interfacial water to
597 maintain adhesion in humid conditions. *Nature Communications* **9**, 3048-3057
598 (2018).
- 599 15. Jeffery, F. *et al.* Microdissection of black widow spider silk-producing glands.
600 *Journal of Visualized Experiments* (2011).
- 601 16. Chaw, R.C. & Hayashi, C.Y. Dissection of silk glands in the western black
602 widow *Latrodectus hesperus*. *Journal of Arachnology* **46**, 159-161 (2018).
- 603 17. Vollrath, F. & Knight, D.P. Liquid crystalline spinning of spider silk. *Nature* **410**,
604 541-548 (2001).
- 605 18. Chaw, R.C. *et al.* Intragenic homogenization and multiple copies of prey-
606 wrapping silk genes in *Argiope* garden spiders. *Bmc Evolutionary Biology*
607 **14**(2014).

- 608 19. Casem, M.L., Collin, M.A., Ayoub, N.A. & Hayashi, C.Y. Silk gene transcripts
609 in the developing tubuliform glands of the Western black widow, *Latrodectus*
610 *hesperus*. *Journal of Arachnology* **38**, 99-103 (2010).
- 611 20. Garb, J.E. & Hayashi, C.Y. Modular evolution of egg case silk genes across orb-
612 weaving spider superfamilies. *Proceedings of the National Academy of Sciences*
613 *of the United States of America* **102**, 11379-11384 (2005).
- 614 21. Hu, X.Y. *et al.* Spider egg case core fibers: Trimeric complexes assembled from
615 TuSp1, ECP-1, and ECP-2. *Biochemistry* **45**, 3506-3516 (2006).
- 616 22. Jiang, P. *et al.* Structure, composition and mechanical properties of the silk fibres
617 of the egg case of the Joro spider, *Nephila clavata* (Araneae, Nephilidae). *Journal*
618 *of Biosciences* **36**, 897-910 (2011).
- 619 23. Blasingame, E. *et al.* Pyriform spidroin 1, a novel member of the silk gene family
620 that anchors dragline silk fibers in attachment discs of the black widow spider,
621 *Latrodectus hesperus*. *Journal of Biological Chemistry* **284**, 29097-29108 (2009).
- 622 24. Geurts, P. *et al.* Synthetic spider silk fibers spun from pyriform spidroin 2, a glue
623 silk protein discovered in orb-weaving spider attachment discs.
624 *Biomacromolecules* **11**, 3495-3503 (2010).
- 625 25. Norton, R.S. & Pallaghy, P.K. The cystine knot structure of ion channel toxins
626 and related polypeptides. *Toxicon* **36**, 1573-1583 (1998).
- 627 26. Huang, L., Wang, Z., Yu, N., Li, J. & Liu, Z. Toxin diversity revealed by the
628 venom gland transcriptome of *Pardosa pseudoannulata*, a natural enemy of
629 several insect pests. *Comparative Biochemistry and Physiology D-Genomics &*
630 *Proteomics* **28**, 172-182 (2018).
- 631 27. Pineda, S.S. *et al.* ArachnoServer 3.0: an online resource for automated discovery,
632 analysis and annotation of spider toxins. *Bioinformatics* **34**, 1074-1076 (2018).
- 633 28. ArachnoServer spider toxin database.
634 <http://www.arachnoserver.org/mainMenu.html>.
- 635 29. Anangi, R., Rash, L.D., Mobli, M. & King, G.F. Functional expression in
636 *Escherichia coli* of the disulfide-rich sea anemone peptide APETx2, a potent
637 blocker of acid-sensing ion channel 3. *Marine Drugs* **10**, 1605-1618 (2012).
- 638 30. de Lima, M.E. *et al.* The toxin Tx4(6-1) from the spider *Phoneutria nigriventer*
639 slows down Na⁺ current inactivation in insect CNS via binding to receptor site 3.
640 *Journal of Insect Physiology* **48**, 53-61 (2002).

