

2018-01-29

## The genome of the Hi5 germ cell line from *Trichoplusia ni*, an agricultural pest and novel model for small RNA biology

Yu Fu  
*Boston University*

*Et al.*

### Let us know how access to this document benefits you.

Follow this and additional works at: <https://escholarship.umassmed.edu/oapubs>



Part of the [Biology Commons](#), [Computational Biology Commons](#), [Ecology and Evolutionary Biology Commons](#), [Genetics Commons](#), and the [Genomics Commons](#)

---

### Repository Citation

Fu Y, Yang Y, Zhang H, Farley G, Wang J, Quarles KA, Weng Z, Zamore PD. (2018). The genome of the Hi5 germ cell line from *Trichoplusia ni*, an agricultural pest and novel model for small RNA biology. Open Access Articles. <https://doi.org/10.7554/eLife.31628>. Retrieved from <https://escholarship.umassmed.edu/oapubs/3274>

Creative Commons License



This work is licensed under a [Creative Commons Attribution 4.0 License](#).

This material is brought to you by eScholarship@UMMS. It has been accepted for inclusion in Open Access Articles by an authorized administrator of eScholarship@UMMS. For more information, please contact [Lisa.Palmer@umassmed.edu](mailto:Lisa.Palmer@umassmed.edu).

1 **The genome of the Hi5 germ cell line**  
2 **from *Trichoplusia ni*, an agricultural pest**  
3 **and novel model for small RNA biology**

4 Yu Fu,<sup>1,2</sup> Yujing Yang,<sup>3†</sup> Han Zhang,<sup>3</sup> Gwen Farley,<sup>3</sup> Junling Wang,<sup>3</sup> Kaycee A.  
5 Quarles,<sup>3</sup> Zhiping Weng,<sup>2,4\*</sup> and Phillip D. Zamore<sup>3\*</sup>

6 <sup>1</sup>Bioinformatics Program, Boston University, 44 Cummington Mall, Boston, MA 02215,  
7 USA

8 <sup>2</sup>Program in Bioinformatics and Integrative Biology, University of Massachusetts Medical  
9 School, 368 Plantation Street, Worcester, MA 01605, USA

10 <sup>3</sup>RNA Therapeutics Institute and Howard Hughes Medical Institute, University of  
11 Massachusetts Medical School, 368 Plantation Street, Worcester, MA 01605, USA

12 <sup>4</sup>Department of Biochemistry and Molecular Pharmacology, University of Massachusetts  
13 Medical School, 368 Plantation Street, Worcester, MA 01605, USA

14 **\*For correspondence:** zhiping.weng@umassmed.edu (ZW),  
15 phillip.zamore@umassmed.edu (PDZ)

16 **†Present addresses:** Yujing Yang, Laboratory of Pathology, State Key Laboratory of  
17 Biotherapy and Department of Pathology, West China Hospital, West China Medical  
18 School, Sichuan University, Chengdu, China

19

20 **Abstract**

21 We report a draft assembly of the genome of Hi5 cells from the lepidopteran insect pest,  
22 *Trichoplusia ni*, assigning 90.6% of bases to one of 28 chromosomes and predicting  
23 14,037 protein-coding genes. Chemoreception and detoxification gene families reveal *T.*  
24 *ni*-specific gene expansions that may explain its widespread distribution and rapid  
25 adaptation to insecticides. Transcriptome and small RNA data from thorax, ovary, testis,  
26 and the germline-derived Hi5 cell line show distinct expression profiles for 295  
27 microRNA- and >393 piRNA-producing loci, as well as 39 genes encoding small RNA  
28 pathway proteins. Nearly all of the W chromosome is devoted to piRNA production, and  
29 *T. ni* siRNAs are not 2'-O-methylated. To enable use of Hi5 cells as a model system, we  
30 have established genome editing and single-cell cloning protocols. The *T. ni* genome  
31 provides insights into pest control and allows Hi5 cells to become a new tool for  
32 studying small RNAs ex vivo.

33

34

35

36

37

38

39

40 **Keywords:** genome assembly; piRNA; siRNA; miRNA; sex determination; cultured cell;  
41 High Five; CRISPR; cabbage looper; *Trichoplusia ni*; Lepidoptera

42

## 43 Introduction

44 Lepidoptera (moths and butterflies), one of the most species-rich orders of insects,  
45 comprises more than 170,000 known species (*Mallet, J, 2007; Chapman, AD, 2009*),  
46 including many agricultural pests. One of the largest lepidopteran families, the  
47 Noctuidae diverged over 100 million years ago (mya) from the Bombycidae—best-  
48 known for the silkworm, *Bombyx mori* (*Rainford, JL et al., 2014*). The Noctuidae family  
49 member cabbage looper (*Trichoplusia ni*) is a widely distributed generalist pest that  
50 feeds on cruciferous crops such as broccoli, cabbage, and cauliflower (*Capinera, J,*  
51 *2001*). *T. ni* has evolved resistance to the chemical insecticide  
52 Dichlorodiphenyltrichloroethane (DDT; [*McEwen, FL, Hervey, GER, 1956*]) and the  
53 biological insecticide *Bacillus thuringiensis* toxin (*Janmaat, AF, Myers, J, 2003*),  
54 rendering pest control increasingly difficult. A molecular understanding of insecticide  
55 resistance requires a high-quality *T. ni* genome and transcriptome.

56 Hi5 cells derive from *T. ni* ovarian germ cells (*Granados, RR et al., 1986;*  
57 *Granados, RR et al., 1994*). Hi5 cells are a mainstay of recombinant protein production  
58 using baculoviral vectors (*Wickham, TJ et al., 1992*) and hold promise for the  
59 commercial-scale production of recombinant adeno-associated virus for human gene  
60 therapy (*Kotin, RM, 2011; van Oers, MM et al., 2015*). Hi5 cells produce abundant  
61 microRNAs (miRNAs), small interfering RNAs (siRNAs), and PIWI-interacting RNAs  
62 (*Kawaoka, S et al., 2009*) (piRNAs), making them one of just a few cell lines suitable for  
63 the study of all three types of animal small RNAs. The most diverse class of small  
64 RNAs, piRNAs protect the genome of animal reproductive cells by silencing  
65 transposons (*Saito, K et al., 2006; Vagin, VV et al., 2006; Brennecke, J et al., 2007;*  
66 *Houwing, S et al., 2007; Aravin, AA et al., 2007; Kawaoka, S et al., 2008*). The piRNA  
67 pathway has been extensively studied in the dipteran insect *Drosophila melanogaster*  
68 (fruit fly), but no piRNA-producing, cultured cell lines exist for dipteran germline cells. *T.*  
69 *ni* Hi5 cells grow rapidly without added hemolymph (*Hink, WF, 1970*), are readily

70 transfected, and—unlike *B. mori* BmN4 cells (Iwanaga, M *et al.*, 2014), which also  
71 express germline piRNAs—remain homogeneously undifferentiated even after  
72 prolonged culture. In contrast to *B. mori*, no *T. ni* genome sequence is available, limiting  
73 the utility of Hi5 cells.

74 To further understand this agricultural pest and its Hi5 cell line, we combined  
75 diverse genomic sequencing data to assemble a chromosome-level, high-quality *T. ni*  
76 genome. Half the genome sequence resides in scaffolds >14.2 megabases (Mb), and  
77 >90% is assembled into 28 chromosome-length scaffolds. Automated gene prediction  
78 and subsequent manual curation, aided by extensive RNA-seq data, allowed us to  
79 examine gene orthology, gene families such as detoxification proteins, sex  
80 determination genes, and the miRNA, siRNA, and piRNA pathways. Our data allowed  
81 assembly of the gene-poor, repeat-rich W chromosome, which remarkably produces  
82 piRNAs across most of its length. To enable the use of cultured *T. ni* Hi5 cells as a  
83 novel insect model system, we established methods for efficient genome editing using  
84 the CRISPR/Cas9 system (Ran, FA *et al.*, 2013) as well as single-cell cloning. With  
85 these new tools, *T. ni* promises to become a powerful companion to flies to study gene  
86 expression, small RNA biogenesis and function, and mechanisms of insecticide  
87 resistance in vivo and in cultured cells.

## 88 **Results**

### 89 **Genome sequencing and assembly**

90 We combined Pacific Biosciences long reads and Illumina short reads (Figure 1A, Table  
91 1, and Materials and methods) to sequence genomic DNA from Hi5 cells and *T. ni* male  
92 and female pupae. The initial genome assembly from long reads (46.4× coverage with  
93 reads >5 kb) was polished using paired-end (172.7× coverage) and mate-pair reads  
94 (172.0× coverage) to generate 1,976 contigs spanning 368.2 megabases (Mb). Half of  
95 genomic bases reside in contigs >621.9 kb (N50). Hi-C long-range scaffolding (186.5×

96 coverage) produced 1,031 scaffolds (N50 = 14.2 Mb), with >90% of the sequences  
97 assembled into 28 major scaffolds. Karyotyping of metaphase Hi5 cells revealed that  
98 these cells have  $112 \pm 5$  chromosomes (Figure 1B, Figure 1—figure supplement 1).  
99 Because lepidopteran cell lines are typically tetraploid (*Hink WF, 1972*), we conclude  
100 that the ~368.2 Mb *T. ni* genome comprises 28 chromosomes: 26 autosomes plus W  
101 and Z sex chromosomes (see below).

102 To evaluate the completeness of the assembled *T. ni* genome, we compared it to  
103 the Arthropoda data set of the Benchmark of Universal Single-Copy Orthologs (*Simão,*  
104 *FA et al., 2015*) (BUSCO v3). The *T. ni* genome assembly captures 97.5% of these  
105 gene orthologs, more than either the silkworm (95.5%) or monarch butterfly (*D.*  
106 *plexippus*; 97.0%) genomes (Supplementary file 1A). All 79 ribosomal proteins  
107 conserved between mammals and *D. melanogaster* (*Yoshihama, M et al., 2002;*  
108 *Marygold, SJ et al., 2007*) have orthologs in *T. ni*, further evidence of the completeness  
109 of the genome assembly (Supplementary file 1B). Finally, a search for genes in the  
110 highly conserved nuclear oxidative phosphorylation (OXPHOS) pathway (*Porcelli, D et*  
111 *al., 2007*) uncovered *T. ni* orthologs for all known *D. melanogaster* OXPHOS genes  
112 (Supplementary file 1C).

113 The genomes of wild insect populations are typically highly heterogeneous,  
114 which poses a significant impediment to assembly (*Keeling, CI et al., 2013; You, M et*  
115 *al., 2013*). We were unable to generate an isogenic *T. ni* strain by inbreeding.  
116 Therefore, our *T. ni* sequence reflects the genome of Hi5 cells, not cabbage looper  
117 itself. Hi5 cells presumably derive from a single immortalized, germline founder cell,  
118 which should reduce genomic variation among the cell line's four sets of chromosomes.  
119 To test this supposition, we identified the sequence variants in the Hi5 genome. In total,  
120 we called variants at 165,370 genomic positions (0.0449% of the genome assembly),  
121 with 2,710 in predicted coding regions (0.0132% of coding sequence), indicating that  
122 the genome of Hi5 cells is fairly homogenous. For the majority (88.8%) of these

123 genomic positions (covering 0.0399% of the genome), only one copy of the  
124 chromosome has the variant allele while the other three chromosomal copies match the  
125 reference genome. We can make three conclusions. First, Hi5 cells originated from a  
126 single founder cell or a homogenous population of cells. Second, the founder cells were  
127 haploid. Third, most sequence variants were acquired after the original derivation of the  
128 line from *T. ni* eggs.

129 We also assembled de novo *T. ni* genomes using paired-end DNA-seq data  
130 obtained from male and female pupae, but the resulting assemblies are fragmented  
131 (scaffold N50  $\leq$  2.4 kb, Supplementary file 1D), likely due to the limitations of short-insert  
132 libraries and the high levels of heterozygosity commonly observed for genomes of wild  
133 insect populations (Keeling, *CI et al.*, 2013; You, *M et al.*, 2013). The animal genome  
134 contigs are highly concordant with the Hi5 genome, with  $\leq$  1.37% of animal contigs  
135 misassembled (Supplementary file 1D). Although we cannot determine scaffold-level  
136 differences between the animal and Hi5 cells, at the contig-level the Hi5 genome  
137 assembly is representative of the *T. ni* animal genome.

### 138 **Gene orthology**

139 We annotated 14,034 protein-coding genes in the *T. ni* genome (Supplementary file  
140 1E), similar to other Lepidoptera (Challis, *RJ et al.*, 2016). Analysis of the homology of  
141 *T. ni* genes to genes in 20 species that span the four common insect orders  
142 (Lepidoptera, Diptera, Coleoptera, Hymenoptera), non-insect arthropods, and mammals  
143 defines 30,448 orthology groups each containing orthologous proteins from two or more  
144 species (Hirose, *Y, Manley, JL*, 1997); 9,112 groups contain at least one *T. ni* gene. In  
145 all, 10,936 *T. ni* protein-coding genes are orthologous to at least one gene among the  
146 20 reference species (Figure 1C, Figure 1—figure supplement 2).

147 *T. ni* contains 2,287 Lepidoptera-specific orthology groups (*T. ni*, *B. mori*, *D.*  
148 *plexippus*, and *P. xylostella* [diamondback moth]). Far fewer orthology groups are

149 unique to Diptera (404), Coleoptera (371), or Hymenoptera (1,344), suggesting that the  
150 lepidopteran lifestyle requires more order-specific genes. The *T. ni* genome additionally  
151 contains 3,098 orphan protein-coding genes for which we could detect no orthologous  
152 sequences in the 20 reference species. Of these orphan genes, 14.5% are present as  
153 two or more copies in the genome (“in-paralogs”), suggesting they evolved recently.  
154 Some of these in-paralogs may have arisen by gene duplication after the divergence of  
155 *T. ni* and *B. mori* ~111 mya (Gaunt, MW, Miles, MA, 2002; Rota-Stabelli, O et al., 2013;  
156 Wheat, CW, Wahlberg, N, 2013; Rainford, JL et al., 2014).

