

1 **Pan-arthropod analysis reveals somatic piRNAs as an ancestral TE defence**

2

3 Samuel H. Lewis^{1,10,11}

4 Kaycee A. Quarles²

5 Yujing Yang²

6 Melanie Tanguy^{1,3}

7 Lise Frézal^{1,3,4}

8 Stephen A. Smith⁵

9 Prashant P. Sharma⁶

10 Richard Cordaux⁷

11 Clément Gilbert^{7,8}

12 Isabelle Giraud⁷

13 David H. Collins⁹

14 Phillip D. Zamore^{2*}

15 Eric A. Miska^{1,3*}

16 Peter Sarkies^{10,11*}

17 Francis M. Jiggins^{1*}

18 *These authors contributed equally to this work

19

20 Correspondence should be addressed to F.M.J. (fmj1001@cam.ac.uk) or S.H.L.

21 (sam.lewis@gen.cam.ac.uk).

22

23

24

25 **Abstract**

26 In animals, PIWI-interacting RNAs (piRNAs) silence transposable elements (TEs),
27 protecting the germline from genomic instability and mutation. piRNAs have been
28 detected in the soma in a few animals, but these are believed to be specific
29 adaptations of individual species. Here, we report that somatic piRNAs were likely
30 present in the ancestral arthropod more than 500 million years ago. Analysis of 20
31 species across the arthropod phylum suggests that somatic piRNAs targeting TEs
32 and mRNAs are common among arthropods. The presence of an RNA-dependent
33 RNA polymerase in chelicerates (horseshoe crabs, spiders, scorpions) suggests that
34 arthropods originally used a plant-like RNA interference mechanism to silence TEs.
35 Our results call into question the view that the ancestral role of the piRNA pathway
36 was to protect the germline and demonstrate that small RNA silencing pathways
37 have been repurposed for both somatic and germline functions throughout arthropod
38 evolution.

39

40 In animals, 23–31 nucleotide (nt) PIWI-interacting RNAs (piRNAs) protect the
41 germline from double-strand DNA breaks and insertion mutagenesis by silencing
42 transposons^{1–3}. In *Drosophila*, piRNAs also function in gonadal somatic cells that
43 support oogenesis^{4,5}. Although the role of piRNAs in the germline appears to be
44 deeply conserved across animals, they have also been reported to function outside
45 the germline. In the mosquito *Aedes aegypti*, there are abundant non-gonadal
46 somatic piRNAs that defend against viruses^{6,7}. In other species, piRNAs are
47 produced in specific cell lineages. For example, somatic piRNAs silence transposons
48 in *D. melanogaster* fat body⁸ and brain^{9,10}, they are important for stem cell
49 maintenance and regeneration in the planarian *Schmidtea mediterranea*^{11,12}, and
50 they contribute to memory in the central nervous system of the mollusc *Aplisia*
51 *californica*¹³.

52 piRNA pathway genes in *Drosophila* species evolve rapidly, likely reflecting an
53 evolutionary arms race with TEs^{14,15}. Expansion and loss of key genes in the piRNA
54 pathway has occurred in platyhelminths¹⁶, nematodes¹⁷, and arthropods^{18–20}. This
55 gene turnover is accompanied by a wide variety of functions for piRNAs, such as sex
56 determination in the silkworm *Bombyx mori* and epigenetic memory formation in the
57 nematode *C. elegans*²¹. There is also considerable divergence in downstream
58 pathways linked to piRNA silencing—for example, in *C. elegans* where piRNAs act
59 upstream of an RNA-dependent RNA polymerase (RdRP) pathway that generates
60 secondary siRNAs antisense to piRNA targets. Moreover, many nematode species
61 have lost the piRNA pathway altogether, with RNAi-related mechanisms assuming
62 the role of TE suppression²². These examples highlight the need for further
63 characterisation across animals to better understand the diversity of the piRNA
64 pathway.

65 To reconstruct the evolutionary history of small RNA pathways, we sampled
66 20 arthropod species with sequenced genomes: three chelicerates, one myriapod,
67 one crustacean, and 15 insects. For each species, we sequenced long and small
68 RNAs from somatic and germline adult (Extended Data Table 1). Our results

69 highlight the rapid diversification of small RNA pathways in animals, challenging
70 previous assumptions based on model organisms. First, we find that RdRP was an
71 integral part of an ancestral siRNA pathway in early arthropods that has been lost in
72 insects. Second, we demonstrate that somatic piRNAs are an ancestral trait of
73 arthropods. Intriguingly, the somatic piRNA pathway is predominantly targeted to
74 transposable elements, suggesting that the piRNA pathway was active in the soma
75 of the last common ancestor of the arthropods to keep mobile genetic elements in
76 check.

77 **Extensive turnover in arthropod small RNA pathways**

78 The duplication or loss of small RNA pathway genes can lead to the gain or loss of
79 small RNA functions. To identify expansions of small RNA genes throughout the
80 arthropods, we identified homologs of key small RNA pathway genes and used
81 Bayesian phylogenetics to reconstruct the timing of duplication and loss (Fig. 1a).
82 Small RNAs bind to Argonaute proteins and guide them to their RNA targets. siRNAs
83 are associated with Ago2-family Argonautes, and these have been extensively
84 duplicated across the arthropods, with an ancient duplication in the arachnid (spider
85 and scorpion) ancestor, and lineage-specific duplications in the scorpion
86 *Centruroides sculpturatus*, the spider *Parasteatoda tepidariorum*, the locust *Locusta*
87 *migratoria*, and the beetle *Tribolium castaneum*²³. piRNAs are associated with PIWI-
88 family Argonautes, which have undergone similar duplications. *Piwi* has duplicated in
89 *L. migratoria*, the centipede *Strigamia maritima*, the pea aphid *Acyrtosiphon*
90 *pisum*¹⁸, the mosquito *Aedes aegypti*²⁴, and flies (generating *piwi* and *aubergine*¹⁹).
91 All species harbour a single copy of *ago3*, which encodes the other PIWI-family
92 Argonaute associated with piRNAs, except for *A. pisum* which has two *ago3* genes.
93 RdRPs amplify an siRNA signal by generating double-stranded RNA (dsRNA)
94 from single-stranded RNA²⁵, but *Drosophila* and other insects lack RdRP genes.
95 RdRP is present in some ticks²⁶, and similarly, we identified RdRP genes across the
96 chelicerates, frequently in multiple copies (Fig. 1a). In each species, one or more

97 RdRPs are expressed in at least one tissue (Extended Data Fig. 1). We also
98 identified an RdRP in the centipede *S. maritima*; however, phylogenetic analysis
99 provides strong evidence that this is not an orthologue of the ancestral arthropod
100 RdRP, but is more closely related to RdRP from fungi (*Neurospora crassa* and
101 *Schizosaccharomyces pombe*; Fig. 1b). In contrast, the chelicerate RdRP is most
102 closely related to other animal RdRPs. Given that RdRPs are present in nematodes
103 and *Nematostella vectensis*, the most parsimonious explanation is that RdRP was
104 present in the common ancestor of arthropods and has been retained in the
105 chelicerates. It was then lost in all other arthropods ~500 MYA, and subsequently
106 regained by *S. maritima* by horizontal gene transfer from a fungus (Figs. 1a,b).

107 The RdRPs expressed in the chelicerates and *S. maritima* may generate
108 dsRNA precursors which can then be processed by Dicer to generate siRNAs,
109 similar to RdRPs in basal nematodes²², while species lacking an RdRP would
110 require bidirectional transcription by RNA polymerase II to generate dsRNA. To test
111 this idea, we sequenced long RNA (RNA-Seq) and small RNA from all species
112 (Extended Data Table 1). Within each species, we identified TEs that were
113 expressed and targeted by siRNAs, and estimated the difference between their
114 sense and antisense expression. Compared to species lacking RdRPs, we find that
115 species with RdRPs have less antisense transcription of these TEs (Mann-Whitney
116 *U* test, animal RdRP versus no RdRP: $p = 0.0381$; Fig. 1c). This pattern is also
117 apparent when comparing antisense transcription and siRNA production across the
118 15 most highly-expressed TEs within a single species. For example, in *H.*
119 *melpomene*, which does not have an RdRP, there is a significant positive correlation
120 between the proportion of antisense transcripts and siRNA production (Spearman
121 rank correlation $\rho = 0.52$, $p = 2 \times 10^{-5}$). Furthermore, none of the TEs with low
122 antisense transcription are among the top siRNA targets (Fig. 1d). These results
123 suggest that *H. melpomene* requires bidirectional transcription to generate siRNAs.
124 In contrast, in *P. tepidariorum* (six RdRPs) there is no correlation between the
125 proportion of antisense transcripts and siRNA production (Spearman rank correlation

126 $\rho = 0.09$, $p = 0.512$), and several TEs with very few antisense transcripts generate
127 abundant siRNAs (Fig. 1d). Together, our results suggest that chelicerates are less
128 dependent on bidirectional transcription to provide the precursors for siRNA
129 production, and may use RdRP to generate dsRNA from TEs, similar to plants and
130 some nematodes. However, we note that the antisense enrichment for siRNA targets
131 in *S. maritima* is more similar to species lacking an RdRP, making it unclear whether
132 its horizontally-transferred RdRP acts in this way.