- 641 31. Kuhn-Nentwig, L., Stoecklin, R. & Nentwig, W. Venom composition and
642 strategies in spiders: Is everything possible? *Advances in Insect Physiology* **40**, 1-
643 86 (2011).
- 644 32. Nentwig, W. & Wissel, C. A comparison of prey lengths among spiders.
645 *Oecologia* **68**, 595-600 (1986).
- 646 33. Sambrook, J., Russell, D.W., Sambrook, J. & Russell, D.W. *Molecular cloning:*
647 *A laboratory manual*, (Cold Spring Harbor Laboratory Press , 10 Skyline Drive,
648 Plainview, NY, 11803-2500, USA, 2001).
- 649 34. Li, R. *et al.* The sequence and de novo assembly of the giant panda genome.
650 *Nature* **463**, 1106-1106 (2010).
- 651 35. Liu, B. *et al.* Estimation of genomic characteristics by analyzing k-mer frequency
652 in *de novo* genome projects. *Quantitative Biology* **35**, 62-67 (2013).
- 653 36. Kajitani, R. *et al.* Efficient de novo assembly of highly heterozygous genomes
654 from whole-genome shotgun short reads. *Genome Research* **24**, 1384-1395
655 (2014).
- 656 37. English, A.C. *et al.* Mind the gap: Upgrading genomes with Pacific Biosciences
657 RS long-read sequencing technology. *Plos One* **7**(2012).
- 658 38. English, A.C., Salerno, W.J. & Reid, J.G. PBHoney: identifying genomic variants
659 via long-read discordance and interrupted mapping. *Bmc Bioinformatics* **15**(2014).
- 660 39. Adey, A. *et al.* In vitro, long-range sequence information for de novo genome
661 assembly via transposase contiguity. *Genome Research* **24**, 2041-2049 (2014).
- 662 40. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-
663 Wheeler transform. *Bioinformatics* **25**, 1754-1760 (2009).
- 664 41. Grabherr, M.G. *et al.* Full-length transcriptome assembly from RNA-Seq data
665 without a reference genome. *Nature Biotechnology* **29**, 644-652 (2011).
- 666 42. Parra, G., Bradnam, K. & Korf, I. CEGMA: a pipeline to accurately annotate core
667 genes in eukaryotic genomes. *Bioinformatics* **23**, 1061-1067 (2007).
- 668 43. Simao, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. & Zdobnov,
669 E.M. BUSCO: assessing genome assembly and annotation completeness with
670 single-copy orthologs. *Bioinformatics* **31**, 3210-3212 (2015).
- 671 44. Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive
672 elements in genomic sequences. *Current protocols in bioinformatics* **Chapter 4**,
673 Unit 4.10-Unit 4.10 (2009).
- 674 45. RepeatModeler. <http://www.repeatmasker.org/RepeatModeler/>.