### 157 **Opsins**

158 The ability of insects to respond to light is crucial to their survival. Opsins, members of  
159 the G-protein-coupled receptor superfamily, play important roles in vision. Covalently  
160 bound to light-sensing chromophores, opsins absorb photons and activate the  
161 downstream visual transduction cascade (Terakita, A, 2005). The *T. ni* genome  
162 encodes ultraviolet, blue, and long-wavelength opsins. Thus, this nocturnal insect  
163 retains the full repertoire of insect opsins and has color vision (Zimyanin, VL et al.,  
164 2008) (Figure 1—figure supplement 3). *T. ni* also encodes an ortholog of the non-visual  
165 Rh7 opsin, which is found in a variety of insects (Initiative, IGG, 2014; Futahashi, R et  
166 al., 2015). In the *D. melanogaster* brain, Rh7 opsin participates in the entrainment of  
167 circadian rhythms by sensing violet light (Ni, JD et al., 2017). *T. ni* also encodes an  
168 ortholog of the vertebrate-like opsin, pteropsin, which was first detected in the honeybee  
169 (*A. mellifera*) brain and is found widely among insects except for *Drosophilid* flies  
170 (Velarde, RA et al., 2005).

### 171 **Sex determination**

172 Understanding the *T. ni* sex-determination pathway holds promise for engineering  
173 sterile animals for pest management. ZW and ZO chromosome systems determine sex  
174 in lepidopterans: males are ZZ and females are either ZW or ZO (Traut, W et al., 2007).



175 To determine which system *T. ni* uses and to identify which contigs belong to the sex  
176 chromosomes, we sequenced genomic DNA from male and female pupae and  
177 calculated the male:female coverage ratio for each contig. We found that 175  
178 presumably Z-linked contigs (20.0 Mb) had approximately twice the coverage in male  
179 compared to female DNA (median male:female ratio = 1.92; Figure 2A, Figure 2—figure  
180 supplement 1A). Another 276 contigs (11.1 Mb) had low coverage in males (median  
181 male:female ratio = 0.111), suggesting they are W-linked. We conclude that sex is  
182 determined in *T. ni* by a ZW system in which males are homogametic (ZZ) and females  
183 are heterogametic (ZW).

184 For some lepidopteran species, dosage compensation has been reported to  
185 equalize Z-linked transcript abundance between ZW females and ZZ males in the soma,  
186 while other species show higher expression of Z-linked genes in males (*Walters, JR et*  
187 *al., 2015; Gu, L et al., 2017*). In the soma, *T. ni* compensates for Z chromosome  
188 dosage: transcripts from Z-linked genes are approximately equal in male and female  
189 thoraces ( $Z \approx ZZ$ , Figure 2B). In theory, somatic dosage compensation could reflect  
190 increased transcription of the single female Z chromosome, reduced transcription of  
191 both male Z chromosomes, or silencing of one of the two male Z chromosomes.

192 To distinguish among these possibilities, we compared the abundance of Z-  
193 linked and autosomal transcripts (Z/AA in female and ZZ/AA in male, Figure 2—figure  
194 supplement 1B and 1C). Z-linked transcripts in the male thorax are expressed at lower  
195 levels than autosomal transcripts, but not as low as half ( $ZZ \approx 70\% AA$ ). These data  
196 support a dosage compensation mechanism that decreases transcription from each Z  
197 chromosome in the *T. ni* male soma, but does not fully equalize Z-linked transcript  
198 levels between the sexes ( $Z \approx ZZ \approx 70\% AA$ ). In contrast, *T. ni* lacks germline dosage  
199 compensation: in the ovary, Z-linked transcript abundance is half that of autosomal  
200 transcripts ( $Z \approx 50\% AA$ ), whereas in testis, Z-linked and autosomal transcripts have  
201 equal abundance ( $ZZ \approx AA$ ). We conclude that *T. ni*, like *B. mori* (*Walters, JR,*

202 *Hardcastle, TJ, 2011*), *Cydia pomonella* (*Gu, L et al., 2017*), and *Heliconius* butterflies  
203 (*Walters, JR et al., 2015*), compensates for Z chromosome dosage in the soma by  
204 reducing gene expression in males, but does not decrease Z-linked gene expression in  
205 germline tissues.

206 Little is known about lepidopteran W chromosomes. The W chromosome is not  
207 included in the genome assembly of *Manduca sexta* (*Kanost, MR et al., 2016*) or *B.*  
208 *mori* (*The, ISG, 2008*), and earlier efforts to assemble the silkworm W resulted in  
209 fragmented sequences containing transposons (*Abe, H et al., 2005; Abe, H et al., 2008;*  
210 *Kawaoka, S et al., 2011*). The monarch genome scaffold continuity (N50 = 0.207 Mb  
211 versus N50 = 14.2 Mb for *T. ni*; [*Zhan, S et al., 2011*]) is insufficient to permit assembly  
212 of a W chromosome. Our genome assembly includes the 2.92 Mb *T. ni* W chromosome  
213 comprising 32 contigs (contig N50=101 kb). In *T. ni*, W-linked contigs have higher  
214 repeat content, lower gene density, and lower transcriptional activity than autosomal or  
215 Z-linked contigs (Figure 2B). Other lepidopteran W chromosomes are similarly enriched  
216 in repeats and depleted of genes (*Abe, H et al., 2005; Fuková, I et al., 2005; Traut, W et*  
217 *al., 2007*).

218 A search for *T. ni* genes that are homologous to insect sex determination  
219 pathway genes detected *doublesex (dsx)*, *masculinizer (masc)*, *vitellogenin*, *transformer*  
220 *2*, *intersex*, *sex lethal*, *ovarian tumor*, *ovo*, and *sans fille*. *T. ni* males produce a four-  
221 exon isoform of *dsx*, while females generate a six-exon *dsx* isoform (Figure 2—figure  
222 supplement 1D). The Lepidoptera-specific gene *masc* encodes a CCCH zinc finger  
223 protein. *masc* is associated with the expression of the sex-specific isoforms of *dsx* in  
224 lepidopterans, including silkworm (*Katsuma, S et al., 2015*). As in *B. mori*, *T. ni masc*  
225 lies next to the *scap* gene, supporting our annotation of *T. ni masc*. Lepidopteran *masc*  
226 genes are rapidly diverging and have low sequence identity with one another (30.1%).  
227 Figure 2C shows the multiple sequence alignment of the CCCH zinc finger domain of  
228 Masc proteins from several lepidopteran species.

## 229 **Telomeres and centromeres**

230 Like many non-dipteran insects, *T. ni* has a single telomerase gene and telomeres  
231 containing TTAGG repeats (*Sahara, K et al., 1999*). We found 40 (TTAGG)<sub>n</sub> stretches  
232 longer than 100 nt (mean ± S.D. = 600 ± 800 nt), nine at and 31 near contig boundaries  
233 (Supplementary file 1F; distance between (TTAGG)<sub>n</sub> and contig boundary = 5,000 ±  
234 6,000 nt for the 40 stretches), indicating that our assembly captures the sequences of  
235 many telomeres. More than half (59%) of the sequences flanking the (TTAGG)<sub>n</sub> repeats  
236 are transposons, and ~49% of these belong to the non-long-terminal-repeat LINE/R1  
237 family (Supplementary file 1G). These telomeric and subtelomeric characteristics of *T.*  
238 *ni* resemble those of *B. mori* (*Fujiwara, H et al., 2005*).

239 Lepidopteran chromosomes generally lack a coherent, monocentric centromere  
240 and are instead holocentric or diffuse (*Labbé, R et al., 2011*), and the silkworm,  
241 monarch butterfly, and diamondback moth genomes do not encode CenH3, a protein  
242 associated with monocentric chromosomes. The *T. ni* genome similarly does not  
243 contain a gene for CenH3, suggesting that its chromosomes are also holocentric.

## 244 **CpG content and DNA methylation**

245 The *T. ni* genome is 35.6% GC, slightly less than *B. mori* (37.3%). The distributions of  
246 observed/expected CpG ratios in genes and across the genome (Figure 2—supplement  
247 2A) reveal that *T. ni* is similar to other lepidopterans (silkworm, monarch butterfly,  
248 diamondback moth) and a coleopteran species (red flour beetle, *T. castaneum*), but  
249 different from honeybee and fruit fly. The honeybee genome has a high CpG content in  
250 genes and exhibits a bimodal CpG distribution across the genome as a whole; the fruit  
251 fly genome is uniformly depleted of CpG dinucleotides. The differences in CpG patterns  
252 reflect the presence of both the DNMT1 and DNMT3 DNA methyltransferases in the  
253 honeybee, the absence of either in fruit fly, and the presence of only DNMT1 in *T. ni*, *B.*  
254 *mori*, *D. plexippus*, *P. xylostella*, and *T. castaneum*. Thus, like many other insects, the

255 *T. ni* genome likely has low levels of DNA methylation (Xiang, H et al., 2010; Glastad,  
256 KM et al., 2011).

### 257 **Transposons and repeats**

258 The *T. ni* genome contains 75.3 Mb of identifiable repeat elements (20.5% of the  
259 assembly), covering 458 repeat families (Figure 2—figure supplement 2B,  
260 Supplementary file 1H). With this level of repeat content, *T. ni* fits well with the positive  
261 correlation between genome size and repeat content among lepidopteran genomes  
262 (Figure 2—figure supplement 2C).

263 The DNA transposon piggyBac was originally isolated from a *T. ni* cell line  
264 (Fraser, MJ et al., 1983) and transposes effectively in a variety of species (Lobo, N et  
265 al., 1999; Bonin, CP, Mann, RS, 2004; Wang, W et al., 2008). We identified 262 copies  
266 of piggyBac in the Hi5 cell genome assembly. The family divergence rate of piggyBac is  
267 ~0.17%, substantially lower than other transposon families in the genome  
268 (Supplementary file 1I provides divergence rates for all transposon families). Among the  
269 individual piggyBac elements in the *T. ni* genome, 71 are specific to Hi5 cells.  
270 Compared to the 191 piggyBac insertions shared between *T. ni* and Hi5 cells  
271 (divergence rate = 0.22%), the Hi5 cell-specific elements are more highly conserved  
272 (divergence rate = 0.04%). We conclude that the piggyBac transposon entered the *T. ni*  
273 genome more recently than other transposons and, likely driven by the presence of  
274 many active piggyBac elements, expanded further during the immortalization of Hi5  
275 cells in culture.

### 276 **microRNAs**

277 microRNAs (miRNAs) are ~22 nt non-coding RNAs that regulate mRNA stability and  
278 translation (He, L, Hannon, GJ, 2004; Gao, G et al., 2005). In insects, miRNA targets  
279 function in metamorphosis, reproduction, diapause, and other pathways of insect  
280 physiology and development (Lucas, K, Raikhel, AS, 2013). To characterize the *T. ni*

281 miRNA pathway, we sequenced RNA and small RNA from ovary, testis, thorax, and Hi5  
282 cells. Then, we manually identified miRNA biogenesis genes such as *dcr-1*, *pasha*,  
283 *drosha*, and *ago2* (Supplementary file 2A) and computationally predicted 295 miRNA  
284 genes (Figure 3, Supplementary file 3A and Supplementary file 4), including 77  
285 conserved, 31 Lepidoptera-specific, and 187 novel, *T. ni*-specific miRNAs.

286 In thorax, 222 of 270 miRNAs had comparable abundance in males and females  
287 ( $\leq 2$ -fold difference or false discovery rate [FDR]  $\geq 0.1$ ; Figure 3A). Of the 48 miRNAs  
288 having significantly different abundances in female and male thorax ( $> 2$ -fold difference  
289 and FDR  $< 0.1$ ; Figure 3A), miR-1a, let-7, and miR-278 were highly abundant ( $> 1000$   
290 parts per million [ppm]) in either female or male thorax. miR-1a, a miRNA thought to be  
291 expressed in all animal muscle, was the most abundant miRNA in thorax in both sexes,  
292 but was 2.2-fold more abundant in males. miR-1 was previously shown to regulate  
293 muscle development in fruit flies (*Sokol, NS, Ambros, V, 2005*) and to increase when  
294 locusts transition from solitary to swarming (*Wei, Y et al., 2009*). *T. ni* let-7, which has  
295 the same mature miRNA sequence as its *D. melanogaster*, *C. elegans*, and mammalian  
296 counterparts (*Lagos-Quintana, M et al., 2001*) was also more abundant in males,  
297 whereas miR-278 was 2.6-fold more abundant in females. let-7 may act in sex-specific  
298 pathways in metamorphosis (*Caygill, EE, Johnston, LA, 2008*), whereas miR-278 may  
299 play a sex-specific role in regulating energy homeostasis (*Teleman, AA et al., 2006*).

300 A subset of less well conserved miRNAs was also differentially expressed  
301 between male and female thorax. In general, poorly conserved miRNAs were less  
302 abundant: the median expression level for conserved miRNAs was 316 ppm, but only  
303 161 ppm for Lepidoptera-specific and 4.22 ppm for *T. ni*-specific miRNAs. However,  
304 mir-2767, a Lepidoptera-specific miRNA, and three *T. ni*-specific miRNAs (mir-novel1,  
305 mir-novel4, mir-novel11) were both abundant ( $> 1000$  ppm) and differentially expressed  
306 in males and female thorax. We speculate that these recently evolved miRNAs may  
307 prove useful as targets for pest management.

308 Ovary, testis, and Hi5 cells have distinct miRNA expression profiles. We  
309 analyzed the expression patterns of the 44 most abundant miRNAs (Figure 3B and 3C),  
310 which explain 90% of miRNA reads in a tissue or cell line. Thirteen were expressed in  
311 ovaries, testes, and Hi5 cells. Of these 13, 11 were significantly more abundant in testis,  
312 5 in ovary, and 3 in Hi5 cells (Figure 3B), suggesting that these miRNAs have important  
313 tissue- or cell-type-specific roles. miR-31 and miR-375, highly expressed in *T. ni* testis,  
314 are both mammalian tumor suppressors (*Creighton, CJ et al., 2010; Kinoshita, T et al.,*  
315 *2012*). miR-989, the most abundant miRNA in *T. ni* ovaries, plays an important role in  
316 border cell migration during *Drosophila* oogenesis (*Kugler, J-M et al., 2013*). miR-10, a  
317 miRNA in the Hox gene cluster, was preferentially expressed in Hi5 cells; its orthologs  
318 have been implicated in development and cancer (*Lund, AH, 2009*), suggesting miR-10  
319 played a role in the immortalization of the germline cells from which Hi5 cells derive.

## 320 siRNAs

321 siRNAs, typically 20–22 nt long, regulate gene expression, defend against viral  
322 infection, and silence transposons (*Agrawal, N et al., 2003; van Rij, RP et al., 2006;*  
323 *Sanchez-Vargas, I et al., 2009; Tyler, DM et al., 2008; Tam, OH et al., 2008; Zambon,*  
324 *RA et al., 2006; Chung, WJ et al., 2008; Okamura, K et al., 2008b; Czech, B et al.,*  
325 *2008; Okamura, K et al., 2008b; Flynt, A et al., 2009*). They are processed by Dicer  
326 from double-stranded RNAs or hairpins into short double-stranded fragments bearing  
327 two-nucleotide, overhanging 3' ends, which are subsequently loaded into Argonaute  
328 proteins (*Bernstein, E et al., 2001; Elbashir, SM et al., 2001; Siomi, H, Siomi, MC,*  
329 *2009*). siRNAs require extensive sequence complementarity to their targets to elicit  
330 Argonaute-catalyzed target cleavage.