133 **Germline piRNAs are found across arthropods**

134 Current evidence supports the view that the piRNA pathway is a germline-specific
135 defence against transposon mobilization. As expected, we found piRNAs derived
136 from the genome in the female germline of all 20 arthropod species (Extended Data
137 Table 1, Extended Data Fig. 2), consistent with deep conservation of this function
138 from the last common ancestor of mammals and arthropods. Germline piRNAs target
139 TEs in a wide variety of animals, including nematodes, fish, birds, and mammals, as
140 was the case in all our species (Extended Data Fig. 3); moreover, TE abundance
141 and piRNA abundance were positively correlated as previously found in *D.*
142 *melanogaster* (Extended Data Fig. 4). In 10 species, we also sampled the male
143 germline. Male germline piRNAs were found in all species except the bumblebee
144 *Bombus terrestris*, which lacked detectable piRNAs in both testis and mature sperm-
145 containing vas deferens, even when using a protocol that specifically enriches for
146 piRNAs by depleting miRNAs⁹ (Fig. 2a; Extended Data Figs. 5 and 6). In contrast,
147 piRNAs were abundant in *B. terrestris* ovary (Fig. 2b; Extended Data Fig. 2).
148 Moreover, mRNAs encoding the core piRNA pathway proteins Piwi and Vasa were
149 10-fold less abundant in testis compared to ovary (Extended Data Fig. 7), suggesting
150 that the piRNA pathway is not active in the *B. terrestris* male germline. To our
151 knowledge, this is the first report of sex-specific absence of piRNAs in the germline,
152 and suggests that other processes may have taken on the function of TE
153 suppression in *B. terrestris* males. Male bumblebees are haploid and produce sperm

154 by mitosis rather than meiosis²⁷, unlike males from the other eight species analysed.
155 However, in the testis of the haplodiploid honey bee *Apis mellifera* piRNAs are
156 detectable by their characteristic Ping-Pong signature, albeit at low levels (Extended
157 Data Figs. 5 and 8).

158 **Somatic piRNAs are widespread across arthropods**

159 Among the 20 arthropods we surveyed, somatic piRNAs were readily detected in 16
160 species: three chelicerates (*L. polyphemus*, *C. sculpturatus*, and *P. tepidariorum*),
161 the myriapod *S. maritima*, and 12 insect species (Figs. 1a and 3c,d; Extended Data
162 Fig. 9). We did not detect piRNAs in the somatic tissues of the crustacean
163 *Armadillidium vulgare* or the insects *N. vespilloides*, *B. terrestris*, and *D.*
164 *melanogaster* (Extended Data Fig. 9). Although somatic piRNAs have been detected
165 previously in *D. melanogaster* heads^{9,10}, we detected no piRNAs in *D. melanogaster*
166 thorax. Somatic expression of the piRNA pathway genes *vasa*, *ago3*, *Hen1*, and *Piwi*
167 was strongly associated with the presence of somatic piRNAs (Fig. 3a). We conclude
168 that an active somatic piRNA pathway is widespread throughout the arthropods.

169 The phylogenetic distribution of somatic piRNAs suggests that they were
170 either ancestral to all arthropods or have been independently gained in different
171 lineages. To distinguish between these possibilities, we used ancestral state
172 reconstruction to infer the presence or absence of somatic piRNAs on the internal
173 branches of the arthropod phylogeny. Our results indicate that somatic piRNAs are
174 ancestral to all arthropods (posterior probability = 1), and have been independently
175 lost at least four times (Fig. 1a).

176 **Functions of somatic piRNAs**

177 In all but one species with somatic piRNAs, at least 2% of piRNAs mapped to TEs
178 (Fig. 3c, Extended Data Fig. 3), suggesting that their anti-transposon role is
179 conserved in the soma. The exception to this pattern was *O. fasciatus*, where only
180 0.009% of somatic and 0.074% of germline piRNAs were derived from annotated
181 TEs. Moreover, somatic piRNAs from all species displayed the hallmark features of

182 piRNA biogenesis and amplification: a 5' uracil bias, 5' ten-nucleotide
183 complementarity between piRNAs from opposite genomic strands (“Ping-Pong”
184 signature), and resistance to oxidation by sodium periodate, consistent with their
185 bearing a 2'-O-methyl modification at their 3' ends (e.g., Fig. 3d). Given the ubiquity
186 of TE-derived somatic piRNAs, we wondered whether there was a relationship
187 between the TE content of a species' genome and the presence of somatic piRNAs.
188 However, although species with somatic piRNAs tend to have a higher TE content,
189 this difference is not significant ($p = 0.18$, Extended Data Fig. 10).

190 In *Drosophila*, piRNAs derived from protein-coding genes are thought to play
191 a role in regulating gene expression²⁸. Somatic piRNAs derived from protein-coding
192 sequences and untranslated regions (UTRs) were present in all species possessing
193 somatic piRNAs except *A. mellifera*, *D. virilis* and *M. domestica*, which lack both a
194 distinct peak of 25-29nt sRNAs and a Ping-Pong signature (Extended Data Fig. 3,
195 Extended Data Fig. 11). When scaled to the genome content of each feature, there
196 is no consistent difference in the abundance of piRNAs from protein-coding
197 sequence and UTRs (Extended Data Fig. 12), suggesting that somatic piRNAs target
198 genes across the entire length of the transcript, rather than just UTRs.

199 In the mosquito *A. aegypti*, somatic piRNAs target viruses^{6,7}. To test whether
200 somatic piRNAs derive from viruses in other species, we reconstructed partial viral
201 genomes from each species using somatic RNA-Seq data, then mapped small RNAs
202 from these tissues to these viral contigs. In *A. aegypti*, we recovered the partial
203 genome of a positive-sense, single-stranded RNA virus that was targeted by both
204 siRNAs (21 nt) and 5' U-biased, 25–30 nt piRNAs bearing the signature of Ping-Pong
205 amplification (Fig. 4a). These data recapitulate previous results showing that both
206 the siRNA and piRNA pathways mount an antiviral response in *A. aegypti*⁶, and thus
207 validate our approach. In eight additional species, we could similarly reconstruct
208 viruses that generated antiviral siRNAs (Fig. 4c, Extended Data Fig. 13). Four of
209 these species also produced 25–30 nt, 5' U-biased RNAs derived from viruses
210 including negative- and positive-sense RNA viruses and DNA viruses (Fig. 4b,

211 Extended Data Fig. 13). There was no evidence of Ping-Pong amplification of viral
212 piRNAs in any of these species—in *C. sculpturatus* somatic piRNAs were of low
213 abundance (Fig. 4b), and in *T. castaneum*, *D. virgifera* and *P. xylostella* piRNAs
214 mapped to only one strand (Extended Data Fig. 13), a feature reminiscent of the
215 somatic piRNAs present in *Drosophila* follicle cells^{4,5}. Despite removing sequencing
216 reads that map to the reference genome, we cannot exclude the possibility that
217 these piRNAs come from viruses integrated in the host genome²⁹. Together these
218 results suggest that although some viruses may be targeted by somatic piRNAs,
219 siRNAs likely remain the primary antiviral defence against most viruses across the
220 arthropods.

221 **Conclusions**

222 The rapid evolution of small RNA pathways makes inferences drawn from detailed
223 studies of individual model organisms misleading²². Our results suggest that the best
224 studied arthropods, concentrated in a small region of the phylogenetic tree, are not
225 representative of the entire phylum (Fig. 5). First, ancestral arthropods likely used an
226 RdRP to generate siRNAs from transposable elements. RdRPs likely expand the
227 range of substrates that can generate siRNAs, because these RNA-copying
228 enzymes provide an alternative to the generation of dsRNA precursors by RNA
229 polymerase II. Second, and more surprising, somatic piRNAs are ubiquitous across
230 arthropods, where they target transposable elements and mRNAs. The rapid and
231 dynamic evolution of somatic and germline piRNA pathways across the arthropods
232 highlights the need for a deeper examination of the origins and adaptations of the
233 piRNA pathway in other phyla.

234

235 **Methods**

236 **Tissue dissection**

237 To sample germline tissue from each species, we dissected the female germline of
238 all 20 arthropods (ovary and accessory tissue). For *Limulus polyphemus*,
239 *Centruroides sculpturatus*, *Parasteatoda tepidariorum*, *Armadillidium vulgare*,
240 *Locusta migratoria*, *Bombus terrestris*, *Apis mellifera*, *Nicrophorus vespilloides*,
241 *Heliconius melpomene* and *Trichoplusia ni*, we also dissected the male germline
242 (testes, vas deferens, and accessory tissue). We were unable to isolate sufficient
243 germline tissue for *Strigamia maritima*.