- 675 46. Price, A.L., Jones, N.C. & Pevzner, P.A. De novo identification of repeat families
676 in large genomes. *Bioinformatics* **21**, I351-I358 (2005).
- 677 47. Xu, Z. & Wang, H. LTR_FINDER: an efficient tool for the prediction of full-
678 length LTR retrotransposons. *Nucleic Acids Research* **35**, W265-W268 (2007).
- 679 48. Bao, W., Kojima, K.K. & Kohany, O. Repbase Update, a database of repetitive
680 elements in eukaryotic genomes. *Mobile DNA* **6**(2015).
- 681 49. RepBase. <https://www.girinst.org/server/RepBase/index.php>.
- 682 50. Benson, G. Tandem repeats finder: a program to analyze DNA sequences.
683 *Nucleic Acids Research* **27**, 573-580 (1999).
- 684 51. Tandem Repeats Finder. <http://tandem.bu.edu/trf/trf.download.html>.
- 685 52. Wang, Z. *et al.* The draft genomes of soft-shell turtle and green sea turtle yield
686 insights into the development and evolution of the turtle-specific body plan.
687 *Nature Genetics* **45**, 701-706 (2013).
- 688 53. Altschul, S.F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of
689 protein database search programs. *Nucleic Acids Research* **25**, 3389-3402 (1997).
- 690 54. TBLASTN. <https://blast.ncbi.nlm.nih.gov/Blast.cgi>.
- 691 55. Darriba, D., Taboada, G.L., Doallo, R. & Posada, D. jModelTest 2: more models,
692 new heuristics and parallel computing. *Nature Methods* **9**, 772-772 (2012).
- 693 56. Birney, E., Clamp, M. & Durbin, R. GeneWise and genomewise. *Genome*
694 *Research* **14**, 988-995 (2004).
- 695 57. GeneWise. <https://www.ebi.ac.uk/~birney/wise2/>.
- 696 58. Stanke, M., Schoffmann, O., Morgenstern, B. & Waack, S. Gene prediction in
697 eukaryotes with a generalized hidden Markov model that uses hints from external
698 sources. *Bmc Bioinformatics* **7**(2006).
- 699 59. Stanke, M. & Waack, S. Gene prediction with a hidden Markov model and a new
700 intron submodel. *Bioinformatics* **19**, II215-II225 (2003).
- 701 60. Salamov, A.A. & Solovyev, V.V. Ab initio gene finding in Drosophila genomic
702 DNA. *Genome Research* **10**, 516-522 (2000).
- 703 61. GeneScan. <http://genome.dkfz-heidelberg.de/cgi-bin/GENSCAN/genscan.cgi>.
- 704 62. Parra, G., Blanco, E. & Guigo, R. GeneID in Drosophila. *Genome Research* **10**,
705 511-515 (2000).
- 706 63. GeneID. <http://genome.crg.es/software/geneid/index.html>.

- 707 64. Majoros, W.H., Pertea, M. & Salzberg, S.L. TigrScan and GlimmerHMM: two
708 open source ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878-2879
709 (2004).
- 710 65. GlimmerHMM. <http://ccb.jhu.edu/software/glimmerhmm/>.
- 711 66. Korf, I. Gene finding in novel genomes. *BMC Bioinformatics* **5**(2004).
- 712 67. Semi-HMM-based Nucleic Acid Parser. <http://korflab.ucdavis.edu/software.html>.
- 713 68. Trapnell, C., Pachter, L. & Salzberg, S.L. TopHat: discovering splice junctions
714 with RNA-Seq. *Bioinformatics* **25**, 1105-1111 (2009).
- 715 69. TopHat. <https://ccb.jhu.edu/software/tophat/index.shtml>.
- 716 70. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals
717 unannotated transcripts and isoform switching during cell differentiation. *Nature*
718 *Biotechnology* **28**, 511-515 (2010).
- 719 71. Cufflinks. <http://cole-trapnell-lab.github.io/cufflinks/>.
- 720 72. Trinity. <https://github.com/trinityrnaseq/trinityrnaseq/wiki>.
- 721 73. Haas, B.J. *et al.* Improving the Arabidopsis genome annotation using maximal
722 transcript alignment assemblies. *Nucleic Acids Research* **31**, 5654-5666 (2003).
- 723 74. Program to Assemble Spliced Alignments.
724 <https://github.com/PASApipeline/PASApipeline/wiki>.
- 725 75. Haas, B.J. *et al.* Automated eukaryotic gene structure annotation using
726 EVIDENCEModeler and the program to assemble spliced alignments. *Genome*
727 *Biology* **9** (2008).
- 728 76. EVIDENCEModeler. <https://github.com/EVIDENCEModeler/EVIDENCEModeler/>.
- 729 77. National Center for Biotechnology Information.
730 <https://www.ncbi.nlm.nih.gov/protein/>.
- 731 78. Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its
732 supplement TrEMBL in 2000. *Nucleic Acids Research* **28**, 45-48 (2000).
- 733 79. UniProtKB. <http://www.uniprot.org/uniprot/>.
- 734 80. Mulder, N. & Apweiler, R. InterPro and InterProScan: tools for protein sequence
735 classification and comparison. in *Comparative Genomics* (ed. Bergman, N.H.)
736 59-70 (Humana Press, Totowa, NJ, 2007).
- 737 81. Bateman, A. *et al.* The Pfam protein families database. *Nucleic Acids Research*
738 **28**, 263-266 (2000).
- 739 82. Pfam. <http://pfam.xfam.org/>.