## 331 Endogenous siRNAs from transposons and *cis*-NATs

332 Endogenous siRNAs (endo-siRNAs) can derive from transposon RNAs, *cis*-natural  
333 antisense transcripts (*cis*-NATs), and long hairpin RNAs (*Czech, B et al., 2008;*

334 Ghildiyal, M et al., 2008; Okamura, K et al., 2008a; Chung, WJ et al., 2008; Kawamura,  
335 Y et al., 2008; Okamura, K et al., 2008a; Tam, OH et al., 2008; Watanabe, T et al.,  
336 2008) (hpRNAs). In *T. ni* ovary, testis, thorax, and Hi5 cells, 20.7–52.4% of siRNAs map  
337 to transposons, suggesting *T. ni* endogenous siRNAs suppress transposons in both the  
338 soma and the germline. Among the non-transposon siRNAs, <4.6% map to predicted  
339 hairpins, while 11.6–31.3% siRNAs map to *cis*-NATs (Supplementary file 3B).

#### 340 Exogenous siRNAs against a virus

341 Hi5 cells are latently infected with a positive-sense, bipartite alphanodavirus, TNCL  
342 virus (Li et al., 2007, #97210; Miller and Ball, 2012, #84273) (Tn5 Cell Line virus). We  
343 asked if TNCL virus RNA is present in our *T. ni* samples and whether the RNAi pathway  
344 provides anti-viral defense via TNCL virus-derived siRNAs. We detected no viral RNA in  
345 the *T. ni* ovary, testis, or thorax transcriptome, but both TNCL virus RNA1 (5,010  
346 fragments per kilobase of transcript per million mapped reads [FPKM]) and RNA2  
347 (8,280 FPKM) were readily found in the Hi5 transcriptome (Figure 4A). To test whether  
348 Hi5 cells mount an RNAi defense to TNCL virus infection, we mapped small RNA-seq  
349 reads that were not mappable to the *T. ni* genome to the two TNCL virus genomic  
350 segments. TNCL virus-mapping small RNAs showed a median length of 21 nt (modal  
351 length = 20 nt; Figure 4A), typical for siRNAs, suggesting that the Hi5 RNAi pathway  
352 actively combats the virus. The TNCL virus-mapping small RNAs bear the two-  
353 nucleotide, 3' overhanging ends that are the hallmark of siRNAs (Figure 4B) (Elbashir,  
354 SM et al., 2001; Elbashir, SM et al., 2001; Elbashir, SM et al., 2001). Moreover, the  
355 phased pattern of TNCL virus-mapping siRNAs suggests they are made one-after-  
356 another starting at the end of a dsRNA molecule: the distance between siRNA 5' ends  
357 shows a periodicity of 20 nt, the length of a typical TNCL virus-mapping siRNA (Figure  
358 4C). In *D. melanogaster*, Dicer-2 processively produces siRNAs, using ATP energy to  
359 translocate along a dsRNA molecule (Cenik, ES et al., 2011). The phasing of anti-viral

360 siRNAs in Hi5 cells suggests that *T. ni* Dicer-2 similarly generates multiple siRNAs from  
361 each molecule of dsRNA before dissociating.

362 In addition to siRNAs, the TNCL-mapping small RNAs include some 23–32 nt  
363 RNAs. These are unlikely to be anti-viral piRNAs, because they lack the characteristic  
364 first-nucleotide uridine bias and show no significant ping-pong signal ( $Z$ -score =  
365  $-0.491$ ). We conclude that Hi5 cells do not use piRNAs for viral defense.

### 366 **Lepidopteran siRNAs are not 2'-O-methylated**

367 The discovery that the 3' ends of *D. melanogaster* siRNAs, but not miRNAs, are 2'-O-  
368 methylated (*Pelisson, A et al., 2007*) led to the idea that insects in general methylate  
369 both siRNAs and piRNAs. Resistance to oxidation by NaIO<sub>4</sub> is the hallmark of 3'  
370 terminal, 2'-O-methylation, and the enrichment of a small RNA in a high-throughput  
371 sequencing library prepared from NaIO<sub>4</sub>-treated RNA suggests 2'-O-methylation.  
372 Conversely, depletion of small RNAs, such as miRNAs, from such an oxidized RNA  
373 library is strong evidence for unmodified 2',3' vicinal hydroxyl groups. Surprisingly,  
374 TNCL virus-mapping siRNAs were 130-fold depleted from our oxidized small RNA-seq  
375 library (22.0 ppm) compared to the unoxidized library (2,870 ppm), suggesting that they  
376 are unmethylated. Sequencing of oxidized and unoxidized small RNA from *T. ni* ovary,  
377 testis, and thorax detected 20–22 nt peaks in unoxidized libraries; such peaks were  
378 absent from oxidized libraries (Figure 4D), suggesting that *T. ni* genome-mapping,  
379 endogenous siRNAs also lack 2'-O-methylation. We conclude that both *T. ni* exo- and  
380 endo-siRNAs are not 2'-O-methyl modified.

381 Are siRNAs unmethylated in other Lepidopteran species? We sequenced  
382 oxidized and unoxidized small RNAs from two additional Lepidoptera: *P. xylostella* and  
383 *B. mori*. Like *T. ni*, siRNAs from these Lepidoptera were abundant in libraries prepared  
384 from unoxidized small RNA but depleted from oxidized libraries (Figure 4—figure  
385 supplement 1A). The ratio of siRNAs in the oxidized library to siRNAs in the



386 corresponding unoxidized library (ox/unox) provides a measure of siRNA 2',3'  
387 modification. For *D. melanogaster* siRNAs, the median ox/unox ratio was 1.00, whereas  
388 the three Lepidoptera species had median ox/unox ratios between 0.17 and 0.22  
389 (Figure 4E), indicating their siRNAs were depleted from oxidized libraries and therefore  
390 bear unmodified 2',3' hydroxyl groups. We conclude that the last common ancestor of *T.*  
391 *ni*, *B. mori*, and *P. xylostella*, which diverged 170 mya, lacked the ability to 2'-*O*-  
392 methylate siRNA 3' ends. We do not currently know whether the last common ancestor  
393 of Lepidoptera lost the capacity to methylate siRNAs or if some or all members of  
394 Diptera, the sister order of Lepidoptera, acquired this function, which is catalyzed by the  
395 piRNA-methylating enzyme Hen1 (*Saito, K et al., 2007; Horwich, MD et al., 2007;*  
396 *Kirino, Y, Mourelatos, Z, 2007*).

397       Terminal 2' methylation of *D. melanogaster* siRNAs is thought to protect them  
398 from non-templated nucleotide addition (tailing), 3'-to-5' trimming, and wholesale  
399 degradation (*Ameres, SL et al., 2010*). Since *T. ni* siRNAs lack a 2'-*O*-methyl group at  
400 their 3' ends, we first asked if we could observe frequent trimming by examining shorter  
401 TNCL-mapping siRNA (18–19 nt). These siRNAs account for 1.05% of all TNCL-  
402 mapping siRNAs. They did not possess the typical siRNA one-after-another pattern ( $Z_1$   
403 = -0.674,  $p = 0.500$ ), yet more than 97.5% of these were prefixes of longer, phased  
404 siRNAs, indicating that these were trimmed siRNAs. We conclude that TNCL siRNA  
405 trimming is rare in Hi5 cells. We next asked whether *T. ni* and other lepidopteran  
406 siRNAs have higher frequencies of tailing. Despite the lack of 2'-*O*-methylation, most  
407 TNCL virus siRNAs were not tailed: just 6.69% of all virus-mapping small RNA reads  
408 contained 3' non-templated nucleotides (Figure 4—figure supplement 1B). Among the 3'  
409 non-templated nucleotides, the most frequent addition was one or more uridines  
410 (49.6%) as observed previously for miRNAs and siRNAs in other animals (*Ameres, SL*  
411 *et al., 2010; Chou, MT et al., 2015*). Endogenous siRNA tailing frequencies for the  
412 lepidopterans *T. ni* (10.2%, ovary), *B. mori* (5.97%, eggs), and *P. xylostella* (8.58%,

413 ovary) were also similar to *D. melanogaster* (6.71%, ovary). We speculate that  
414 lepidopterans have other mechanisms to maintain siRNA stability or that trimming and  
415 tailing in lepidopterans are less efficient than in flies.

416 siRNAs are non-randomly loaded into Argonaute proteins: the guide strand, the  
417 strand with the more weakly base paired 5' end, is favored for loading (*Khvorova, A et*  
418 *al., 2003; Schwarz, DS et al., 2003*); the disfavored passenger strand is destroyed.  
419 Thus, loading skews the abundance of the two siRNA strands. To test if non-methylated  
420 siRNAs are loaded into Argonaute, we computationally paired single-stranded siRNAs  
421 that compose an siRNA duplex bearing two-nucleotide overhanging 3' ends and  
422 calculated the relative abundance of the two siRNA strands. For TNCL-mapping  
423 siRNAs, 72.3% of siRNA duplexes had guide/passenger strand ratios  $\geq 2$  (median =  
424 3.90; mean = 10.2; Figure 4—figure supplement 2). Among genome-mapping 20–22 nt  
425 small RNAs 78.5% of duplexes had guide/passenger strand ratios  $\geq 2$  (median 5.44;  
426 average 56.2). We conclude that the majority of exogenous and endogenous siRNAs  
427 are loaded, presumably into Ago2.

#### 428 piRNAs

429 In animals, piRNAs, ~23–32 nt long, protect the germline genome by suppressing the  
430 transcription or accumulation of transposon and repetitive RNA (*Girard, A et al., 2006;*  
431 *Lau, NC et al., 2006; Vagin, VV et al., 2006; Brennecke, J et al., 2007; Aravin, AA et al.,*  
432 *2007*). In *D. melanogaster*, dedicated transposon-rich loci (piRNA clusters) give rise to  
433 piRNA precursor transcripts, which are processed into piRNAs loaded into one of three  
434 PIWI proteins, Piwi, Aubergine (Aub), or Argonaute3 (Ago3). Piwi acts in the nucleus to  
435 direct tri-methylation of histone H3 on lysine 9 on transposon and repetitive genomic  
436 sequences (*Sienski, G et al., 2012; Le Thomas, A et al., 2014; Le Thomas, A et al.,*  
437 *2014*). In fly cytoplasm, piRNAs guide the Piwi paralog Aub to cleave transposon  
438 mRNAs. The mRNA cleavage products can then produce more piRNAs, which are

439 loaded into Ago3. In turn, these sense piRNAs direct Ago3 to cleave transcripts from  
440 piRNA clusters, generating additional piRNAs bound to Aub. The resulting “Ping-Pong”  
441 feed-forward loop both amplifies piRNAs and represses transposon activity (*Brennecke,*  
442 *J et al., 2007; Gunawardane, LS et al., 2007*). Finally, Ago3 cleavage not only produces  
443 Aub-bound piRNAs, but also initiates the production of Piwi-bound, phased piRNAs that  
444 diversify the piRNA pool (*Mohn, F et al., 2015; Han, BW et al., 2015*).

#### 445 piRNA pathway proteins

446 The *T. ni* genome contains a full repertoire of genes encoding piRNA pathway proteins  
447 (Supplementary file 2B). These genes were expressed in both germline and somatic  
448 tissues, but were higher in ovary, testis, and Hi5 cells compared to thorax (median  
449 ratios: ovary/thorax = 14.2, testis/thorax = 2.9, and Hi5/thorax = 4.9; Figure 5A).  
450 Expression of piRNA pathway genes in the Hi5 cell line suggests that it recapitulates the  
451 germline piRNA pathway. Although most *T. ni* piRNA pathway genes correspond  
452 directly to their *D. melanogaster* orthologs, *T. ni* encodes only two PIWI proteins, TnPiwi  
453 and TnAgo3. The fly proteins Aub and Piwi are paralogs that arose from a single  
454 ancestral PIWI protein after the divergence of flies and mosquitos (*Lewis, SH et al.,*  
455 *2016*). We do not yet know whether TnPiwi functions more like *Drosophila* Aub or Piwi.  
456 In *D. melanogaster*, piRNA clusters—the genomic sources of most transposon-silencing  
457 germline piRNAs—are marked by the proteins Rhino, Cutoff, and Deadlock, which allow  
458 transcription of these heterochromatic loci (*Klattenhoff, C et al., 2009; Pane, A et al.,*  
459 *2011; Mohn, F et al., 2014; Zhang, Z et al., 2014*). *T. ni* lacks detectable Rhino, Cutoff,  
460 and Deadlock orthologs. In fact, this trio of proteins is poorly conserved, and the  
461 mechanism by which they mark fly piRNA source loci may be unique to Drosophilids. In  
462 this regard, *T. ni* likely provides a more universal insect model for the mechanisms by  
463 which germ cells distinguish piRNA precursor RNAs from other protein-coding and non-  
464 coding transcripts.

## 465 piRNA cluster architecture

466 In both the germline and the soma, *T. ni* piRNAs originate from discrete genomic loci.  
467 To define these piRNA source loci, we employed an expectation-maximization algorithm  
468 that resolves piRNAs mapping to multiple genomic locations. Applying this method to  
469 multiple small RNA-seq datasets, we defined piRNA-producing loci comprising 10.7 Mb  
470 (348 clusters) in ovary, 3.1 Mb (79 clusters) in testis, 3.0 Mb (71 clusters) in Hi5 cells,  
471 and 2.4 Mb (65 clusters) in thorax (Figure 5B). For each tissue or cell-type, these 393  
472 clusters explain >70% of uniquely mapped piRNAs and >70% of all piRNAs when using  
473 expectation-maximization mapping. A core set of piRNA-producing loci comprising 1.5  
474 Mb is active in both germline and somatic tissues.

475 *T. ni* piRNA clusters vary substantially in size and expression level. In ovary, half  
476 the bases in piRNA clusters are in just 67 loci, with a median length of 53 kb. Among  
477 these, five span >200 kb, while the smallest is just 38 kb. The most productive piRNA  
478 source is a 264 kb locus on chromosome 13 (Figure 5—figure supplement 1); 7.8% of  
479 uniquely mapped piRNAs—50,000 distinct piRNA sequences—reside in this locus.  
480 Collectively, the top 20 ovary piRNA loci explain half the uniquely mapped piRNAs, yet  
481 constitute only 0.7% of the genome. Globally, 61.9% of bases in piRNA clusters are  
482 repetitive, and 74.5% transposon-mapping piRNAs are antisense, suggesting that *T. ni*  
483 uses antisense piRNAs to suppress transposon transcripts.