244 To isolate somatic tissue, we used different dissection approaches depending
245 on the anatomy of the species. In each case, we minimized the risk of germline
246 contamination by selecting tissue from either a body region that was separate (e.g.,
247 thorax) or physically distant from the germline. For insects, thorax served as a
248 representative somatic tissue. For *Oncopeltus fasciatus*, *Acyrtosiphon pisum*, *Apis*
249 *mellifera*, *Tribolium castaneum*, *Diabrotica virgifera*, *Plutella xylostella*, *Aedes*
250 *aegypti*, *Musca domestica* and *Drosophila melanogaster* we used female thorax; for
251 *Locusta migratoria*, *Bombus terrestris*, *Nicrophorus vespilloides*, *Heliconius*
252 *melpomene*, and *Trichoplusia ni* we used female and male thorax separately. For
253 non-insect species, we took mixed tissue from either the mesosoma (*Parasteatoda*
254 *tepidariorum*), prosoma (*Centruroides sculpturatus*), pereon and pleon
255 (*Armadillidium vulgare*) or muscle, heart, and liver (*Limulus polyphemus*). For these
256 non-insect species, we isolated somatic tissue from males and females separately.
257 For *Strigamia maritima*, we pooled female and male fat body.

258 **RNA extraction and library preparation: Protocol 1**

259 For *Limulus polyphemus*, *Centruroides sculpturatus*, *Parasteatoda tepidariorum*,
260 *Strigamia maritima*, *Armadillidium vulgare*, *Locusta migratoria*, *Bombus terrestris*,
261 *Nicrophorus vespilloides* and *Heliconius melpomene* we extracted total RNA and
262 constructed sequencing libraries using Protocol 1. Following dissection, each sample

263 was homogenized in Trizol (Invitrogen, Carlsbad, CA, USA) and stored at -80°C .
264 RNA from each sample was extracted with isopropanol/chloroform (2.5:1), and RNA
265 integrity was checked using the Bioanalyzer RNA Nano kit (Agilent, Santa Clara, CA,
266 USA).

267 For small RNA sequencing, each sample was initially spiked with *C. elegans*
268 RNA (N2 strain) at 1/10th mass of the input RNA (e.g., 0.1 μg *C. elegans* RNA with 1
269 μg sample RNA). This allowed us to quantify the efficiency of sRNA library
270 production. To sequence all small RNAs in a 5'-independent manner, we removed 5'
271 triphosphates by treating each sample with 5' polyphosphatase (Epicentre/Illumina,
272 Madison, WI, USA) for 30 min. We used the TruSeq Small RNA Library Preparation
273 Kit (Illumina, San Diego, CA, USA) according to the manufacturer's instructions to
274 produce libraries from total RNA. We sequenced each sRNA library on a HiSeq 1500
275 (Illumina) to generate 36 nt single-end reads.

276 piRNAs are typically 2'-O-methylated at their 3' ends, which makes them
277 resistant to sodium periodate oxidation. To test for the presence of modified 3' ends,
278 we resuspended RNA in 5x borate buffer, treated with sodium periodate (25 mM f.c.,
279 e.g., 5 μl 200 mM sodium periodate in 40 μl reaction) for 10 min, recovered the
280 treated RNA by ethanol precipitation³⁰ and constructed and sequenced libraries as
281 above.

282 For transcriptome and virus RNA-Seq, each sample was initially spiked with
283 *C. elegans* RNA (N2 strain) at 1/10th mass of the input RNA. To remove ribosomal
284 RNA, we treated each sample with the Ribo-Zero rRNA Removal Kit
285 (Human/Mouse/Rat; Illumina) according to manufacturer's instructions, then
286 prepared strand-specific RNA-Seq libraries using the NEBNext Ultra Directional RNA
287 Library Prep kit (New England Biolabs, Ipswich, MA, USA), with the optional User
288 Enzyme step to selectively degrade the 2nd strand before PCR amplification. RNA-
289 Seq libraries were sequenced on a HiSeq 4000 to generate 150 nt paired-end reads
290 (*C. sculpturatus* and *S. maritima*), or a HiSeq 2500 to generate 125 nt paired-end
291 reads (all other species).

292 **RNA extraction and library preparation: Protocol 2**

293 For *Oncopeltus fasciatus*, *Acyrtosiphon pisum*, *Apis mellifera*, *Tribolium castaneum*,
294 *Diabrotica virgifera*, *Plutella xylostella*, *Trichoplusia ni*, *Aedes aegypti*, *Musca*
295 *domestica*, *Drosophila virilis* and *Drosophila melanogaster* we extracted total RNA
296 and constructed sequencing libraries using Protocol 2. Following dissection, we
297 washed each sample in PBS, proceeded directly to RNA extraction using the
298 *mirVana* miRNA Isolation kit (Ambion, Life Technologies, CA, USA) according to the
299 manufacturer's protocol, and precipitated RNA with ethanol. We prepared RNA-Seq
300 libraries for each sample from 5 µg total RNA as described³¹, after first depleting
301 rRNA using the Ribo-Zero rRNA Removal Kit (Human/Mouse/Rat; Illumina). We
302 sequenced each library on a NextSeq 500 (Illumina) to generate 79 nt paired-end
303 reads.

304 Small RNA sequencing libraries were generated as described³². First, we
305 purified 16–35 nt RNA from 10–20 µg total RNA by 15% denaturing urea-
306 polyacrylamide gel electrophoresis. Half of each sample was then treated with
307 sodium periodate (above). We then ligated 3' pre-adenylated adapter to treated or
308 untreated RNA using homemade, truncated mutant K227Q T4 RNA ligase 2 (amino
309 acids 1–249) and purified the 3'-ligated product by 15% denaturing urea-
310 polyacrylamide gel electrophoresis. To exclude 2S rRNA from sequencing libraries,
311 2S blocker oligo³³ was added to all samples before the 5'-adapter was appended
312 using T4 RNA ligase (Ambion). cDNA was synthesized using AMV reverse
313 transcriptase (New England Biolabs) and the reverse transcription primer 5'-
314 CCTTGGCACCCGAGAATTCCA-3'. The small RNA library was amplified using
315 AccuPrime Pfx DNA polymerase (ThermoFisher, USA) and forward (5'-
316 AATGATACGGCGACCACCGAGATCTACACGTTTCAGAGTTCTACAGTCCGA-3')
317 and barcoded reverse (5'-CAAGCAGAAGACGGCATAACGAGAT-barcode(N₆)-
318 GTGACTGGAGTTCCTTGGCACCCGAGAATTCCA-3') primers, purified from a 2%
319 agarose gel, and sequenced on a NextSeq 500 to generate 50 nt single-end reads.

320 **Bioinformatics analysis**

321 **Gene family evolution**

322 To reconstruct duplications and losses of sRNA pathway components, we searched
323 for homologs of *Ago1*, *Ago2*, *Ago3*, *Piwi*, *Dcr1*, *Dcr2*, *Drosha*, *Hen1* and *Vasa*. For
324 each species, we took the annotated protein set and used DIAMOND³⁴ to perform
325 reciprocal all-versus-all BLASTp searches against all proteins in *D. melanogaster*,
326 and retained only the top hit in each case. Accession numbers for the genome
327 assemblies and annotated protein sets are detailed in Extended Data Table 2. To
328 find homologs of RdRP, which is absent from *D. melanogaster*, we took the
329 annotated protein set for each species and used DIAMOND to perform BLASTp
330 searches against the RdRP from *Ixodes scapularis* (ISCW018089). For proteins in
331 the Argonaute and Dicer families, we identified domains in hits using
332 InterProScan5³⁵ with the Pfam database, and retained only those hits containing at
333 least one of the conserved domains in these families (PAZ and Piwi for the
334 Argonaute family, PAZ, Dicer, Ribonuclease and Helicase for the Dicer family). For
335 each protein, partial BLAST hits were manually curated into complete proteins if the
336 partial hits were located adjacent to each other on the same scaffold or contig. To
337 establish the evolutionary relationships between homologs, we aligned each set of
338 homologs as amino acid sequences using MAFFT³⁶ with default settings, screened
339 out poorly aligned regions using Gblocks³⁷ with the least stringent settings, and
340 inferred a gene tree using the Bayesian approach implemented in MrBayes v3.2.6³⁸.
341 We specified a GTR substitution model with gamma-distributed rate variation and a
342 proportion of invariable sites. We ran the analysis for 10 million generations,
343 sampling from the posterior every 1000 generations.

344 **Transposable element annotation**

345 To annotate transposable elements (TEs) in each genome, we used RepeatMasker
346 v4.0.6³⁹ with the “Metazoa” library to identify homologs to any previously-identified
347 metazoan TEs. In addition, we used RepeatModeler v1.0.8⁴⁰ to generate a *de novo*

348 Hidden Markov Model for TEs in each genome, and ran RepeatMasker using this
349 HMM to identify TEs without sufficient homology to previously-identified metazoan
350 TEs. We combined these two annotations to generate a single, comprehensive TE
351 annotation file for each species. We then screened out all annotations <100 nt long.
352 The source code for this analysis is accessible on GitHub
353 (<https://github.com/SamuelHLewis/TEAnnotator>), and the TE annotation files are
354 available from the Cambridge Data Archive (<https://doi.org/10.17863/CAM.10266>).