- 740 83. Attwood, T.K., Beck, M.E., Bleasby, A.J. & Parrysmith, D.J. PRINTS - a
741 database of protein motif fingerprints. *Nucleic Acids Research* **22**, 3590-3596
742 (1994).
- 743 84. PRINTS. <http://www.bioinf.man.ac.uk/dbbrowser/PRINTS/>
- 744 85. Hulo, N. *et al.* The PROSITE database. *Nucleic Acids Research* **34**, D227-D230
745 (2006).
- 746 86. PROSITE Database of protein domains, families and functional sites.
747 <http://www.expasy.ch/prosite/>.
- 748 87. Bru, C. *et al.* The ProDom database of protein domain families: more emphasis
749 on 3D. *Nucleic Acids Research* **33**, D212-D215 (2005).
- 750 88. ProDom. <http://prodom.prabi.fr/>.
- 751 89. Letunic, I., Doerks, T. & Bork, P. SMART 7: recent updates to the protein
752 domain annotation resource. *Nucleic Acids Research* **40**, D302-D305 (2012).
- 753 90. Simple Modular Architecture Research Tool. <http://smart.embl-heidelberg.de/>.
- 754 91. Mi, H.Y. *et al.* The PANTHER database of protein families, subfamilies,
755 functions and pathways. *Nucleic Acids Research* **33**, D284-D288 (2005).
- 756 92. PANTHER Classification System. <http://www.pantherdb.org/>.
- 757 93. InterPro: protein sequence analysis & classification.
758 <http://www.ebi.ac.uk/interpro/>.
- 759 94. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nature*
760 *Genetics* **25**, 25-29 (2000).
- 761 95. The Gene Ontology Resource. <http://geneontology.org/>.
- 762 96. Ogata, H. *et al.* KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic*
763 *Acids Research* **27**, 29-34 (1999).
- 764 97. KEGG: Kyoto Encyclopedia of Genes and Genomes.
765 <https://www.genome.jp/kegg/>.
- 766 98. Li, L., Stoeckert, C.J. & Roos, D.S. OrthoMCL: Identification of ortholog groups
767 for eukaryotic genomes. *Genome Research* **13**, 2178-2189 (2003).
- 768 99. OrthoMCL. <https://orthomcl.org/orthomcl/>.
- 769 100. Edgar, R.C. MUSCLE: multiple sequence alignment with high accuracy and high
770 throughput. *Nucleic Acids Research* **32**, 1792-1797 (2004).
- 771 101. Stamatakis, A. RAxML-VI-HPC: Maximum likelihood-based phylogenetic
772 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688-2690
773 (2006).