484 In the fly ovary germline, most piRNA clusters generate precursor RNAs from  
485 both DNA strands. These dual-strand clusters fuel the ‘Ping-Pong’ amplification cycle  
486 (Brennecke, J et al., 2007; Gunawardane, LS et al., 2007). Other fly piRNA clusters,  
487 such as the paradigmatic *flamenco* gene (Prud’homme, N et al., 1995; Brennecke, J et  
488 al., 2007; Pelisson, A et al., 2007; Malone, CD et al., 2009; Goriaux, C et al., 2014) are  
489 transcribed from one strand only and are organized to generate antisense piRNAs  
490 directly, without further Ping-Pong amplification (Malone, CD et al., 2009). These uni-  
491 strand clusters are the only sources of piRNAs in the follicle cells, somatic cells that

492 support fly oocyte development and express only a single PIWI protein, Piwi (*Malone,*  
493 *CD et al., 2009*).

494 The *T. ni* genome contains both dual- and uni-strand piRNA clusters. In ovary, 62  
495 of 348 piRNA-producing loci are dual-strand (Watson/Crick >0.5 or Watson/Crick < 2).  
496 These loci produce 35.9% of uniquely mapped piRNAs and 22.8% of all piRNAs; 71.6%  
497 of transposon-mapping piRNA reads from these loci are antisense. The remaining 286  
498 uni-strand loci account for 54.8% of uniquely mapped piRNAs and 36.7% of all piRNAs.  
499 Most piRNAs (74.8% of reads) from uni-strand clusters are antisense to transposons,  
500 the orientation required for repressing transposon mRNA accumulation. At least part of  
501 the piRNA antisense bias reflects positive selection for antisense insertions in uni-strand  
502 clusters: 57.1% of transposon insertions—79.7% of transposon-mapping nucleotides—  
503 are opposite the direction of piRNA precursor transcription, significantly different from  
504 dual-strand clusters, in which transposons are inserted randomly: 49.5% of transposon  
505 insertions in dual-strand clusters are in the antisense direction (Figure 5—figure  
506 supplement 2A). For one 77 kb uni-strand cluster on chromosome 20, 99.0% of piRNA  
507 reads (96% of piRNA sequences) that can be uniquely assigned are from the Crick  
508 strand, while 67.6% of transposon insertions and 79.7% of transposon-mapping  
509 nucleotides at this locus lie on the Watson strand.

#### 510 **Nearly the entire W chromosome produces piRNAs**

511 The largest ovary cluster is a 462 kb W-linked region, consistent with our finding that the  
512 W chromosome is a major source of piRNAs (Figure 5B and 5C and Figure 5—figure  
513 supplement 2B). Our data likely underestimates the length of this large piRNA cluster,  
514 as it is difficult to resolve reads mapping to its flanking regions: 70.8% of bases in the  
515 flanking regions do not permit piRNAs to map uniquely to the genome. In fact, 85.1% of  
516 the sequences between clusters on the W chromosome are not uniquely mappable.

517 These gaps appear to reflect low mappability and not boundaries between discrete  
518 clusters. We propose that the W chromosome itself is a giant piRNA cluster.

519 To further test this idea, we identified piRNA reads that uniquely map to one  
520 location among all contigs and measured their abundance per kilobase of the genome.  
521 W-linked contigs had a median piRNA abundance of 14.4 RPKM in ovaries, 379-fold  
522 higher than the median of all autosomal and Z-linked contigs, consistent with the view  
523 that almost the entire W chromosome produces piRNAs. In *B. mori* females, a plurality  
524 of piRNAs come from the W chromosome: ovary-enriched piRNAs often map to W-  
525 linked sequences, but not autosomes (*Kawaoka, S et al., 2011*). Similarly, for *T. ni*,  
526 27.2% of uniquely mapping ovary piRNAs derive from W-linked sequences, even  
527 though these contigs compose only 2.8% of the genome (Figure 5C). The W  
528 chromosome may produce more piRNAs than our estimate, as the unassembled  
529 repetitive portions of the W chromosome likely also produce piRNAs. Thus, the entire W  
530 chromosome is a major source of piRNAs in *T. ni* ovaries (Figure 5B). To our  
531 knowledge, the *T. ni* W chromosome is the first example of an entire chromosome  
532 devoted to piRNA production.

533 To determine if there are W-linked regions devoid of piRNAs, we mapped all  
534 piRNAs to the W-linked contigs and found that 11.0% of the W-linked bases were not  
535 covered by any piRNAs, indicating at least part of the W chromosome does not produce  
536 any piRNAs. Next, we manually inspected 74 putative W-linked protein-coding genes  
537 and nine putative W-linked miRNAs. All nine W-linked miRNAs (Figure 5B,  
538 Supplementary file 1J) are *T. ni*-specific, and small RNAs mapping to these predicted  
539 miRNA loci showed significant ping-pong signature ( $Z$ -score = 14.2,  $p = 1.81 \times 10^{-45}$ ),  
540 suggesting that these are likely piRNAs, not authentic miRNAs. For the putative protein-  
541 coding genes, we categorized them into orphan genes (no homologs found),  
542 transposons (good homology to transposons), uncharacterized/hypothetical proteins,  
543 and potential protein-coding genes with homology to the NCBI non-redundant protein

544 sequences. We then asked whether piRNAs were produced from these genes (Figure  
545 5—figure supplement 2C). Among W-linked genes, those with transposon homology on  
546 average produced the most piRNAs (44.9 median ppm) whereas those with homology  
547 to annotated genes produced the fewest (9.81 median ppm). Some putative genes  
548 (such as TNI001015 and TNI005339) produced no piRNAs at all. We conclude that  
549 although some W-linked loci do not produce piRNAs, nearly the entire W chromosome  
550 produces piRNAs.

551 In contrast to the W chromosome, *T. ni* autosomes and the Z chromosome  
552 produce piRNAs from discrete loci—63 autosomal and 11 Z-linked contigs had piRNA  
553 levels >10 rpk. Few piRNAs are produced outside of these loci: for example, the  
554 median piRNA level across all autosomal and Z-linked contigs was ~0 in ovaries (Figure  
555 5—figure supplement 2B).

### 556 **Expression of piRNA clusters**

557 In the *T. ni* germline, piRNA production from individual clusters varies widely, but the  
558 same five piRNA clusters produce the most piRNAs in ovary (34.9% of piRNAs), testis  
559 (49.3%), and Hi5 cells (44.0%), suggesting that they serve as master loci for germline  
560 transposon silencing. Other piRNA clusters show tissue-specific expression, with the W  
561 chromosome producing more piRNAs in ovary than in Hi5 cells, and three Z-linked  
562 clusters producing many more piRNAs in testis than in ovary (15.0–24.7 times more),  
563 even after accounting for the absence of dosage compensation in germline tissues  
564 (Figure 6—figure supplement 1A).

565 Hi5 cells are female, yet many piRNA-producing regions of the W chromosome  
566 that are active in the ovary produce few piRNAs in Hi5 cells (Figure 6—figure  
567 supplement 1A). We do not know whether this reflects a reorganization of cluster  
568 expression upon Hi5 cell immortalization or if Hi5 cells correspond to a specific germ  
569 cell type that is underrepresented in whole ovaries. At least 40 loci produce piRNAs in

570 Hi5 cells but not in ovaries. Comparison of DNA-seq data from *T. ni* and Hi5 identified  
571 74 transposon insertions in 12 of the Hi5-specific piRNA clusters. Older transposons  
572 have more time to undergo sequence drift from the consensus sequence of the  
573 corresponding transposon family. The 74 Hi5-specific transposon insertions, which  
574 include both DNA and LTR transposons, had significantly lower divergence rates than  
575 those common to ovary and Hi5 cells (Figure 6A), consistent with the idea that recent  
576 transposition events generated the novel piRNA clusters in Hi5 cells. We conclude that  
577 the Hi5-specific piRNA-producing loci are quite young, suggesting that *T. ni* and  
578 perhaps other lepidopterans can readily generate novel piRNA clusters.

579 piRNA clusters active in thorax occupy ~0.57% of the genome and explain 86.8%  
580 of uniquely mapped somatic piRNAs in females and 89.5% in males. More than 90% of  
581 bases in clusters expressed in thorax are shared with clusters expressed in ovary  
582 (Figure 6—figure supplement 1B). Such broadly expressed clusters explain 83.7% of  
583 uniquely mapping piRNAs in female thorax and 86.1% in male thorax. Thus, the  
584 majority of piRNAs in the *T. ni* soma come from clusters that are also active in the  
585 germline. In general, autosomal piRNA cluster expression is similar between female and  
586 male thorax, but 12 clusters are differentially expressed between male and female  
587 thorax. Of these, nine are W-linked clusters that produce significantly more piRNAs in  
588 female than in male thorax (Figure 6B).

### 589 piRNA precursor transcripts are rarely spliced

590 In *D. melanogaster*, Rhino suppresses splicing of piRNA precursors transcribed from  
591 dual-strand piRNA clusters (Mohn, F et al., 2014; Zhang, Z et al., 2014). Fly uni-strand  
592 piRNA clusters do not bind Rhino and behave like canonical RNA polymerase II  
593 transcribed genes (Brennecke, J et al., 2007; Goriaux, C et al., 2014). Although *T. ni*  
594 has no *rhino* ortholog, its piRNA precursor RNAs are rarely spliced as observed for  
595 clusters in flies. We identified splicing events in our RNA-seq data, requiring  $\geq 10$  reads



596 that map across exon-exon junctions and a minimum splicing entropy of 2 to exclude  
597 PCR duplicates (Graveley, BR et al., 2011). This approach detected just 27 splice sites  
598 among all piRNA precursor transcripts from ovary, testis, thorax, and Hi5 piRNA  
599 clusters (Figure 6C). Of these 27 splice sites, 19 fall in uni-strand piRNA clusters. We  
600 conclude that, as in flies, transcripts from *T. ni* dual-strand piRNAs clusters are rarely if  
601 ever spliced. Unlike flies (Goriaux, C et al., 2014), RNA from *T. ni* uni-strand piRNA  
602 clusters also undergoes splicing infrequently.

603         The absence of piRNA precursor splicing in dual-strand piRNA clusters could  
604 reflect an active suppression of the splicing machinery or a lack of splice sites. To  
605 distinguish between these two mechanisms, we predicted gene models for piRNA-  
606 producing loci, employing the same parameters used for protein-coding genes. For  
607 piRNA clusters, this approach generated 1,332 gene models encoding polypeptides  
608 >200 amino acids. These models comprise 2,544 introns with consensus splicing  
609 signals (Figure 6—figure supplement 1C). Notably, ~90% of these predicted gene  
610 models had high sequence similarity to transposon consensus sequences (BLAST e-  
611 value<math>10^{-10}</math>), indicating that many transposons in piRNA clusters have intact splice  
612 sites. We conclude that piRNA precursors contain splice sites, but their use is actively  
613 suppressed.

614         To measure splicing efficiency, we calculated the ratio of spliced to unspliced  
615 reads for each predicted splice site in the piRNA clusters. High-confidence splice sites  
616 in protein-coding genes outside piRNA clusters served as a control. Compared to the  
617 control set of genes, splicing efficiency in piRNA loci was 9.67-fold lower in ovary, 2.41-  
618 fold lower in testis, 3.23-fold lower in thorax, and 17.0-fold lower in Hi5 cells (Figure 6D),  
619 showing that *T. ni* piRNA precursor transcripts are rarely and inefficiently spliced. To  
620 test whether uni- and dual-strand piRNA cluster transcripts are differentially spliced in *T.*  
621 *ni*, we evaluated the experimentally supported splice sites from Hi5, ovary, testis, and  
622 thorax collectively. Dual-strand cluster transcripts had 1.71-fold lower splicing efficiency

623 compared to uni-strand clusters (Figure 6D). Thus, *T. ni* suppresses splicing of dual-  
624 and uni-strand piRNA cluster transcripts by a mechanism distinct from the Rhino-  
625 dependent pathway in *D. melanogaster*. That this novel splicing suppression pathway is  
626 active in Hi5 cells should facilitate its molecular dissection.

## 627 **Genome-editing and single-cell cloning of Hi5 cells**

628 The study of arthropod piRNAs has been limited both by a lack of suitable cultured cell  
629 models and by the dominance of *D. melanogaster* as a piRNA model for arthropods  
630 generally. Although Vasa-positive *D. melanogaster* ovarian cells have been isolated and  
631 cultured (Niki, Y et al., 2006), no dipteran germ cell line is currently available. *D.*  
632 *melanogaster* somatic OSS, OSC and Kc167 cells produce piRNAs, but lack key  
633 features of the canonical germline pathway (Lau, NC et al., 2009; Saito, K et al., 2009;  
634 Vrettos, N et al., 2017). In addition to Hi5 cells, lepidopteran cell lines from *Spodoptera*  
635 *frugiperda* (Sf9) and *B. mori* (BmN4) produce germline piRNAs (Kawaoka, S et al.,  
636 2009). The *S. frugiperda* genome remains a draft with 37,243 scaffolds and an N50 of  
637 53.7 kb (Kakumani, PK et al., 2014). Currently, the BmN4 cell line is the only ex vivo  
638 model for invertebrate germline piRNA biogenesis and function. The *B. mori* genome  
639 sequence currently comprises 43,463 scaffolds with an N50 of 4.01 Mb (The, ISG,  
640 2008). Unfortunately, BmN4 cells readily differentiate into two morphologically distinct  
641 cell types (Iwanaga, M et al., 2014). Although genome editing with Cas9 has been  
642 demonstrated in BmN4 cells (Zhu, L et al., 2015), no protocols for cloning individual,  
643 genome-modified BmN4 cells have been reported (Mon, H et al., 2004; Kawaoka, S et  
644 al., 2009; Honda, S et al., 2013). In contrast, Hi5 cells are cultured using commercially  
645 available media, readily transfected, and, we report here, efficiently engineered with  
646 Cas9 and grown from single cells into clonal lines.

647 The bacterial DNA nuclease Cas9, targeted by a single guide RNA (sgRNA),  
648 enables rapid and efficient genome editing in worms, flies, and mice, as well as in a

649 variety of cultured animal cell lines (*Jinek, M et al., 2012; Barrangou, R, Horvath, P,*  
650 *2017; Komor, AC et al., 2017*). The site-specific double-strand DNA breaks catalyzed by  
651 Cas9 can be repaired by error-prone non-homologous end joining (NHEJ), disrupting a  
652 protein-coding sequence or, when two sgRNAs are used, deleting a region of genomic  
653 DNA. Alternatively, homology-directed repair (HDR) using an exogenous DNA template  
654 allows the introduction of novel sequences, including fluorescent proteins or epitope  
655 tags, as well as point mutations in individual genes (*Cong, L et al., 2013*).