355 **Virus identification and genome assembly**

356 To identify viruses, we first mapped RNA-Seq reads to the genome of the host
357 species to exclude genome-derived transcripts, thus filtering out endogenous viral
358 elements. We then used Trinity⁴¹ with default settings to generate a *de novo*
359 assembly of the remaining RNA-Seq data for each tissue, and extracted the protein
360 sequence corresponding to the longest open reading frame for each contig with
361 TransDecoder (<https://transdecoder.github.io/>), excluding all contigs shorter than
362 100nt. To identify contigs that were potentially of viral origin, we used DIAMOND to
363 perform BLASTp searches against all viral proteins in NCBI
364 (<ftp.ncbi.nih.gov/refseq/release/viral/viral.1.protein.faa.gz> and
365 <ftp.ncbi.nih.gov/refseq/release/viral/viral.2.protein.faa.gz>, downloaded 19/10/16). To
366 screen out false-positive hits from those contigs with similarity to a viral protein, we
367 used DIAMOND to perform BLASTp searches against the NCBI non-redundant (nr)
368 database (downloaded 19/10/16) and retained only those contigs which still had a
369 virus as their top hit. The source code for this analysis is accessible on GitHub
370 (<https://github.com/SamuelHLewis/VirusFinder>), and the viral contigs are available
371 from GenBank (accession codes XXXXXX-XXXXXX).

372 **Small RNA analysis**

373 To characterize sRNAs derived from the genome in each tissue of each species, we
374 first used the FASTX Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/) to screen out
375 small RNA reads with >10% positions with a Qphred score <20 and cutadapt⁴² to

376 trim adapter sequences from reads. We then mapped small RNAs to the genome
377 using Bowtie2 v2.2.6⁴³ in "--fast" mode, which reports the best alignment for reads
378 mapping to multiple locations, or a randomly-chosen location if there are multiple
379 equally-good alignments. We quantified the length distribution, base composition,
380 and strand distribution of sRNAs mapping to the genome using a custom Python
381 script (accessible on GitHub <https://github.com/SamuelHLewis/sRNAplot>),
382 considering unique sRNA sequences only.

383 To characterize sRNAs targeting TEs, we used BEDTools getfasta⁴⁴ to extract
384 TE sequences from the genome in a strand-specific manner (according to the TE
385 annotation for each genome, above), mapped sRNAs as detailed above, and
386 quantified their characteristics using the same custom Python script
387 (<https://github.com/SamuelHLewis/sRNAplot>), this time considering all sRNA
388 sequences. To characterize sRNAs targeting viruses, we first screened out genome-
389 derived sRNAs by mapping sRNAs to the genome and retaining unmapped reads.
390 We then used the same mapping procedure as detailed above, applied to each virus
391 separately.

392 To characterize small RNAs mapping to UTRs in each species (except *D.*
393 *virgifera*, *D. virilis* and *D. melanogaster*), we extracted 200 nt upstream (5' UTR) or
394 downstream (3' UTR) of each gene model. To ensure that these UTR sequences did
395 not overlap with TEs, we masked any sequence that we had annotated as a TE
396 using RepeatMasker (see above). We then screened out TE-derived sRNAs by
397 mapping sRNAs to the TE annotations and retaining unmapped reads. These were
398 mapped to our UTR annotations as detailed for TEs (above). For *D. melanogaster*
399 and *D. virilis* we employed the same method but used the curated set of 5' and 3'
400 UTRs from genomes r6.15 and r1.06 respectively. We excluded *D. virgifera* from this
401 analysis as gene models have not been predicted for its genome.

402 For each species, we defined the presence of UTR-derived piRNAs based on
403 the presence of >200 unique 25-29nt sequences with a 5' U nucleotide bias. For
404 species with somatic piRNAs, we used oxidized sRNA data to assay the presence or

405 absence of somatic UTR-derived piRNAs. We excluded *D. virgifera* from this
406 analysis because of a lack of annotated gene models.

407 To test whether piRNAs show evidence of ping-pong amplification, we
408 calculated whether sense and antisense 25–29nt reads tended to overlap by 10 nt
409 using the z-score method of Zhang et al^{45–47}.

410 **Gene expression analysis**

411 To quantify the expression of genes in small RNA pathways in each tissue, we first
412 used Trim Galore (<https://github.com/FelixKrueger/TrimGalore>) with default settings
413 to trim adapters and low-quality ends from each RNA-Seq mate pair. We then
414 mapped these reads to the genome using Tophat2 v2.1.1⁴⁸ with default settings in “--
415 library-type fr-firststrand” mode. To calculate FPKM values for each gene we used
416 DESeq2⁴⁹, specifying strand-specific counts and summing counts for each gene by
417 all exons. We excluded *D. virgifera* from this analysis because a genome annotation
418 file is unavailable. The source code for this analysis is accessible on GitHub
419 (<https://github.com/SamuelHLewis/GeneExpression>).

420 **Species tree reconstruction**

421 To provide a timescale for the evolution of arthropod sRNA pathways, we combined
422 published phylogenies of insects⁵⁰ and arthropods⁵¹ with our own estimates of
423 divergence dates and branch lengths. We first gathered homologs of 163 proteins
424 that are present as 1:1:1 orthologues in each of our focal species. We then
425 generated a concatenated alignment of these proteins using MAFFT³⁶ with default
426 settings, and screened out poorly-aligned regions with Gblocks³⁷ in least stringent
427 mode. We used this alignment to carry out Bayesian phylogenetic analysis as
428 implemented in BEAST⁵², to infer branch lengths for the phylogeny of our sample
429 species. We specified a birth-death speciation process, a strict molecular clock,
430 gamma distributed rate variation with no invariant sites, and fixed the topology and
431 set prior distributions on key internal node dates (Arthropoda = 568 ± 29 , Insecta-
432 Crustacea = 555 ± 33 , Insecta = 386 ± 27 , Hymenoptera-Coleoptera-Lepidoptera-

433 Diptera = 345 ± 27 , Coleoptera-Lepidoptera-Diptera= 327 ± 26 , Lepidoptera-Diptera =
434 290 ± 46 , Diptera = 158 ± 51) based on a previous large-scale phylogenetic analysis
435 of arthropods⁵⁰. We ran the analysis for 1.5 million generations, and generated a
436 maximum clade credibility tree with TreeAnnotator⁵².

437 **TE content analysis**

438 To compare the TE content of species with and without somatic piRNAs, we used
439 the TE annotations derived from RepeatModeler (above) to calculate the TE content
440 of each genome as a proportion of the entire genome size. We then tested for a
441 difference in TE content between species with and without somatic piRNAs using a
442 phylogenetic general linear mixed model to account for non-independence due to the
443 phylogenetic relationships. The model was implemented using a Bayesian approach
444 in the R package MCMCglmm⁵³ based on the time-scaled species phylogeny (see
445 above). The source code for this analysis is accessible on GitHub
446 (<https://github.com/SamuelHLewis/TEContent>).

447 **RdRP signature**

448 In species with an RdRP, siRNAs can be produced from loci that are transcribed
449 from just the sense strand, as the RdRP synthesizes the complementary strand,
450 whereas in species that lack an RdRP, siRNAs can only be produced from loci that
451 have both sense and antisense transcription. To test the association between siRNA
452 production and antisense transcription in each species, we first used Trimmomatic⁵⁴
453 to extract sRNAs corresponding to the median siRNA length in that species. We then
454 used Bowtie v2.2.6⁴³ in "--fast" mode to map siRNAs and RNA-Seq reads to TE
455 sequences in each genome, and generated strand-specific counts of siRNAs and
456 RNA-Seq reads for each TE using BEDTools coverage⁴⁴. We then calculated the
457 enrichment of antisense expression [$\log_2(\text{antisense RNA-Seq reads}) - \log_2(\text{sense}$
458 $\text{RNA-Seq reads})$] at TEs with >5 RNA-Seq reads per million and >100 siRNAs per
459 million sRNA reads in species with and without RdRP (Fig. 1c), and tested for a
460 difference in enrichment between species with and without RdRP (excluding *S.*

461 *maritima*) using a Wilcoxon unpaired test. We also plotted the 60 most highly
462 expressed TEs for *H. melpomene* and *P. tepidariorum* and highlighted which of
463 these loci were among the top 15 siRNA-producing TEs (Fig. 1d). The source code
464 for this analysis is accessible on GitHub (<https://github.com/SamuelHLewis/RdRP>).

465 **Data Availability**

466 Sequence data that support the findings of this study have been deposited in the
467 NCBI Short Read Archive with the accession code XXXXXXXX. Length distributions
468 of TE-mapping small RNAs and raw data used to plot Figures 1c, 1d & 3a and
469 Extended Data Figures 1, 7 & 10 are available on the Cambridge Data Repository
470 (<https://doi.org/10.17863/CAM.10266>).