- 774 102.Stamatakis, A., Hoover, P. & Rougemont, J. A rapid bootstrap algorithm for the
775 RAxML web servers. *Systematic Biology* **57**, 758-771 (2008).
- 776 103.Yang, Z. PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular*
777 *Biology and Evolution* **24**, 1586-1591 (2007).
- 778 104.Phylogenetic Analysis by Maximum Likelihood. Vol. 2017.
- 779 105.Hedges, S.B., Marin, J., Suleski, M., Paymer, M. & Kumar, S. Tree of life reveals
780 clock-like speciation and diversification. *Molecular Biology and Evolution* **32**,
781 835-845 (2015).
- 782 106.TimeTree. <http://www.timetree.org>.
- 783 107.De Bie, T., Cristianini, N., Demuth, J.P. & Hahn, M.W. CAFE: a computational
784 tool for the study of gene family evolution. *Bioinformatics* **22**, 1269-1271 (2006).
- 785 108.Computational Analysis of gene Family Evolution.
786 <https://sourceforge.net/projects/cafehahnlab/>.
- 787 109.Augustus web interface. <http://bioinf.uni-greifswald.de/augustus/submission.php>.
- 788 110.Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular evolutionary genetics
789 analysis version 7.0 for bigger datasets. *Molecular Biology and Evolution* **33**,
790 1870-1874 (2016).
- 791 111.Liu, W. *et al.* IBS: an illustrator for the presentation and visualization of
792 biological sequences. *Bioinformatics* **31**, 3359-3361 (2015).
- 793 112.Vandesompele, J. *et al.* Accurate normalization of real-time quantitative RT-PCR
794 data by geometric averaging of multiple internal control genes. *Genome biology* **3**,
795 research0034 (2002).
- 796 113.Livak, K.J. & Schmittgen, T.D. Analysis of relative gene expression data using
797 real-time quantitative PCR and the 2(T)(-Delta Delta C) method. *Methods* **25**,
798 402-408 (2001).
- 799 114.Nielsen, H. Predicting Secretory Proteins with SignalP. in *Protein Function*
800 *Prediction: Methods and Protocols*, Vol. 1611 (ed. Kihara, D.) 59-73 (2017).
- 801 115.SignalP 4.1 Server. www.cbs.dtu.dk/services/SignalP-4.1/.
- 802 116.KNOTTIN database. <http://www.dsimb.inserm.fr/KNOTTIN/>.
- 803 117.Addgene. <http://www.addgene.org/>.
- 804 118.Liu, S., Ding, Z., Zhang, C., Yang, B. & Liu, Z. Gene knockdown by intro-
805 thoracic injection of double-stranded RNA in the brown planthopper, Nilaparvata
806 lugens. *Insect Biochemistry and Molecular Biology* **40**, 666-671 (2010).

807 119.Garcia, E., Rios, C. & Sotelo, J. Ventricular injection of nerve growth-factor
808 increase dopamine content in the striata of MPTP-treated mice. *Neurochemical*
809 *Research* **17**, 979-982 (1992).

810 120.Coddington, J.A., Agnarsson, I., Hamilton, C.A. & Bond, J.E. Spiders did not
811 repeatedly gain, but repeatedly lost, foraging webs. *Peerj* **7** (2019).

812 **Figure legend**

813 **Fig. 1. *P. pseudoannulata* and gene gain-and-loss analysis among species of**
814 **Arthropoda. (a)** A *P. pseudoannulata* catching prey. **(b)** The divergence time was
815 estimated by PAML mcmctree and is marked with a scale in million years. All
816 internal branches of the tree are 100% bootstrap supported. The numbers next to the
817 branch represent the numbers of expanded (green) and contracted (red) gene families
818 since the split from a most recent common ancestor.

819 **Fig. 2. Characterization of spidroin genes in *P. pseudoannulata*.** **(a)** Phylogeny
820 and gene structure characteristics of the *P. pseudoannulata* spidroins. The phylogenetic
821 tree was constructed with the 130 N-terminal amino acid residues from each putative
822 gene product and transformed cladogram. Gene structures are drawn to scale. The
823 alternated grey and light grey blocks represent only the repeat regions without any
824 sequence preference. Only the N-terminal region and partial repeat region are
825 available for gene MaSp_16938. **(b)** The anatomy of the four types of silk glands in *P.*
826 *pseudoannulata*. Major ampullate gland (Ma), minor ampullate gland (Mi), aciniform
827 gland (Ac) and piriform gland (Pi) were dissected from a female spider. **(c)** The
828 relative expression of each spidroin gene type in female and male *via* qPCR.