656 As a proof-of-concept, we used Cas9 and two sgRNAs to generate a deletion in  
657 the piRNA pathway gene *TnPiwi*. The two sgRNAs, whose target sites lie 881 bp apart  
658 (Figure 7A), were transcribed *in vitro*, loaded into purified, recombinant Cas9 protein,  
659 and the resulting sgRNA/Cas9 ribonucleoprotein complexes (RNPs) transfected into Hi5  
660 cells. PCR of genomic DNA isolated 48 h later was used to detect alterations in the  
661 *TnPiwi* gene. A novel PCR product, ~900 bp smaller than the product amplified using  
662 DNA from control cells, indicated that the desired deletion had been created (Figure  
663 7B). Sanger sequencing of the PCR products confirmed deletion of 881–896 bp from  
664 the *TnPiwi* gene. The presence of indels—short deletions and non-templated nucleotide  
665 additions—at the deletion junction is consistent with a Cas9-mediated dsDNA break  
666 having been repaired by NHEJ (Figure 7A). We note that these cells still contain at least  
667 one wild-type copy of *TnPiwi*. We have not yet obtained cells in which all four copies of  
668 *TnPiwi* are disrupted, perhaps because in the absence of Piwi, Hi5 cells are inviable.

669 To test whether an exogenous donor DNA could facilitate the site-specific  
670 incorporation of protein tag sequences into Hi5 genome, we designed two sgRNAs with  
671 target sites ~90 bp apart, flanking the *vasa* start codon (Figure 7C). As a donor, we  
672 used a single-stranded DNA (ssDNA) encoding EGFP and an HA epitope tag flanked by  
673 genomic sequences 787 bp upstream and 768 bp downstream of the *vasa* start codon  
674 (Figure 7C). Cas9 and the two sgRNAs were cotransfected with the ssDNA donor, and,  
675 one week later, EGFP-positive cells were detected by fluorescence microscopy. PCR

676 amplification of the targeted region using genomic DNA from EGFP-expressing cells  
677 confirmed integration of EGFP and the HA tag into the *vasa* gene (Figure 7D). Sanger  
678 sequencing further confirmed integration of EGFP and the HA tag in-frame with the  
679 *vasa* open reading frame (Supplemental file 9).

680 To establish a clonal line from the EGFP-HA-tagged Vasa-expressing cells,  
681 individual EGFP-positive cells were isolated by FACS and cultured on selectively  
682 permeable filters above a feeder layer of wild-type Hi5 cells (Figure 8A). Growth of the  
683 genome-modified single cells required live Hi5 feeder cells—conditioned media did not  
684 suffice—presumably because the feeder cells provide short-lived growth factors or other  
685 trophic molecules. Single EGFP-positive clones developed one month after seeding and  
686 could be further grown without feeder cells as a clonally derived cell line (Figure 8B).

#### 687 **Hi5 cell Vasa is present in a nuage-like, perinuclear structure**

688 In the germline of *D. melanogaster* and other species, components of the piRNA  
689 biogenesis pathway, including Vasa, Aub, Ago3, and multiple Tudor-domain proteins,  
690 localize to a perinuclear structure called nuage (Eddy, EM, 1976; Findley, SD et al.,  
691 2003; Lim, AK, Kai, T, 2007; Li, C et al., 2009; Liu, L et al., 2011; Webster, A et al.,  
692 2015). Vasa, a germline-specific nuage component, is widely used as a marker for  
693 nuage. In BmN4 cells, transiently transfected Vasa localizes to a perinuclear structure  
694 resembling nuage (Xiol, J et al., 2012; Patil, AA et al., 2017). To determine whether  
695 nuage-like structures are present in Hi5 cells, we examined Vasa localization in the Hi5  
696 cells in which the endogenous *vasa* gene was engineered to fuse EGFP and an HA  
697 epitope tag to the Vasa amino-terminus. We used two different immunostaining  
698 strategies to detect the EGFP-HA-Vasa fusion protein: a mouse monoclonal anti-GFP  
699 antibody and a rabbit monoclonal anti-HA antibody. GFP and HA colocalized in a  
700 perinuclear structure, consistent with Vasa localizing to nuage in Hi5 cells (Figure 8C).

## 701 **Discussion**

702 Using Hi5 cells, we have sequenced and assembled the genome of the cabbage looper,  
703 *T. ni*, a common and destructive agricultural pest that feeds on many plants of economic  
704 importance. Examination of the *T. ni* genome and transcriptome reveals the expansion  
705 of detoxification-related gene families (Table 1 and Supplementary file 6), many  
706 members of which are implicated in insecticide resistance and are potential targets of  
707 pest control. The *T. ni* genome should enable study of the genetic diversity and  
708 population structure of this generalist pest, which adapts to different environmental  
709 niches worldwide. Moreover, as the sister order of Diptera, Lepidoptera like *T. ni* provide  
710 a counterpoint for the well-studied insect model *D. melanogaster*.

711 The use of Hi-C sequencing was an essential step in assembling the final  
712 368.2 Mb *T. ni* genome into high-quality, chromosome-length scaffolds. The integration  
713 of long reads, short reads, and Hi-C provides a rapid and efficient paradigm for  
714 generating chromosome-level assemblies of other animal genomes. This strategy  
715 assembled the gene-poor, repeat-rich *T. ni* W chromosome, which is, to our knowledge,  
716 the first chromosome-level sequence of a lepidopteran W chromosome. Our analysis of  
717 autosomal, Z-linked, and W-linked transcripts provides insights into lepidopteran dosage  
718 compensation and sex determination. Our data show that *T. ni* compensates for Z  
719 chromosome dosage in the soma by reducing transcription of both Z homologs in  
720 males, but Z dosage is uncompensated in the germline.

721 In addition to long RNAs, we characterized miRNAs, siRNAs, and piRNAs in *T. ni*  
722 gonads, soma, and cultured Hi5 cells. miRNAs are widely expressed in *T. ni* tissues,  
723 providing examples of germline-enriched and somatic miRNAs, as well as highly  
724 conserved, lepidopteran-specific, and novel *T. ni* miRNAs. Like flies, *T. ni* possess  
725 siRNAs that map to transposons, *cis*-NATs and hpRNAs. Unexpectedly, *T. ni* siRNAs—  
726 and likely all lepidopteran siRNAs—lack a 2'-*O*-methyl modification at their 3' ends,  
727 unlike siRNAs in *D. melanogaster*. Consistent with siRNA production by a processive

728 Dicer-2 enzyme, Hi5 cells produce phased siRNAs from the RNA genome of a latent  
729 alphanodavirus. The commonalities and differences between *T. ni* and *D. melanogaster*  
730 small RNA pathways will help identify both deeply conserved and rapidly evolving  
731 components.

732         A major motivation for sequencing the *T. ni* genome was the establishment of a  
733 tractable cell culture model for studying small RNAs, especially piRNAs. We believe that  
734 our genome assembly and gene-editing protocols will enable the use of *T. ni* Hi5 cells to  
735 advance our understanding of how piRNA precursors are defined, made into piRNAs  
736 and act to silence transposons in the germline. Hi5 cells express essentially all known  
737 piRNA pathway genes except those specific to Drosophilids. Furthermore, *T. ni* Vasa  
738 localizes to a perinuclear, nuage-like structure in Hi5 cells, making them suitable for  
739 studying the assembly of the subcellular structures thought to organize piRNA  
740 biogenesis. We have defined genomic piRNA-producing loci in Hi5 cells, as well as in  
741 the soma, testis, and ovary. The most productive piRNA clusters are shared among  
742 ovary, testis, and Hi5 cells. In addition, Hi5 cells contain novel piRNA clusters not found  
743 in the moth itself, suggesting that the process of establishing new piRNA-producing loci  
744 can be recapitulated by experimental manipulation of Hi5 cells.

745         As in *D. melanogaster*, splicing of *T. ni* piRNA precursor transcripts is efficiently  
746 suppressed, yet *T. ni* lacks paralogs of the proteins implicated in splicing suppression in  
747 flies. The ability to study the mechanisms by which piRNA clusters form and how  
748 precursor RNAs are transcribed, exported, and marked for piRNA production in *T. ni*  
749 promises to reveal both conserved and lepidopteran-specific features of this pathway.  
750 Notably, the W chromosome not only is a major piRNA source, but also produces  
751 piRNAs from almost its entirety. Future studies are needed to determine whether this is  
752 a common feature of W chromosomes in Lepidoptera and other insects.

753         The establishment of procedures for genome editing and single-cell cloning of  
754 Hi5 cells, combined with the *T. ni* genome sequence, make this germ cell line a

755 powerful tool to study RNA and protein function ex vivo. Our strategy combines  
756 transfection of pre-assembled Cas9/sgRNA complexes with single clone isolation using  
757 a selectable marker (e.g., EGFP) and feeder cells physically separated from the  
758 engineered cells. Compared with nucleic acid-based delivery of Cas9, transfection of  
759 Cas9 RNP minimizes the off-target mutations caused by prolonged Cas9 expression  
760 and eliminates the risk of integration of sgRNA or Cas9 sequences into the genome  
761 (*Lin, S et al., 2014; Kim, S et al., 2014*). Compared to plasmid donors (*Yu, Z et al.,*  
762 *2014; Ge, DT et al., 2016*), ssDNA homology donors similarly reduce the chance of  
763 introducing exogenous sequences at unintended genomic sites. Techniques for  
764 injecting the embryos of other lepidopteran species have already been established  
765 (*Wang, Y et al., 2013; Takasu, Y et al., 2014; Zhang, Z et al., 2015*). In principle, Cas9  
766 RNP injected into cabbage looper embryos could be used to generate genetically  
767 modified *T. ni* strains both to explore lepidopteran biology and to implement novel  
768 strategies for safe and effective pest control.

## 769 **Acknowledgements**

770 We thank members of the Weng and Zamore laboratories for helpful discussions and  
771 comments on the manuscript; UMass Deep Sequencing Core for Pacific Biosciences  
772 sequencing; Zdenka Matijasevic for sharing the karyotyping protocol. This work was  
773 supported in part by National Institutes of Health grants R37GM062862 to PDZ and  
774 HD078253 to ZW.

## 775 **References**

- 776 Abe, H et al. 2008. Identification of the female-determining region of the W chromosome  
777 in *Bombyx mori*. *Genetica* **133**: 269–282. 10.1007/s10709-007-9210-1
- 778 Abe, H, Mita, K, Yasukochi, Y, Oshiki, T, Shimada, T. 2005. Retrotransposable  
779 elements on the W chromosome of the silkworm, *Bombyx mori*. *Cytogenet*

780           *Genome Res* **110**: 144–151. 10.1159/000084946

781 Agrawal, N, Dasaradhi, PVN, Mohmmed, A, Malhotra, P, Bhatnagar, RK, Mukherjee,  
782           SK. 2003. RNA Interference: Biology, Mechanism, and Applications. *Microbiology*  
783           *and Molecular Biology Reviews* **67**: 657–685. 10.1128/MMBR.67.4.657-685.2003

784 Ai, J et al. 2011. Genome-wide analysis of cytochrome P450 monooxygenase genes in  
785           the silkworm, *Bombyx mori*. *Gene* **480**: 42–50. 10.1016/j.gene.2011.03.002

786 Altschul, SF, Gish, W, Miller, W, Myers, EW, Lipman, DJ. 1990. Basic local alignment  
787           search tool. *Journal of Molecular Biology* **215**: 403–410. 10.1016/S0022-  
788           2836(05)80360-2

789 Ameres, SL et al. 2010. Target RNA-directed trimming and tailing of small silencing  
790           RNAs. *Science* **328**: 1534–1539. 10.1126/science.1187058

791 Aravin, AA, Sachidanandam, R, Girard, A, Fejes-Toth, K, Hannon, GJ. 2007.  
792           Developmentally regulated piRNA clusters implicate MILI in transposon control.  
793           *Science* **316**: 744–747.

794 Attrill, H et al. 2016. FlyBase: establishing a Gene Group resource for *Drosophila*  
795           melanogaster. *Nucleic Acids Res* **44**: D786–92. 10.1093/nar/gkv1046

796 Barrangou, R, Horvath, P. 2017. A decade of discovery: CRISPR functions and  
797           applications. *Nature microbiology* **2**: 17092. 10.1038/nmicrobiol.2017.92

798 Belton, J-M, McCord, RP, Gibcus, JH, Naumova, N, Zhan, Y, Dekker, J. 2012. Hi-C: A  
799           comprehensive technique to capture the conformation of genomes. *Methods* **58**:  
800           268–276. 10.1016/j.ymeth.2012.05.001

801 Benton, R, Vannice, KS, Gomez-Diaz, C, Vosshall, LB. 2009. Variant Ionotropic  
802           Glutamate Receptors as Chemosensory Receptors in *Drosophila*. *Cell* **136**: 149–  
803           162. 10.1016/j.cell.2008.12.001

804 Bernstein, E, Caudy, AA, Hammond, SM, Hannon, GJ. 2001. Role for a bidentate  
805           ribonuclease in the initiation step of RNA interference. *Nature* **409**: 363–366.  
806           10.1038/35053110



807 Bernt, M et al. 2013. MITOS: improved de novo metazoan mitochondrial genome  
808 annotation. *Mol Phylogenet Evol* **69**: 313–319. 10.1016/j.ympev.2012.08.023

809 Bonin, CP, Mann, RS. 2004. A piggyBac Transposon Gene Trap for the Analysis of  
810 Gene Expression and Function in *Drosophila*. *Genetics* **167**: 1801–1811.  
811 10.1534/genetics.104.027557

812 Brennecke, J et al. 2007. Discrete Small RNA-Generating Loci as Master Regulators of  
813 Transposon Activity in *Drosophila*. *Cell* **128**: 1089–1103.  
814 10.1016/j.cell.2007.01.043

815 Burton, JN, Adey, A, Patwardhan, RP, Qiu, R, Kitzman, JO, Shendure, J. 2013.  
816 Chromosome-scale scaffolding of de novo genome assemblies based on  
817 chromatin interactions. *Nature Biotechnology*  
818 *Nat Biotech* **31**: 1119–1125. 10.1038/nbt.2727

819 Campbell, MS, Holt, C, Moore, B, Yandell, M. 2014. Genome Annotation and Curation  
820 Using MAKER and MAKER-P. *Current protocols in bioinformatics* **48**: 4.11.1–  
821 4.11.39. 10.1002/0471250953.bi0411s48

822 Capella-Gutiérrez, S, Silla-Martínez, JM, Gabaldón, T. 2009. trimAl: a tool for  
823 automated alignment trimming in large-scale phylogenetic analyses.  
824 *Bioinformatics* **25**: 1972–1973. 10.1093/bioinformatics/btp348

825 Capinera, J. (2001). *Handbook of Vegetable Pests* (Gulf Professional Publishing).

826 Castresana, J. 2000. Selection of Conserved Blocks from Multiple Alignments for Their  
827 Use in *Phylogenetic Analysis*. *Molecular Biology and Evolution*  
828 *Mol Biol Evol* **17**: 540–552. 10.1093/oxfordjournals.molbev.a026334

829 Caygill, EE, Johnston, LA. 2008. Temporal Regulation of Metamorphic Processes in  
830 *Drosophila* by the let-7 and miR-125 Heterochronic MicroRNAs. *Current Biology*  
831 **18**: 943–950. 10.1016/j.cub.2008.06.020

832 Cenik, ES et al. 2011. Phosphate and R2D2 restrict the substrate specificity of Dicer-2,  
833 an ATP-driven ribonuclease. *Molecular cell* **42**: 172–184.