471 **Code Availability**

472 Source code used in this study is accessible on GitHub
473 (<https://github.com/SamuelHLewis>), please see Methods for details of source code
474 used in each analysis.

475 **References**

- 476 1. Aravin, A. A. *et al.* Double-stranded RNA-mediated silencing of genomic
477 tandem repeats and transposable elements in the *D. melanogaster* germline.
478 *Curr. Biol.* **11**, 1017–27 (2001).
- 479 2. Aravin, A., Lagos-Quintana, M. & Yalcin, A. The Small RNA Profile during
480 *Drosophila melanogaster* Development. *Dev. Cell* **5**, 337–350 (2003).
- 481 3. Czech, B. & Hannon, G. J. One Loop to Rule Them All: The Ping-Pong Cycle
482 and piRNA-Guided Silencing. *Trends Biochem. Sci.* **41**, 324–337 (2016).
- 483 4. Li, C. *et al.* Collapse of germline piRNAs in the absence of Argonaute3 reveals
484 somatic piRNAs in flies. *Cell* **137**, 509–21 (2009).
- 485 5. Malone, C. D. *et al.* Specialized piRNA Pathways Act in Germline and Somatic
486 Tissues of the *Drosophila* Ovary. *Cell* **137**, 522–535 (2009).
- 487 6. Morazzani, E. M., Wiley, M. R., Murreddu, M. G., Adelman, Z. N. & Myles, K.

- 488 M. Production of virus-derived ping-pong-dependent piRNA-like small RNAs in
489 the mosquito soma. *PLoS Pathog.* **8**, e1002470 (2012).
- 490 7. Miesen, P., Girardi, E. & van Rij, R. P. Distinct sets of PIWI proteins produce
491 arbovirus and transposon-derived piRNAs in *Aedes aegypti* mosquito cells.
492 *Nucleic Acids Res.* **43**, 6545–56 (2015).
- 493 8. Jones, B. C. *et al.* A somatic piRNA pathway in *Drosophila* fat body
494 suppresses transposable elements ensuring metabolic homeostasis and
495 normal lifespan. *Nat. Commun.* **7**, 13856 (2016).
- 496 9. Ghildiyal, M. *et al.* Endogenous siRNAs Derived from Transposons and
497 mRNAs in *Drosophila* Somatic Cells. *Science* **320**, 1077–1081 (2008).
- 498 10. Perrat, P. N. *et al.* Transposition-driven genomic heterogeneity in the
499 *Drosophila* brain. *Science* **340**, 91–5 (2013).
- 500 11. Palakodeti, D., Smielewska, M., Lu, Y.-C., Yeo, G. W. & Graveley, B. R. The
501 PIWI proteins SMEDWI-2 and SMEDWI-3 are required for stem cell function
502 and piRNA expression in planarians. *RNA* **14**, 1174–1186 (2008).
- 503 12. Reddien, P. W., Oviedo, N. J., Jennings, J. R., Jenkin, J. C. & Sánchez
504 Alvarado, A. SMEDWI-2 is a PIWI-like protein that regulates planarian stem
505 cells. *Science* **310**, 1327–1330 (2005).
- 506 13. Rajasethupathy, P. *et al.* A role for neuronal piRNAs in the epigenetic control
507 of memory-related synaptic plasticity. *Cell* **149**, 693–707 (2012).
- 508 14. Obbard, D. J., Gordon, K. H. J., Buck, A. H. & Jiggins, F. M. The evolution of
509 RNAi as a defence against viruses and transposable elements. *Philos. Trans.
510 R. Soc. London Biol. Sci.* **364**, 99–115 (2009).
- 511 15. Kolaczkowski, B., Hupaló, D. N. & Kern, A. D. Recurrent adaptation in RNA
512 interference genes across the *Drosophila* phylogeny. *Mol. Biol. Evol.* **28**,
513 1033–1042 (2011).
- 514 16. Skinner, D. E., Rinaldi, G., Koziol, U., Brehm, K. & Brindley, P. J. How might
515 flukes and tapeworms maintain genome integrity without a canonical piRNA
516 pathway? *Trends Parasitol.* **30**, 123–129 (2014).
- 517 17. Buck, A. H. & Blaxter, M. Functional diversification of Argonautes in

- 518 nematodes: an expanding universe. *Biochem. Soc. Trans.* **41**, 881–6 (2013).
- 519 18. Dowling, D. *et al.* Phylogenetic Origin and Diversification of RNAi Pathway
520 Genes in Insects. *Genome Biol. Evol.* **8**, 3784–3793 (2017).
- 521 19. Lewis, S. H., Salmela, H. & Obbard, D. J. Duplication and diversification of
522 Dipteran Argonaute genes, and the evolutionary divergence of Piwi and
523 Aubergine. *Genome Biol. Evol.* **8**, 507–518 (2016).
- 524 20. Palmer, W. J. & Jiggins, F. M. Comparative Genomics Reveals the Origins and
525 Diversity of Arthropod Immune Systems. *Mol. Biol. Evol.* **32**, 2111–2129
526 (2015).
- 527 21. Sarkar, A., Volff, J. N. & Vaury, C. piRNAs and their diverse roles: a
528 transposable element-driven tactic for gene regulation? *FASEB J.* **31**, 436–446
529 (2017).
- 530 22. Sarkies, P. *et al.* Ancient and Novel Small RNA Pathways Compensate for the
531 Loss of piRNAs in Multiple Independent Nematode Lineages. *PLoS Biol.* **13**,
532 e1002061 (2015).
- 533 23. Tomoyasu, Y. *et al.* Exploring systemic RNA interference in insects: a
534 genome-wide survey for RNAi genes in *Tribolium*. *Genome Biol.* **9**, R10
535 (2008).
- 536 24. Campbell, C. L., Black, W. C., Hess, A. M. & Foy, B. D. Comparative genomics
537 of small RNA regulatory pathway components in vector mosquitoes. *BMC*
538 *Genomics* **9**, 425 (2008).
- 539 25. Schiebel, W. *et al.* Isolation of an RNA-Directed RNA Polymerase Specific
540 cDNA Clone from Tomato. *Plant Cell* **10**, 2087–2102 (1998).
- 541 26. Zong, J., Yao, X., Yin, J., Zhang, D. & Ma, H. Evolution of the RNA-dependent
542 RNA polymerase (RdRP) genes: Duplications and possible losses before and
543 after the divergence of major eukaryotic groups. *Gene* **447**, 29–39 (2009).
- 544 27. Bull, J. J. Advantage for the evolution of male. *Heredity* **43**, 361–381 (1979).
- 545 28. Robine, N. *et al.* A Broadly Conserved Pathway Generates 3'UTR-Directed
546 Primary piRNAs. *Curr. Biol.* **19**, 2066–2076 (2009).
- 547 29. Palatini, U. *et al.* Comparative genomics shows that viral integrations are

- 548 abundant and express. *BMC Genomics* **18**, 512 (2017).
- 549 30. Alefelder, S., Patel, B. K. & Eckstein, F. Incorporation of terminal
550 phosphorothioates into oligonucleotides. *Nucleic Acids Res.* **26**, 4983–4988
551 (1998).
- 552 31. Zhang, Z., Theurkauf, W. E., Weng, Z. & Zamore, P. D. Strand-specific
553 libraries for high throughput RNA sequencing (RNA-Seq) prepared without
554 poly(A) selection. *Silence* **3**, 9 (2012).
- 555 32. Han, B. W., Wang, W., Li, C. & Weng, Z. piRNA-guided transposon cleavage
556 initiates Zucchini-dependent, phased piRNA production. *Science* **348**, 817–821
557 (2015).
- 558 33. Wickersheim, M. L. & Blumenstiel, J. P. Terminator oligo blocking efficiently
559 eliminates rRNA from Drosophila small RNA sequencing libraries.
560 *Biotechniques* **55**, 269–272 (2013).
- 561 34. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
562 DIAMOND. *Nat. Methods* **12**, 59–60 (2015).
- 563 35. Jones, P. *et al.* InterProScan 5: Genome-scale protein function classification.
564 *Bioinformatics* **30**, 1236–1240 (2014).
- 565 36. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT : a novel method for
566 rapid multiple sequence alignment based on fast Fourier transform. *Nucleic
567 Acids Res.* **30**, 3059–3066 (2002).
- 568 37. Castresana, J. Selection of conserved blocks from multiple alignments for their
569 use in phylogenetic analysis. *Mol. Biol. Evol.* **17**, 540–552 (2000).
- 570 38. Ronquist, F. & Huelsenbeck, J. P. MrBayes 3: Bayesian phylogenetic
571 inference under mixed models. *Bioinformatics* **19**, 1572–1574 (2003).
- 572 39. Smit, A. F. A., Hubley, R. & Green, P. RepeatMasker. (2013).
- 573 40. Smit, A. F. A. & Hubley, R. RepeatModeler. (2008).
- 574 41. Grabherr, M. G. *et al.* Full-length transcriptome assembly from RNA-Seq data
575 without a reference genome. *Nat. Biotechnol.* **29**, 644–52 (2011).
- 576 42. Martin, M. Cutadapt removes adapter sequences from high-throughput
577 sequencing reads. *EMBnet.journal* **17**, 10 (2011).