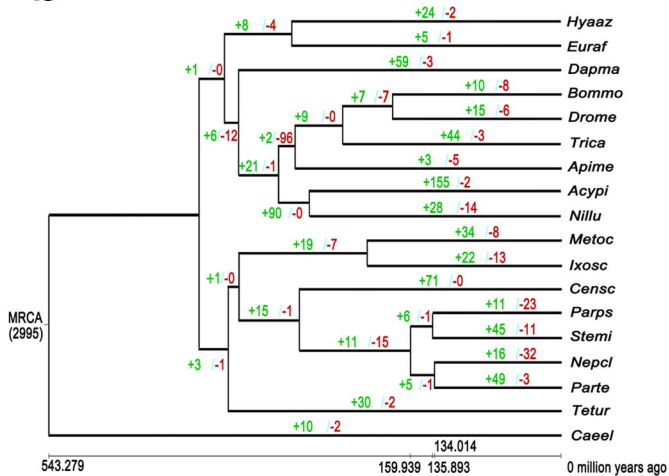
829 **Fig. 3. Evolutionary analysis of neurotoxins from *P. pseudoannulata* and other**
830 **spiders.** Neurotoxins are categorized according to their target organism as
831 invertebrate only (red), vertebrate only (green) and both invertebrate and vertebrate
832 (blue). **(a)** Phylogeny of neurotoxins from *P. pseudoannulata* and other spiders.