834 10.1016/j.molcel.2011.03.002

835 Challis, RJ, Kumar, S, Dasmahapatra, KKK, Jiggins, CD, Blaxter, M. 2016. Lepbase:  
836 the Lepidopteran genome database. *bioRxiv* 056994. 10.1101/056994

837 Chapman, AD. (2009). Numbers of Living Species in Australia and the World: Report for  
838 the Australian Biological Resources Study (Canberra: Department of the  
839 Environment, Water, Heritage and the Arts).

840 Chou, MT, Han, BW, Hsiao, CP, Zamore, PD, Weng, Z, Hung, JH. 2015. Tailor: a  
841 computational framework for detecting non-templated tailing of small silencing  
842 RNAs. *Nucleic Acids Res* 10.1093/nar/gkv537

843 Chung, WJ, Okamura, K, Martin, R, Lai, EC. 2008. Endogenous RNA Interference  
844 Provides a Somatic Defense against *Drosophila* Transposons. *Curr Biol* **18**: 795–  
845 802. 10.1016/j.cub.2008.05.006

846 Conesa, A, Götz, S, García-Gómez, JM, Terol, J, Talón, M, Robles, M. 2005. Blast2GO:  
847 a universal tool for annotation, visualization and analysis in functional genomics  
848 research. *Bioinformatics* **21**: 3674–3676. 10.1093/bioinformatics/bti610

849 Cong, L et al. 2013. Multiplex Genome Engineering Using CRISPR/Cas Systems.  
850 *Science* **339**: 819–823. 10.1126/science.1231143

851 Creighton, CJ et al. 2010. Molecular profiling uncovers a p53-associated role for  
852 microRNA-31 in inhibiting the proliferation of serous ovarian carcinomas and  
853 other cancers. *Cancer research* **70**: 1906–1915. 10.1158/0008-5472.CAN-09-  
854 3875

855 Croset, V et al. 2010. Ancient Protostome Origin of Chemosensory Ionotropic  
856 Glutamate Receptors and the Evolution of Insect Taste and Olfaction. *PLOS*  
857 *Genetics* **6**: e1001064. 10.1371/journal.pgen.1001064

858 Czech, B et al. 2008. An endogenous small interfering RNA pathway in *Drosophila*.  
859 *Nature* **453**: 798–802. 10.1038/nature07007

860 DePristo, MA, Banks, E, Poplin, R, Garimella..., KV. 2011. A framework for variation

861 discovery and genotyping using next-generation DNA sequencing data.  
862 *Nature* ...

863 Dermauw, W, Van Leeuwen, T. 2014. The ABC gene family in arthropods: Comparative  
864 genomics and role in insecticide transport and resistance. *Insect Biochemistry*  
865 *and Molecular Biology* **45**: 89–110. 10.1016/j.ibmb.2013.11.001

866 Eddy, EM. 1976. Germ plasm and the differentiation of the germ cell line. *International*  
867 *review of cytology*

868 Edgar, RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high  
869 throughput. *Nucl Acids Res* **32**: 1792–1797. 10.1093/nar/gkh340

870 Elbashir, SM, Harborth, J, Lendeckel, W, Yalcin, A, Weber, K, Tuschl, T. 2001.  
871 Duplexes of 21-nucleotide RNAs mediate RNA interference in cultured  
872 mammalian cells. *Nature* **411**: 494–498.

873 Findley, SD, Tamanaha, M, Clegg, NJ, Ruohola-Baker, H. 2003. Maelstrom, a  
874 *Drosophila* spindle-class gene, encodes a protein that colocalizes with Vasa and  
875 RDE1/AGO1 homolog, Aubergine, in nuage. *Development* **130**: 859–871.  
876 10.1242/dev.00310

877 Finn, RD et al. 2016. The Pfam protein families database: towards a more sustainable  
878 future. *Nucleic Acids Res* **44**: D279–D285. 10.1093/nar/gkv1344

879 Flynt, A, Liu, N, Martin, R, Lai, EC. 2009. Dicing of viral replication intermediates during  
880 silencing of latent *Drosophila* viruses. *Proc Natl Acad Sci U S A* **106**: 5270–5275.  
881 10.1073/pnas.0813412106

882 Fraser, MJ, Smith, GE, Summers, MD. 1983. Acquisition of Host Cell DNA Sequences  
883 by Baculoviruses: Relationship Between Host DNA Insertions and FP Mutants of  
884 *Autographa californica* and *Galleria mellonella* Nuclear Polyhedrosis Viruses. *J*  
885 *Virol* **47**: 287–300.

886 Friedländer, MR et al. 2008. Discovering microRNAs from deep sequencing data using  
887 miRDeep. *Nat Biotech* **26**: 407–415. 10.1038/nbt1394

888 Friedlander, MR, Mackowiak, SD, Li, N, Chen, W, Rajewsky, N. 2012. miRDeep2  
889 accurately identifies known and hundreds of novel microRNA genes in seven  
890 animal clades. *Nucleic Acids Res* **40**: 37–52. 10.1093/nar/gkr688

891 Fujiwara, H, Osanai, M, Matsumoto, T, Kojima, KK. 2005. Telomere-specific non-LTR  
892 retrotransposons and telomere maintenance in the silkworm, *Bombyx mori*.  
893 *Chromosome Res* **13**: 455–467. 10.1007/s10577-005-0990-9

894 Fuková, I, Nguyen, P, Marec, F. 2005. Codling moth cytogenetics: karyotype,  
895 chromosomal location of rDNA, and molecular differentiation of sex  
896 chromosomes. *Genome* **48**: 1083–1092. 10.1139/g05-063

897 Futahashi, R et al. 2015. Extraordinary diversity of visual opsin genes in dragonflies.  
898 *PNAS* **112**: E1247–E1256. 10.1073/pnas.1424670112

899 Gao, G, Vandenberghe, LH, Wilson, JM. 2005. New recombinant serotypes of AAV  
900 vectors. *Curr Gene Ther* **5**: 285–297. 10.2174/1566523054065057

901 Gaunt, MW, Miles, MA. 2002. An insect molecular clock dates the origin of the insects  
902 and accords with palaeontological and biogeographic landmarks. *Molecular*  
903 *biology and evolution* **19**: 748–761. 10.1093/oxfordjournals.molbev.a004133

904 Ge, DT, Tipping, C, Brodsky, MH, Zamore, PD. 2016. Rapid Screening for CRISPR-  
905 Directed Editing of the *Drosophila* Genome Using white Coconversion. *G3*  
906 (*Bethesda*) **6**: 3197–3206. 10.1534/g3.116.032557

907 Ghildiyal, M et al. 2008. Endogenous siRNAs derived from transposons and mRNAs in  
908 *Drosophila* somatic cells. *Science* **320**: 1077–1081.

909 Girard, A, Sachidanandam, R, Hannon, GJ, Carmell, MA. 2006. A germline-specific  
910 class of small RNAs binds mammalian Piwi proteins. *Nature* **442**: 199–202.  
911 10.1038/nature04917

912 Glastad, KM, Hunt, BG, Yi, SV, Goodisman, MAD. 2011. DNA methylation in insects: on  
913 the brink of the epigenomic era. *Insect Molecular Biology* **20**: 553–565.  
914 10.1111/j.1365-2583.2011.01092.x

915 Gong, D-P, Zhang, H-J, Zhao, P, Xia, Q-Y, Xiang, Z-H. 2009. The Odorant Binding  
916 Protein Gene Family from the Genome of Silkworm, *Bombyx mori*. *BMC*  
917 *Genomics* **10**: 332. 10.1186/1471-2164-10-332

918 Goodman, WG, Granger, NA (2005). The Juvenile Hormones. Gilbert, LI, ed.  
919 (Amsterdam: Elsevier), pp. 319-408.

920 Goriaux, C, Desset, S, Renaud, Y, Vaury, C, Brasset, E. 2014. Transcriptional  
921 properties and splicing of the flamenco piRNA cluster. *EMBO Rep*  
922 10.1002/embr.201337898

923 Granados, RR, Derksen, ACG, Dwyer, KG. 1986. Replication of the *Trichoplusia ni*  
924 granulosis and nuclear polyhedrosis viruses in cell cultures. *Virology* **152**: 472–  
925 476. 10.1016/0042-6822(86)90150-9

926 Granados, RR, Guoxun, L, Derksen, ACG, McKenna, KA. 1994. A new insect cell line  
927 from *Trichoplusia ni* (BTI-Tn-5B1-4) susceptible to *Trichoplusia ni* single  
928 enveloped nuclear polyhedrosis virus. *Journal of Invertebrate Pathology* **64**: 260–  
929 266. 10.1016/S0022-2011(94)90400-6

930 Graveley, BR et al. 2011. The developmental transcriptome of *Drosophila*  
931 *melanogaster*. *Nature* **471**: 473–479. 10.1038/nature09715

932 Gu, L, Walters, JR, Knipple, DC. 2017. Conserved Patterns of Sex Chromosome  
933 Dosage Compensation in the Lepidoptera (WZ/ZZ): Insights from a Moth Neo-Z  
934 Chromosome. *Genome Biol Evol* **9**: 802–816. 10.1093/gbe/evx039

935 Gunawardane, LS et al. 2007. A Slicer-Mediated Mechanism for Repeat-Associated  
936 siRNA 5' End Formation in *Drosophila*. *Science* **315**: 1587–1590.  
937 10.1126/science.1140494

938 Hahn, C, Bachmann, L, Chevreux, B. 2013. Reconstructing mitochondrial genomes  
939 directly from genomic next-generation sequencing reads—a baiting and iterative  
940 mapping approach. *Nucl Acids Res* gkt371. 10.1093/nar/gkt371

941 Han, BW, Wang, W, Zamore, PD, Weng, Z. 2014. piPipes: a set of pipelines for piRNA

942 and transposon analysis via small RNA-seq, RNA-seq, degradome- and CAGE-  
943 seq, ChIP-seq and genomic DNA sequencing. *Bioinformatics* btu647.  
944 10.1093/bioinformatics/btu647

945 Han, BW, Wang, W, Li, C, Weng, Z, Zamore, PD. 2015. Noncoding RNA. piRNA-guided  
946 transposon cleavage initiates Zucchini-dependent, phased piRNA production.  
947 *Science* **348**: 817–821. 10.1126/science.aaa1264

948 He, L, Hannon, GJ. 2004. MicroRNAs: small RNAs with a big role in gene regulation.  
949 *Nat Rev Genet* **5**: 522–531.

950 Hekmat-Safe, DS, Scape, CR, McKinney, AJ, Tanouye, MA. 2002. Genome-Wide  
951 Analysis of the Odorant-Binding Protein Gene Family in *Drosophila*  
952 *melanogaster*. *Genome Res* **12**: 1357–1369. 10.1101/gr.239402

953 Hink, WF. 1972. A Catalog of Invertebrate Cell Lines. In *Invertebrate Tissue Culture*,  
954 Vol. 2, C., V, ed. (New York, NY: Academic Press), pp. 363-387.

955 Hink, WF. 1970. Established insect cell line from the cabbage looper, *Trichoplusia ni*.  
956 *Nature* **226**: 466–467. 10.1038/226466b0

957 Hirose, Y, Manley, JL. 1997. Creatine phosphate, not ATP, is required for 3' end  
958 cleavage of mammalian pre-mRNA in vitro. *J Biol Chem* **272**: 29636–29642.  
959 10.1074/jbc.272.47.29636

960 Honda, S et al. 2013. Mitochondrial protein BmPAPI modulates the length of mature  
961 piRNAs. *RNA* **19**: 1405–1418. 10.1261/rna.040428.113

962 Horwich, MD et al. 2007. The *Drosophila* RNA methyltransferase, DmHen1, modifies  
963 germline piRNAs and single-stranded siRNAs in RISC. *Curr Biol* **17**: 1265–1272.  
964 10.1016/j.cub.2007.06.030

965 Houwing, S et al. 2007. A role for Piwi and piRNAs in germ cell maintenance and  
966 transposon silencing in Zebrafish. *Cell* **129**: 69–82.