- 578 43. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2.
579 *Nat Methods* **9**, 357–359 (2012).
- 580 44. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing
581 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 582 45. Brennecke, J. *et al.* Discrete small RNA-generating loci as master regulators of
583 transposon activity in *Drosophila*. *Cell* **128**, 1089–103 (2007).
- 584 46. Zhang, Z. *et al.* Heterotypic piRNA Ping-Pong requires qin, a protein with both
585 E3 ligase and Tudor domains. *Mol. Cell* **44**, 572–84 (2011).
- 586 47. Antoniewski, C. in *Animal Endo-siRNAs: Methods and Protocols* (ed. Werner,
587 A.) **1173**, 135–146 (2014).
- 588 48. Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence
589 of insertions, deletions and gene fusions. *Genome Biol.* **14**, R36 (2013).
- 590 49. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
591 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 592 50. Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect
593 evolution. *Science* **346**, 763–767 (2014).
- 594 51. Giribet, G. & Edgecombe, G. D. Reevaluating the Arthropod Tree of Life.
595 *Annu. Rev. Entomol.* **57**, 167–186 (2012).
- 596 52. Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian
597 phylogenetics with BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969–1973
598 (2012).
- 599 53. Hadfield, J. D. MCMC methods for multi-response generalized linear mixed
600 models: the MCMCglmm R package. *J. Stat. Softw.* **33**, 1–22 (2010).
- 601 54. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: A flexible trimmer for
602 Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).

603

604

605 **Supplementary Information** is linked to the online version of the paper at

606 www.nature.com/nature.

607 **Acknowledgements**

608 We thank A. McGregor, D. Leite, M. Akam, R. Jenner, R. Kilner, A. Duarte, C.
609 Jiggins, R. Wallbank, A. Bourke, T. Dalmay, N. Moran, K. Warchol, R. Callahan, G.
610 Farley, and T. Livdahl for providing arthropods. This research was supported by a
611 Leverhulme Research Project Grant (RPG-2016-210 to F.M.J., E.A.M. and P.S.), a
612 European Research Council grant (281668 Drosophila Infection to F.M.J.), a Medical
613 Research Council grant (MRC MC-A652-5PZ80 to P.S.) an Imperial College
614 Research Fellowship (to P.S.), Cancer Research UK (C13474/A18583,
615 C6946/A14492 to E.A.M.), the Wellcome Trust (104640/Z/14/Z, 092096/Z/10/Z to
616 E.A.M.), and an NIH R37 grant (GM62862 to P.D.Z.).

617 **Author contributions**

618 S.H.L. and K.A.Q. performed the experiments with assistance from Y.Y., M.T., L.F.,
619 S.A.S., P.P.S., R.C., C.G., I.G., D.H.C.; S.H.L., K.A.Q. & P.S. carried out
620 computational analysis; P.D.Z., E.A.M., P.S. & F.M.J. supervised the project; S.H.L.,
621 K.A.Q., P.D.Z., E.A.M., P.S. & F.M.J. wrote the manuscript.

622 **Competing financial interests**

623 The authors declare no competing financial interests.

624 **Materials and Correspondence**

625 Correspondence and requests for materials should be addressed to F.M.J.
626 (fmj1001@cam.ac.uk) or S.H.L. (sam.lewis@gen.cam.ac.uk).

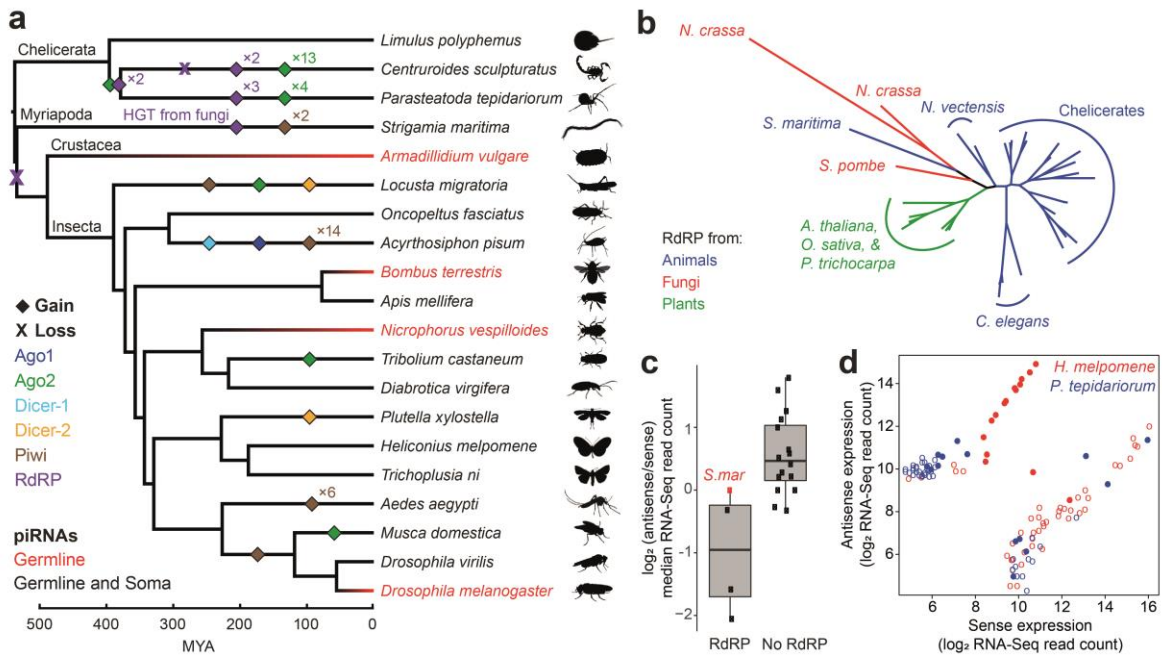
627

628 **Author affiliations**

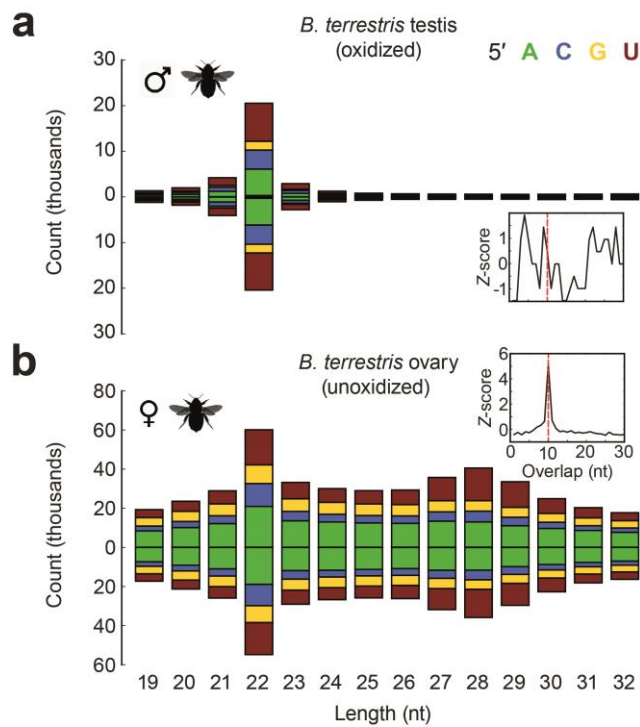
- 629 1. Department of Genetics, University of Cambridge, Downing Street, Cambridge,
630 CB2 3EH, United Kingdom
- 631 2. Howard Hughes Medical Institute, RNA Therapeutics Institute, University of
632 Massachusetts Medical School, 368 Plantation Street, Worcester, MA 01605,
633 USA
- 634 3. Wellcome Trust/Cancer Research UK Gurdon Institute, Cambridge, UK

- 635 4. Institut de Biologie de l'Ecole Normale Supérieure, CNRS, Inserm, ENS, PSL
636 Research University, Paris, France
- 637 5. Dept. Biomedical Sciences and Pathobiology, Virginia Maryland College of
638 Veterinary Medicine, 205 Duck Pond Drive, Virginia Tech, Blacksburg, VA, USA
- 639 6. University of Wisconsin-Madison, Department of Zoology, 352 Birge Hall, 430
640 Lincoln Drive, Madison, WI 53706, USA
- 641 7. Université de Poitiers, Laboratoire Ecologie et Biologie des Interactions, Equipe
642 Ecologie Evolution Symbiose, 5 Rue Albert Turpain, TSA 51106, 86073 Poitiers
643 Cedex 9, France
- 644 8. Laboratoire Evolution, Génomes, Comportement, Écologie, UMR 9191 CNRS,
645 UMR 247 IRD, Université Paris-Sud, Université Paris-Saclay, 91198 Gif-sur-
646 Yvette, France
- 647 9. School of Biological Sciences, University of East Anglia, Norwich Research Park,
648 Norwich NR4 7TJ, UK
- 649 10. MRC London Institute of Medical Sciences, Du Cane Road, London, W12 0NN,
650 UK
- 651 11. Institute for Clinical Sciences, Imperial College London, Du Cane Road, London,
652 W12 0NN, UK
653