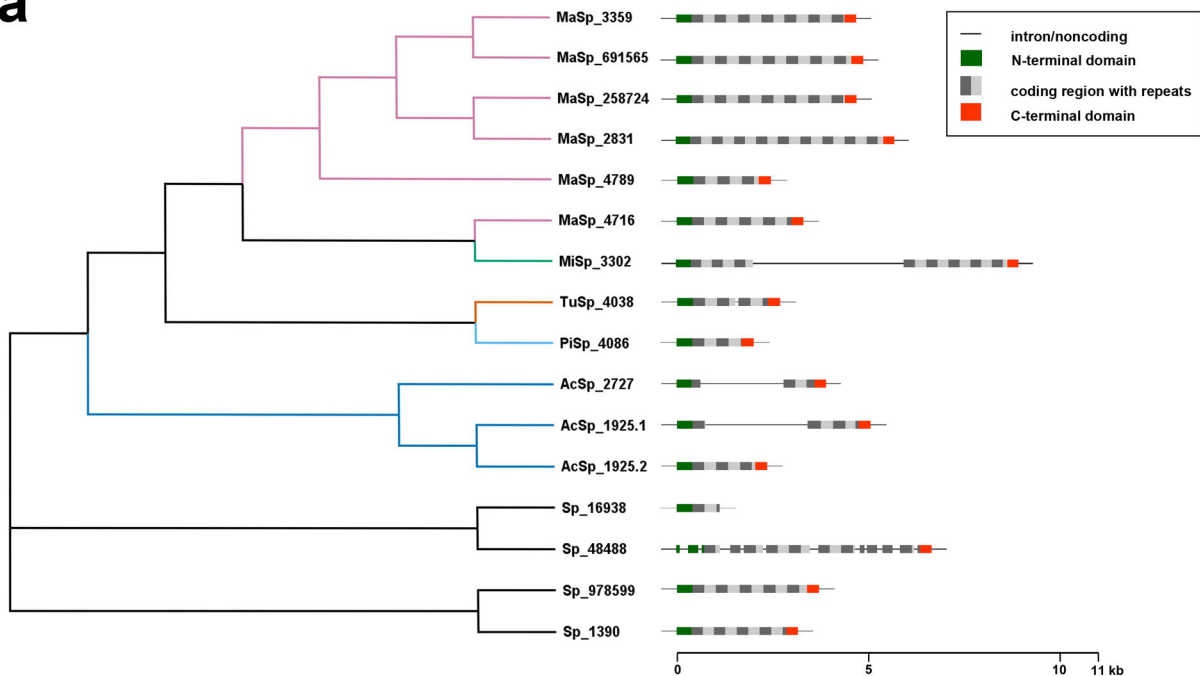
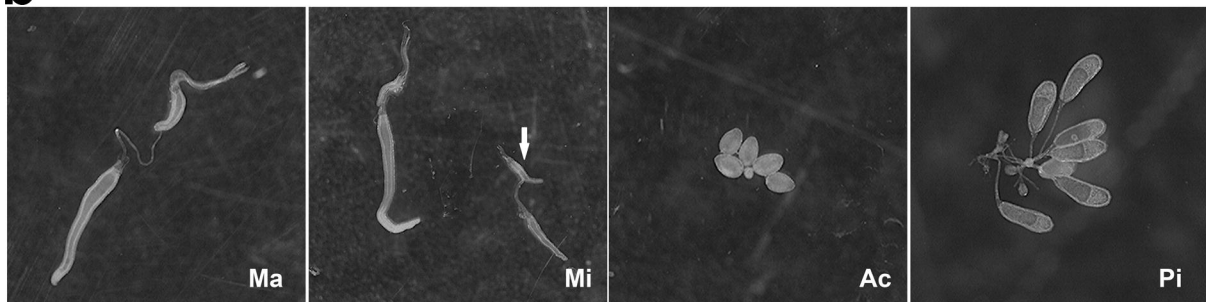
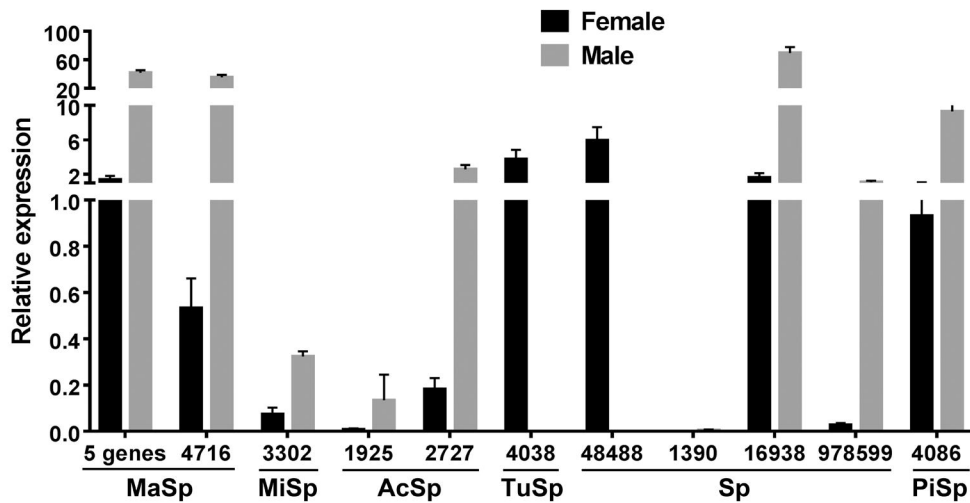
833 Neurotoxins with known target organisms were used as a reference for target
834 organism prediction and are marked with red squares, green circles and blue
835 pentagons. Branches containing these neurotoxins were shaded accordingly. An
836 asterisk marks U1-lycotoxin-Pp1b. **(b)** Distribution of *P. pseudoannulata* neurotoxin
837 subtypes in terms of gene numbers. **(c)** Expression of different *P. pseudoannulata*
838 neurotoxin subtypes as the FPKM value in the venom gland transcriptome. **(d)** The
839 percentage of invertebrate-selective neurotoxins negatively correlated with spider
840 prosoma length. Horizontal bars represent the percentage of neurotoxins based on
841 their target organism selectivity. The black data points indicate prosoma length of the
842 predator spider. **(e)** Diversity of toxin composition in four spiders.
843

a



b



a**b****c**

■ Invertebrate ■ Vertebrate ■ Both