967 Hsu, PD et al. 2013. DNA targeting specificity of RNA-guided Cas9 nucleases. *Nat*  
968 *Biotechnol* **31**: 827–832. 10.1038/nbt.2647

969 Initiative, IGG. 2014. Genome Sequence of the Tsetse Fly (*Glossina morsitans*): Vector  
970 of African Trypanosomiasis. *Science* **344**: 380–386. 10.1126/science.1249656

971 Iwanaga, M, Adachi, Y, Uchiyama, K, Tsukui, K, Katsuma, S, Kawasaki, H. 2014. Long-  
972 term adaptation of the *Bombyx mori* BmN4 cell line to grow in serum-free culture.  
973 *In Vitro CellDevBiol—Animal* **50**: 792–796. 10.1007/s11626-014-9781-y

974 Janmaat, AF, Myers, J. 2003. Rapid evolution and the cost of resistance to *Bacillus*  
975 *thuringiensis* in greenhouse populations of cabbage loopers, *Trichoplusia ni*.  
976 *Proceedings of the Royal Society B* **270**: 2263–2270. 10.1098/rspb.2003.2497

977 Jinek, M, Chylinski, K, Fonfara, I, Hauer, M, Doudna, JA, Charpentier, E. 2012. A  
978 programmable dual-RNA-guided DNA endonuclease in adaptive bacterial  
979 immunity. *Science* **337**: 816–821. 10.1126/science.1225829

980 Jones, P et al. 2014. InterProScan 5: genome-scale protein function classification.  
981 *Bioinformatics* **30**: 1236–1240. 10.1093/bioinformatics/btu031

982 Kakumani, PK, Malhotra, P, Mukherjee, SK, Bhatnagar, RK. 2014. A draft genome  
983 assembly of the army worm, *Spodoptera frugiperda*. *Genomics* **104**: 134–143.  
984 10.1016/j.ygeno.2014.06.005

985 Kanost, MR et al. 2016. Multifaceted biological insights from a draft genome sequence  
986 of the tobacco hornworm moth, *Manduca sexta*. *Insect Biochem Mol Biol* **76**:  
987 118–147. 10.1016/j.ibmb.2016.07.005

988 Katsuma, S, Sugano, Y, Kiuchi, T, Shimada, T. 2015. Two Conserved Cysteine  
989 Residues Are Required for the Masculinizing Activity of the Silkworm Masc  
990 Protein. *J Biol Chem* **290**: 26114–26124. 10.1074/jbc.M115.685362

991 Kawamura, Y et al. 2008. *Drosophila* endogenous small RNAs bind to Argonaute2 in  
992 somatic cells. *Nature* **453**: 793–797. 10.1038/nature06938

993 Kawaoka, S et al. 2008. *Bombyx* small RNAs: genomic defense system against  
994 transposons in the silkworm, *Bombyx mori*. *Insect Biochem Mol Biol* **38**: 1058–  
995 1065. 10.1016/j.ibmb.2008.03.007

996 Kawaoka, S et al. 2009. The Bombyx ovary-derived cell line endogenously expresses  
997 PIWI/PIWI-interacting RNA complexes. *RNA* **15**: 1258–1264.  
998 10.1261/rna.1452209

999 Kawaoka, S et al. 2011. The silkworm W chromosome is a source of female-enriched  
1000 piRNAs. *RNA* **17**: 2144–2151. 10.1261/rna.027565.111

1001 Keeling, CI et al. 2013. Draft genome of the mountain pine beetle, *Dendroctonus*  
1002 *ponderosae* Hopkins, a major forest pest. *Genome Biology* **14**: R27. 10.1186/gb-  
1003 2013-14-3-r27

1004 Khvorova, A, Reynolds, A, Jayasena, SD. 2003. Functional siRNAs and miRNAs exhibit  
1005 strand bias. *Cell* **115**: 209–216.

1006 Kim, S, Kim, D, Cho, SW, Kim, J, Kim, JS. 2014. Highly efficient RNA-guided genome  
1007 editing in human cells via delivery of purified Cas9 ribonucleoproteins. *Genome*  
1008 *Res* **24**: 1012–1019. 10.1101/gr.171322.113

1009 Kinoshita, T, Hanazawa, T, Nohata, N, Okamoto, Y, Seki, N. 2012. The functional  
1010 significance of microRNA-375 in human squamous cell carcinoma: Aberrant  
1011 expression and effects on cancer pathways. *Journal of Human Genetics* **57**:  
1012 556–563. 10.1038/jhg.2012.75

1013 Kirino, Y, Mourelatos, Z. 2007. The mouse homolog of HEN1 is a potential methylase  
1014 for Piwi-interacting RNAs. *RNA* **13**: 1397–1401. 10.1261/rna.659307

1015 Klattenhoff, C et al. 2009. The *Drosophila* HP1 homolog Rhino is required for  
1016 transposon silencing and piRNA production by dual-strand clusters. *Cell* **138**:  
1017 1137–1149. 10.1016/j.cell.2009.07.014

1018 Komor, AC, Badran, AH, Liu, DR. 2017. CRISPR-Based Technologies for the  
1019 Manipulation of Eukaryotic Genomes. *Cell* **168**: 20–36.  
1020 10.1016/j.cell.2016.10.044

1021 Koren, S, Walenz, BP, Berlin, K, Miller, JR, Bergman, NH, Phillippy, AM. 2017. Canu:  
1022 scalable and accurate long-read assembly via adaptive k-mer weighting and



1023 repeat separation. *bioRxiv* 071282. 10.1101/071282

1024 Korf, I. 2004. Gene finding in novel genomes. *BMC Bioinformatics* **5**: 59. 10.1186/1471-  
1025 2105-5-59

1026 Kotin, RM. 2011. Large-scale recombinant adeno-associated virus production. *Hum Mol*  
1027 *Genet* **20**: R2–6. 10.1093/hmg/ddr141

1028 Kugler, J-M, Chen, Y-W, Weng, R, Cohen, SM. 2013. miR-989 Is Required for Border  
1029 Cell Migration in the Drosophila Ovary. *PLOS ONE* **8**: e67075.  
1030 10.1371/journal.pone.0067075

1031 Kurtz, S et al. 2004. Versatile and open software for comparing large genomes.  
1032 *Genome Biol* **5**: R12. 10.1186/gb-2004-5-2-r12

1033 Labbé, R, Caveney, S, Donly, C. 2011. Genetic analysis of the xenobiotic resistance-  
1034 associated ABC gene subfamilies of the Lepidoptera. *Insect Molecular Biology*  
1035 **20**: 243–256. 10.1111/j.1365-2583.2010.01064.x

1036 Lagesen, K, Hallin, P, Rødland, EA, Stærfeldt, H-H, Rognes, T, Ussery, DW. 2007.  
1037 RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucl Acids*  
1038 *Res* **35**: 3100–3108. 10.1093/nar/gkm160

1039 Lagos-Quintana, M, Rauhut, R, Lendeckel, W, Tuschl, T. 2001. Identification of Novel  
1040 Genes Coding for Small Expressed RNAs. *Science* **294**: 853–858.

1041 Lau, NC et al. 2009. Abundant primary piRNAs, endo-siRNAs, and microRNAs in a  
1042 *Drosophila* ovary cell line. *Genome Res* **19**: 1776–1785. 10.1101/gr.094896.109

1043 Lau, NC et al. 2006. Characterization of the piRNA complex from rat testes. *Science*  
1044 **313**: 363–367.

1045 Le Thomas, A et al. 2014. Transgenerationally inherited piRNAs trigger piRNA  
1046 biogenesis by changing the chromatin of piRNA clusters and inducing precursor  
1047 processing. *Genes & Development* **28**: 1667–1680. 10.1101/gad.245514.114

1048 Le Thomas, A, Marinov, G, Aravin, A. 2014. A Transgenerational Process Defines  
1049 piRNA Biogenesis in *Drosophila virilis*. *Cell Reports* 10.1016/j.celrep.2014.08.013

1050 Lee, E et al. 2013. Web Apollo: a web-based genomic annotation editing platform.  
1051 *Genome Biology* **14**: R93. 10.1186/gb-2013-14-8-r93

1052 Lewis, SH, Salmela, H, Obbard, DJ. 2016. Duplication and Diversification of Dipteran  
1053 Argonaute Genes, and the Evolutionary Divergence of Piwi and Aubergine.  
1054 *Genome Biol Evol* **8**: 507–518. 10.1093/gbe/evw018

1055 Li, C et al. 2009. Collapse of Germline piRNAs in the Absence of Argonaute3 Reveals  
1056 Somatic piRNAs in Flies. *Cell* **137**: 509–521. 10.1016/j.cell.2009.04.027

1057 Lim, AK, Kai, T. 2007. Unique germ-line organelle, nuage, functions to repress selfish  
1058 genetic elements in *Drosophila melanogaster*. *Proc Natl Acad Sci U S A* **104**:  
1059 6714–6719. 10.1073/pnas.0701920104

1060 Lin, S, Staahl, BT, Alla, RK, Doudna, JA. 2014. Enhanced homology-directed human  
1061 genome engineering by controlled timing of CRISPR/Cas9 delivery. *eLife* **3**:  
1062 e04766. 10.7554/eLife.04766

1063 Liu, L, Qi, H, Wang, J, Lin, H. 2011. PAPI, a novel TUDOR-domain protein, complexes  
1064 with AGO3, ME31B and TRAL in the nuage to silence transposition.  
1065 *Development* **138**: 1863–1873. 10.1242/dev.059287

1066 Liu, S et al. 2011. Genome-wide identification and characterization of ATP-binding  
1067 cassette transporters in the silkworm, *Bombyx mori*. *BMC Genomics* **12**: 491.  
1068 10.1186/1471-2164-12-491

1069 Lobo, N, Li, X, Fraser, MJ. 1999. Transposition of the piggyBac element in embryos of  
1070 *Drosophila melanogaster*, *Aedes aegypti* and *Trichoplusia ni*. *Mol Gen Genet*  
1071 **261**: 803–810. 10.1007/s004380050024

1072 Lucas, K, Raikhel, AS. 2013. Insect MicroRNAs: Biogenesis, expression profiling and  
1073 biological functions. *Insect Biochemistry and Molecular Biology* **43**: 24–38.  
1074 10.1016/j.ibmb.2012.10.009

1075 Lund, AH. 2009. miR-10 in development and cancer. *Cell Death Differ* **17**: 209–214.  
1076 10.1038/cdd.2009.58

1077 Luo, R et al. 2012. SOAPdenovo2: an empirically improved memory-efficient short-read  
1078 de novo assembler. *Gigascience* **1**: 18. 10.1186/2047-217X-1-18

1079 Mallet, J. 2007. Taxonomy of Lepidoptera: the scale of the problem. *The Lepidoptera*  
1080 *Taxome Project*

1081 Malone, CD et al. 2009. Specialized piRNA Pathways Act in Germline and Somatic  
1082 Tissues of the *Drosophila* Ovary. *Cell* **137**: 522–535. 10.1016/j.cell.2009.03.040

1083 Marchler-Bauer, A et al. 2015. CDD: NCBI's conserved domain database. *Nucl Acids*  
1084 *Res* **43**: D222–D226. 10.1093/nar/gku1221

1085 Marygold, SJ et al. 2007. The ribosomal protein genes and Minute loci of *Drosophila*  
1086 *melanogaster*. *Genome Biology* **8**: R216. 10.1186/gb-2007-8-10-r216

1087 Matijasevic, Z, Krzywicka-Racka, A, Sluder, G, Jones, SN. 2008. MdmX regulates  
1088 transformation and chromosomal stability in p53-deficient cells. *Cell Cycle* **7**:  
1089 2967–2973. 10.4161/cc.7.19.6797

1090 McEwen, FL, Hervey, GER. 1956. An Evaluation of Newer Insecticides for Control of  
1091 DDT-Resistant Cabbage Loopers. *Journal of Economic Entomology* **49**: 385–  
1092 387. 10.1093/jee/49.3.385

1093 McKenna, A, Hanna, M, Banks, E, Sivachenko..., A. 2010. The Genome Analysis  
1094 Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
1095 data. *Genome ...*

1096 Mitchell, A et al. 2015. The InterPro protein families database: the classification  
1097 resource after 15 years. *Nucl Acids Res* **43**: D213–D221. 10.1093/nar/gku1243

1098 Mohn, F, Handler, D, Brennecke, J. 2015. piRNA-guided slicing specifies transcripts for  
1099 Zucchini-dependent, phased piRNA biogenesis. *Science* **348**: 812–817.  
1100 10.1126/science.aaa1039

1101 Mohn, F, Sienski, G, Handler, D, Brennecke, J. 2014. The Rhino-Deadlock-Cutoff  
1102 Complex Licenses Noncanonical Transcription of Dual-Strand piRNA Clusters in  
1103 *Drosophila*. *Cell* **157**: 1364–1379. 10.1016/j.cell.2014.04.031

1104 Mon, H, Kusakabe, T, Lee, JM, Kawaguchi, Y, Koga, K. 2004. In vivo DNA double-  
1105 strand breaks enhance gene targeting in cultured silkworm cells. *Comparative*  
1106 *Biochemistry and Physiology Part B: Biochemistry and Molecular Biology* **139**:  
1107 99–106. 10.1016/j.cbpc.2004.06.013

1108 Nelson, DR. 2009. The Cytochrome P450 Homepage. *Human Genomics* **4**: 59.  
1109 10.1186/1479-7364-4-1-59

1110 Ni, JD, Baik, LS, Holmes, TC, Montell, C. 2017. A rhodopsin in the brain functions in  
1111 circadian photoentrainment in *Drosophila*. *Nature* **545**: 340–344.  
1112 10.1038/nature22325

1113 Niki, Y, Yamaguchi, T, Mahowald, AP. 2006. Establishment of stable cell lines of  
1114 *Drosophila* germ-line stem cells. *Proc Natl Acad Sci U S A* **103**: 16325–16330.  
1115 10.1073/pnas.0607435103

1116 Okamura, K, Balla, S, Martin, R, Liu, N, Lai, EC. 2008a. Two distinct mechanisms  
1117 generate endogenous siRNAs from bidirectional transcription in *Drosophila*  
1118 *melanogaster*. *Nat Struct Mol Biol* **15**: 581–590. 10.1038/nsmb.1438

1119 Okamura, K, Chung, WJ, Ruby, JG, Guo, H, Bartel, DP, Lai, EC. 2008b. The *Drosophila*  
1120 hairpin RNA pathway generates endogenous short interfering RNAs. *Nature* **453**:  
1121 803–806. 10.1038/nature07015

1122 Pace, JK, Feschotte, C. 2007. The evolutionary history of human DNA transposons:  
1123 Evidence for intense activity in the primate lineage. *Genome Res* **17**: 422–432.  
1124 10.1101/gr.5826307

1125 Pane, A, Jiang, P, Zhao, DY, Singh, M, Schupbach, T. 2011. The Cutoff protein  
1126 regulates piRNA cluster expression and piRNA production in the *Drosophila*  
1127 germline. *EMBO J* 10.1038/emboj.2011.334

1128 Patil, AA et al. 2017. Characterization of Armitage and Yb containing granules and their  
1129 relationship to nuage in ovary-derived cultured silkworm cell. *Biochem Biophys*  
1130 *Res Commun* **490**: 134–140. 10.1016/j.bbrc.2017.06.008