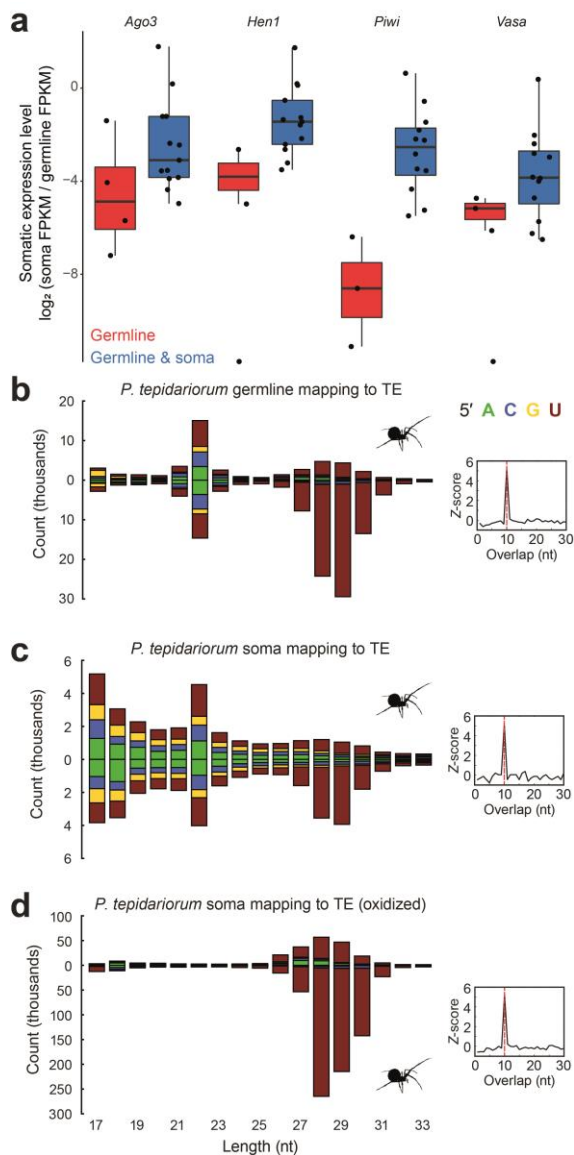
654 **Figures**



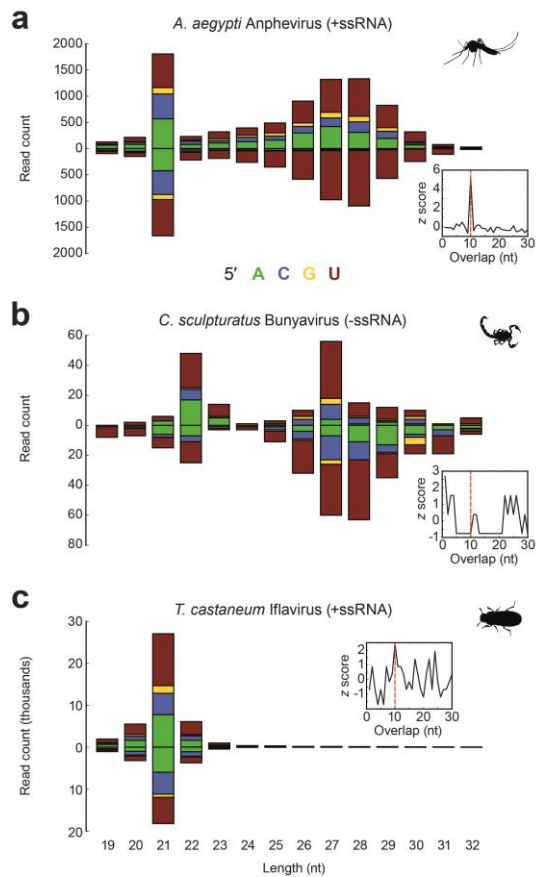
655 **Figure 1: Genes in small RNA pathways evolve rapidly throughout the**
 656 **arthropods. a**, The gain and loss of genes encoding the components of different
 657 sRNA pathways during arthropod evolution. Taxa with somatic piRNAs are shown in
 658 black, and the colour of the branches is a Bayesian reconstruction of whether
 659 somatic piRNAs were present. The posterior probability that the ancestral arthropod
 660 had somatic piRNAs is 0.9956. **b**, Phylogenetic analysis of RdRP genes from
 661 arthropods, other animals, plants and fungi. Note *S. maritima* is more closely related
 662 to fungal than animal RdRP (posterior probability at *N. crassa* - *S. maritima* node is
 663 1). **c**, The antisense enrichment (measured as $\log_2(\text{antisense/sense})$ median RNA-
 664 Seq read counts) for TEs that produce siRNAs. Species are classified by possession
 665 of an RdRP. Note *S. maritima* (red) lacks an animal RdRP. **d**, Among the 60 most
 666 highly expressed TEs in *H. melpomene* (no RdRP; red), the 15 TE that generate
 667 the most siRNAs (filled circles) have high rates of antisense transcription. In *P.*
 668 *tepidariorum* (six RdRPs; blue), TE with little antisense transcription generate
 669 siRNAs.
 670



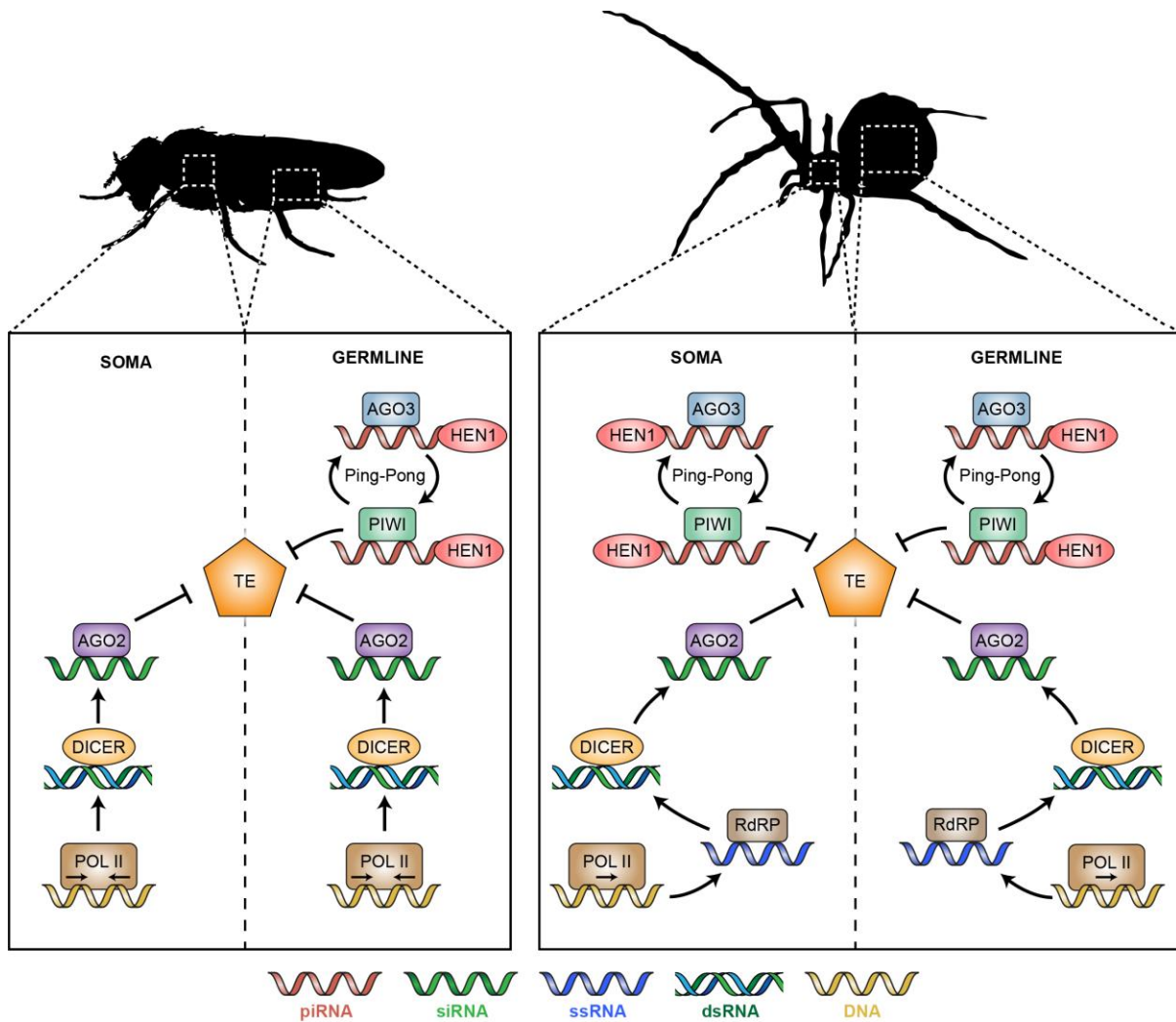
671 **Figure 2: piRNAs are absent in *B. terrestris* male germline.** The size and 5'
672 nucleotide of sRNAs from testis (a) and ovary (b). Reads derived from the sense
673 strand are shown above zero, antisense reads below. Plots show unique reads that
674 map to the genome (where the same sequence occurred more than once, all but one
675 read was eliminated). The inset shows the overlap between sense and antisense 25-
676 29nt sRNAs.
677



678 **Figure 3: Somatic piRNAs are widespread, and target TEs throughout the**
679 **arthropods. a**, Genes in the piRNA pathway have higher somatic expression in
680 species with somatic piRNAs. For species with multiple copies of a gene, the mean
681 scaled somatic expression level of each duplicate is displayed. The box shows the
682 median and interquartile range (IQR), and the bottom and top whiskers show the
683 range of points no further than 1.5×IQR away from the first and third quartiles
684 respectively. **b-d**, The size and 5' nucleotide of sRNAs mapping to the TEs from *P.*
685 *tepidariorum*, showing 10 bp overlap between sense and antisense 25–29nt sRNAs.
686 piRNAs targeting TEs are evident in the germline (**b**) and soma (**c**), and these
687 somatic piRNAs are resistant to sodium periodate oxidation, indicating that they are
688 3' methylated (**d**). Plots show all reads that map to the TEs.



689 **Figure 4: Virally-derived sRNAs in three arthropod species.** The size and 5'
690 nucleotide of sRNAs mapping to viral transcripts and genomes reconstructed from
691 RNA-Seq data. Virally-derived piRNAs are evident in *A. aegypti* (a) and *C.*
692 *sculpturatus* (b), and virally-derived siRNAs are found in *T. castaneum* (c). Only *A.*
693 *aegypti* shows the 10 bp overlap between sense and antisense 25–29 nt sRNAs that
694 is diagnostic of Ping-Pong amplification (insets). Reads derived from the sense
695 strand are shown above zero, antisense reads below.
696



697 **Figure 5: A model of the divergent sRNA pathways silencing TEs in different**
698 **arthropods.** Our data suggest that the mechanisms of sRNA pathways have
699 diverged in two key areas. In some lineages, the piRNA pathway is restricted to the
700 germline (e.g., flies), whereas in most others it is active in the soma and the germline
701 (e.g., spiders). Additionally, in some lineages (e.g., spiders), RdRP may synthesize
702 dsRNA from transcripts produced by RNA polymerase II, amplifying the siRNA
703 response.
704

705 **Extended Data Figure 1: Expression of RdRP in chelicerates and myriapods.**

706 The expression (fragments per kilobase per million reads) of *RdRP* homologues in
707 germline and somatic tissue of females (*L. polyphemus*, *C. sculpturatus* and *P.*
708 *tepidariorum*) or mixed-sex somatic tissue (*S. maritima*).

709 **Extended Data Figure 2: Genome-derived sRNAs in the female germline.** Size

710 distribution and 5' base composition of unoxidized sRNAs from the female germline
711 mapped to the entire genome. The Y-axis is the number of unique reads mapping
712 across the genome (where the same sequence occurred more than once, all but one
713 read was eliminated).

714 **Extended Data Figure 3: Proportion of somatic and germline piRNAs derived**

715 **from CDS, 5' UTR, 3' UTR, TE or elsewhere in the genome.** Total counts of sense
716 and antisense piRNAs for each feature type were extracted after mapping unique
717 sRNA sequences to the genome, and UTR counts were calculated after excluding
718 any UTR annotation regions that overlapped with TE annotations.

719 **Extended Data Figure 4: The correlation between TE expression and piRNA**

720 **abundance in the germline.** RNA-Seq and sRNA reads from the germline were
721 mapped to the genome, and read counts mapping to the sense and antisense
722 strands were totalled for each TE. Unexpressed TEs were screened out, and RNA-
723 Seq and sRNA counts for each TE family were calculated from the remaining TEs.

724 **Extended Data Figure 5: Genome-derived sRNAs in the male germline.** Size

725 distribution and 5' base composition of unoxidized sRNAs from the male germline
726 mapped to the entire genome. The Y-axis is the number of unique reads mapping
727 across the genome (where the same sequence occurred more than once, all but one
728 read was eliminated).

729 **Extended Data Figure 6: Genome-derived sRNAs in the *B. terrestris* germline.**

730 Size distribution and 5' base composition of sRNAs from *B. terrestris* testis, vas

731 deferens and ovary mapped to the entire genome. The Y axis is the number of
732 unique reads mapping across the genome (where the same sequence occurred
733 more than once, all but one read was eliminated).

734 **Extended Data Figure 7: Expression of piRNA pathway genes in *B. terrestris*.**

735 The expression (fragments per kilobase per million reads) of *Ago3*, *Hen1*, *piwi* and
736 *vasa* in *B. terrestris* ovary, testis and vas deferens.

737 **Extended Data Figure 8: Genome-derived sRNAs in the *A. mellifera* germline.**

738 Size distribution and 5' base composition of sRNAs from *A. mellifera* testis and ovary
739 mapped to the entire genome. The Y-axis is the number of unique reads mapping
740 across the genome (where the same sequence occurred more than once, all but one
741 read was eliminated). Samples in the right column were oxidised by sodium
742 periodate to exclude sRNA lacking 2'-O-methyl modification at their 3' end.

743 **Extended Data Figure 9: Genome-derived sRNAs in the soma of 20 arthropod**

744 **species.** Size distribution, 5' base composition, and strand distribution of sRNAs
745 from mixed-sex somatic tissue (*S. maritima*) or female somatic tissue, mapped to the
746 entire genome. The Y axis is the number of unique reads mapping across the
747 genome (where the same sequence occurred more than once, all but one read was
748 eliminated).

749 **Extended Data Figure 10: The TE content of 20 arthropod species. A**

750 comparison of the percentage of the genome made up of different TE classes for
751 species with somatic and germline piRNAs ('Germline and Soma') or just germline
752 piRNAs ('Germline'). The box shows the median and interquartile range (IQR), and
753 the bottom and top whiskers show the range of points no further than 1.5×IQR away
754 from the first and third quartiles respectively.

755 **Extended Data Figure 11: The phylogenetic distribution of somatic UTR-**

756 **derived piRNAs across 15 arthropod species. (a)** The phylogenetic distribution of

757 somatic UTR-derived piRNAs across 15 arthropods. For each species, we defined
758 the presence of UTR-derived piRNAs based on the presence of >200 unique 25–
759 29nt sequences with a 5' U nucleotide bias after oxidation treatment. **(b)**
760 Representative size distribution and 5' base composition of somatic UTR-derived
761 sRNAs in *C. sculpturatus*, *S. maritima*, and *P. xylostella*. sRNA derived from the
762 sense strand are above zero, antisense below. The Y axis is the number of unique
763 reads mapping to UTR (where the same sequence occurred more than once, all but
764 one read was eliminated). Samples were oxidised by sodium periodate to exclude
765 sRNAs lacking 2'-O-methyl modification at their 3' end.

766 **Extended Data Figure 12: Proportion of somatic piRNAs derived from CDS, 5'**
767 **UTR, and 3' UTR, scaled to the genome content of each feature.** Counts for each
768 feature type were extracted after mapping unique sRNA sequences to the genome,
769 scaled to the total number of unique sRNA sequences mapping to CDS, 5' UTR and
770 3' UTR, and then scaled to the total number of bases annotated in the genome.

771 **Extended Data Figure 13: Virus-derived small RNAs in the soma of arthropods.**
772 **(a)** Counts of viral contigs assembled from somatic RNA-Seq reads that did not map
773 to the host genome, together with whether they were targeted by somatic siRNAs
774 and/or piRNAs. A virus contig was categorised as being targeted by siRNAs and/or
775 piRNAs on the basis of visual inspection of the size distribution of all mapped
776 unoxidized reads. Note that a single virus may produce more than one contig. **(b)**
777 Unoxidized sRNA size distribution, 5' base composition, and strand distribution for
778 representative viruses targeted by siRNAs (piRNA viruses in main text). Reads
779 above zero are from the sense strand, below are antisense reads. The Y axis is the
780 number of all reads mapping across the contig. **(c)** Oxidized sRNA size distribution,
781 5' base composition, and strand distribution for piRNA-generating viruses.

782 **Extended Data Table 1: Pan-arthropod sequencing of sRNAs and total RNA in**
783 **the soma and germline of 20 arthropods.** The sex and tissue of samples used for

784 sRNA sequencing and RNA-Seq, and the read counts of the resulting sRNA
785 libraries. In *S. maritima* a mixed-sex somatic tissue sample was used.

786 **Extended Data Table 2: Metadata for genome assemblies and gene models.**

787 The accession numbers and URLs for the genome assemblies and gene model
788 annotations of each species.