- 1131 Pelisson, A, Sarot, E, Payen-Groschene, G, Bucheton, A. 2007. A novel repeat-  
1132 associated small interfering RNA-mediated silencing pathway downregulates  
1133 complementary sense *gypsy* transcripts in somatic cells of the *Drosophila* ovary.  
1134 *J Virol* **81**: 1951–1960. 10.1128/JVI.01980-06
- 1135 Porcelli, D, Barsanti, P, Pesole, G, Caggese, C. 2007. The nuclear OXPHOS genes in  
1136 insecta: a common evolutionary origin, a common cis-regulatory motif, a common  
1137 destiny for gene duplicates. *BMC Evolutionary Biology* **7**: 215. 10.1186/1471-  
1138 2148-7-215
- 1139 Prud'homme, N, Gans, M, Masson, M, Terzian, C, Bucheton, A. 1995. *Flamenco*, a  
1140 gene controlling the gypsy retrovirus of *Drosophila melanogaster*. *Genetics* **139**:  
1141 697–711.
- 1142 Rainford, JL, Hofreiter, M, Nicholson, DB, Mayhew, PJ. 2014. Phylogenetic Distribution  
1143 of Extant Richness Suggests Metamorphosis Is a Key Innovation Driving  
1144 Diversification in Insects. *PLOS ONE* **9**: e109085. 10.1371/journal.pone.0109085
- 1145 Ran, FA, Hsu, PD, Wright, J, Agarwala, V, Scott, DA, Zhang, F. 2013. Genome  
1146 engineering using the CRISPR-Cas9 system. *Nat Protocols* **8**: 2281–2308.  
1147 10.1038/nprot.2013.143
- 1148 Robertson, HM, Gordon, KHJ. 2006. Canonical TTAGG-repeat telomeres and  
1149 telomerase in the honey bee, *Apis mellifera*. *Genome Res* **16**: 1345–1351.  
1150 10.1101/gr.5085606
- 1151 Rota-Stabelli, O, Daley, AC, Pisani, D. 2013. Molecular timetrees reveal a Cambrian  
1152 colonization of land and a new scenario for ecdysozoan evolution. *Current*  
1153 *Biology* **23**: 392–398. 10.1016/j.cub.2013.01.026
- 1154 Sahara, K, Marec, F, Traut, W. 1999. TTAGG telomeric repeats in chromosomes of  
1155 some insects and other arthropods. *Chromosome Res* **7**: 449–460.  
1156 10.1023/A:1009297729547
- 1157 Saito, K et al. 2009. A regulatory circuit for piwi by the large Maf gene *traffic jam* in

1158 *Drosophila*. *Nature* **461**: 1296–1299. 10.1038/nature08501

1159 Saito, K et al. 2006. Specific association of Piwi with rasiRNAs derived from  
1160 retrotransposon and heterochromatic regions in the *Drosophila* genome. *Genes*  
1161 *Dev* **20**: 2214–2222.

1162 Saito, K, Sakaguchi, Y, Suzuki, T, Suzuki, T, Siomi, H, Siomi, MC. 2007. Pimet, the  
1163 *Drosophila* homolog of HEN1, mediates 2'-O-methylation of Piwi-interacting  
1164 RNAs at their 3' ends. *Genes Dev* **21**: 1603–1608. 10.1101/gad.1563607

1165 Sanchez-Vargas, I et al. 2009. Dengue virus type 2 infections of *Aedes aegypti* are  
1166 modulated by the mosquito's RNA interference pathway. *PLoS Pathog* **5**:  
1167 e1000299. 10.1371/journal.ppat.1000299

1168 Schneider, I. 1979. Giemsa staining of insect chromosomes for karyotype analysis.  
1169 *Methods in Cell Science* **5**: 1027–1028.

1170 Schwarz, DS, Hutvagner, G, Du, T, Xu, Z, Aronin, N, Zamore, PD. 2003. Asymmetry in  
1171 the assembly of the RNAi enzyme complex. *Cell* **115**: 199–208.

1172 Shirayama, M et al. 2012. piRNAs Initiate an Epigenetic Memory of Nonself RNA in the  
1173 *C. elegans* Germline. *Cell* **150**: 65–77. 10.1016/j.cell.2012.06.015

1174 Sienski, G, Donertas, D, Brennecke, J. 2012. Transcriptional silencing of transposons  
1175 by Piwi and maelstrom and its impact on chromatin state and gene expression.  
1176 *Cell* **151**: 964–980. 10.1016/j.cell.2012.10.040

1177 Simão, FA, Waterhouse, RM, Ioannidis, P, Kriventseva, EV, Zdobnov, EM. 2015.  
1178 BUSCO: assessing genome assembly and annotation completeness with single-  
1179 copy orthologs. *Bioinformatics* **btv351**. 10.1093/bioinformatics/btv351

1180 Siomi, H, Siomi, MC. 2009. On the road to reading the RNA-interference code. *Nature*  
1181 **457**: 396–404. 10.1038/nature07754

1182 Slater, GSC, Birney, E. 2005. Automated generation of heuristics for biological  
1183 sequence comparison. *BMC Bioinformatics* **6**: 31. 10.1186/1471-2105-6-31

1184 Smit, AFA, Hubley, R, Green, P. 2017. RepeatMasker Open-3.0.

1185 Sokol, NS, Ambros, V. 2005. Mesodermally expressed *Drosophila* microRNA-1 is  
1186 regulated by Twist and is required in muscles during larval growth. *Genes Dev*  
1187 **19**: 2343–2354.

1188 Stanke, M, Tzvetkova, A, Morgenstern, B. 2006. AUGUSTUS at EGASP: using EST,  
1189 protein and genomic alignments for improved gene prediction in the human  
1190 genome. *Genome Biology* **7**: 1–8. 10.1186/gb-2006-7-s1-s11

1191 Suetsugu, Y et al. 2013. Large Scale Full-Length cDNA Sequencing Reveals a Unique  
1192 Genomic Landscape in a Lepidopteran Model Insect, *Bombyx mori*. *G3* **3**: 1481–  
1193 1492. 10.1534/g3.113.006239

1194 Takasu, Y, Tamura, T, Sajwan, S, Kobayashi, I, Zurovec, M. 2014. The use of TALENs  
1195 for nonhomologous end joining mutagenesis in silkworm and fruitfly. *Methods* **69**:  
1196 46–57. 10.1016/j.ymeth.2014.02.014

1197 Tam, OH et al. 2008. Pseudogene-derived small interfering RNAs regulate gene  
1198 expression in mouse oocytes. *Nature* **453**: 534–538. 10.1038/nature06904

1199 Teleman, AA, Maitra, S, Cohen, SM. 2006. *Drosophila* lacking microRNA miR-278 are  
1200 defective in energy homeostasis. *Genes Dev* **20**: 417–422. 10.1101/gad.374406

1201 Terakita, A. 2005. The opsins. *Genome Biology* **6**: 213. 10.1186/gb-2005-6-3-213

1202 The, ISG. 2008. The genome of a lepidopteran model insect, the silkworm *Bombyx*  
1203 *mori*. *Insect Biochemistry and Molecular Biology* **38**: 1036–1045.  
1204 10.1016/j.ibmb.2008.11.004

1205 Traut, W, Sahara, K, Marec, F. 2007. Sex chromosomes and sex determination in  
1206 Lepidoptera. *Sex Dev* **1**: 332–346. 10.1159/000111765

1207 Tyler, DM et al. 2008. Functionally distinct regulatory RNAs generated by bidirectional  
1208 transcription and processing of microRNA loci. *Genes Dev* **22**: 26–36.  
1209 10.1101/gad.1615208

1210 Vagin, VV, Sigova, A, Li, C, Seitz, H, Gvozdev, V, Zamore, PD. 2006. A distinct small  
1211 RNA pathway silences selfish genetic elements in the germline. *Science* **313**:

1212 320–324. 10.1126/science.1129333

1213 Van der Auwera, GA et al. 2013. From FastQ data to high confidence variant calls: the  
1214 Genome Analysis Toolkit best practices pipeline. *Curr Protoc Bioinformatics* **43**:  
1215 11.10.1–33. 10.1002/0471250953.bi1110s43

1216 van Oers, MM, Pijlman, GP, Vlak, JM. 2015. Thirty years of baculovirus-insect cell  
1217 protein expression: from dark horse to mainstream technology. *J Gen Virol* **96**:  
1218 6–23. 10.1099/vir.0.067108-0

1219 van Rij, RP et al. 2006. The RNA silencing endonuclease Argonaute 2 mediates specific  
1220 antiviral immunity in *Drosophila melanogaster*. *Genes Dev* **20**: 2985–2995.  
1221 10.1101/gad.1482006

1222 van Schooten, B, Jiggins, CD, Briscoe, AD, Papa, R. 2016. Genome-wide analysis of  
1223 ionotropic receptors provides insight into their evolution in *Heliconius* butterflies.  
1224 *BMC Genomics* **17**: 254. 10.1186/s12864-016-2572-y

1225 Vastenhouw, NL et al. 2010. Chromatin signature of embryonic pluripotency is  
1226 established during genome activation. *Nature* 10.1038/nature08866

1227 Velarde, RA, Sauer, CD, O. Walden, KK, Fahrbach, SE, Robertson, HM. 2005.  
1228 Pteropsin: A vertebrate-like non-visual opsin expressed in the honey bee brain.  
1229 *Insect Biochemistry and Molecular Biology* **35**: 1367–1377.  
1230 10.1016/j.ibmb.2005.09.001

1231 Vrettos, N, Maragkakis, M, Alexiou, P, Mourelatos, Z. 2017. Kc167, a widely used  
1232 *Drosophila* cell line, contains an active primary piRNA pathway. *RNA* **23**: 108–  
1233 118. 10.1261/rna.059139.116

1234 Walker, BJ et al. 2014. Pilon: An Integrated Tool for Comprehensive Microbial Variant  
1235 Detection and Genome Assembly Improvement. *PLOS ONE* **9**: e112963.  
1236 10.1371/journal.pone.0112963

1237 Walters, JR, Hardcastle, TJ. 2011. Getting a Full Dose? Reconsidering Sex  
1238 Chromosome Dosage Compensation in the Silkworm, *Bombyx mori*. *Genome*



1239 *Biol Evol* **3**: 491–504. 10.1093/gbe/evr036

1240 Walters, JR, Hardcastle, TJ, Jiggins, CD. 2015. Sex Chromosome Dosage  
1241 Compensation in Heliconius Butterflies: Global yet Still Incomplete. *Genome Biol*  
1242 *Evol* **7**: 2545–2559. 10.1093/gbe/evv156

1243 Wang, W et al. 2008. Chromosomal transposition of PiggyBac in mouse embryonic  
1244 stem cells. *PNAS* **105**: 9290–9295. 10.1073/pnas.0801017105

1245 Wang, Y et al. 2013. The CRISPR/Cas system mediates efficient genome engineering  
1246 in *Bombyx mori*. *Cell research* **23**: 1414. 10.1038/cr.2013.146

1247 Wanner, KW, Robertson, HM. 2008. The gustatory receptor family in the silkworm moth  
1248 *Bombyx mori* is characterized by a large expansion of a single lineage of putative  
1249 bitter receptors. *Insect Molecular Biology* **17**: 621–629. 10.1111/j.1365-  
1250 2583.2008.00836.x

1251 Watanabe, T et al. 2008. Endogenous siRNAs from naturally formed dsRNAs regulate  
1252 transcripts in mouse oocytes. *Nature* **453**: 539–543. 10.1038/nature06908

1253 Webster, A et al. 2015. Aub and Ago3 Are Recruited to Nuage through Two  
1254 Mechanisms to Form a Ping-Pong Complex Assembled by Krimper. *Mol Cell* **59**:  
1255 564–575. 10.1016/j.molcel.2015.07.017

1256 Wei, Y, Chen, S, Yang, P, Ma, Z, Kang, L. 2009. Characterization and comparative  
1257 profiling of the small RNA transcriptomes in two phases of locust. *Genome*  
1258 *Biology* **10**: R6. 10.1186/gb-2009-10-1-r6

1259 Wheat, CW, Wahlberg, N. 2013. Phylogenomic insights into the cambrian explosion, the  
1260 colonization of land and the evolution of flight in arthropoda. *Syst Biol* **62**: 93–  
1261 109. 10.1093/sysbio/sys074

1262 Wickham, TJ, Davis, T, Granados, RR, Shuler, ML, Wood, HA. 1992. Screening of  
1263 Insect Cell Lines for the Production of Recombinant Proteins and Infectious Virus  
1264 in the Baculovirus Expression System. *Biotechnol Progress* **8**: 391–396.  
1265 10.1021/bp00017a003

1266 Xavier, B, David, M, Piulachs, M-D. 2005. The Mevalonate Pathway and the Synthesis  
1267 of Juvenile Hormone in Insects. *Annual Review of Entomology* **50**: 181–199.  
1268 10.1146/annurev.ento.50.071803.130356

1269 Xiang, H et al. 2010. Single base-resolution methylome of the silkworm reveals a sparse  
1270 epigenomic map. *Nat Biotech* **28**: 516–520. 10.1038/nbt.1626

1271 Xiol, J et al. 2012. A Role for Fkbp6 and the Chaperone Machinery in piRNA  
1272 Amplification and Transposon Silencing. *Mol Cell* **47**: 970–979.  
1273 10.1016/j.molcel.2012.07.019

1274 Yoshihama, M et al. 2002. The Human Ribosomal Protein Genes: Sequencing and  
1275 Comparative Analysis of 73 Genes. *Genome Res* **12**: 379–390.  
1276 10.1101/gr.214202

1277 You, M et al. 2013. A heterozygous moth genome provides insights into herbivory and  
1278 detoxification. *Nat Genet* **45**: 220–225. 10.1038/ng.2524

1279 Yu, Q-Y, Lu, C, Li, W-L, Xiang, Z-H, Zhang, Z. 2009. Annotation and expression of  
1280 carboxylesterases in the silkworm, *Bombyx mori*. *BMC Genomics* **10**: 553.  
1281 10.1186/1471-2164-10-553

1282 Yu, Q et al. 2008. Identification, genomic organization and expression pattern of  
1283 glutathione S-transferase in the silkworm, *Bombyx mori*. *Insect Biochemistry and*  
1284 *Molecular Biology* **38**: 1158–1164. 10.1016/j.ibmb.2008.08.002

1285 Yu, Z et al. 2014. Various applications of TALEN-and CRISPR/Cas9-mediated  
1286 homologous recombination to modify the *Drosophila* genome. *Biology Open* **3**:  
1287 271–280. 10.1242/bio.20147682

1288 Zambon, RA, Vakharia, VN, Wu, LP. 2006. RNAi is an antiviral immune response  
1289 against a dsRNA virus in *Drosophila melanogaster*. *Cell Microbiol* **8**: 880–889.

1290 Zhan, S, Merlin, C, Boore, J, Reppert, S. 2011. The Monarch Butterfly Genome Yields  
1291 Insights into Long-Distance Migration. *Cell* **147**: 1171–1185.  
1292 10.1016/j.cell.2011.09.052

1293 Zhang, Z, Theurkauf, WE, Weng, Z, Zamore, PD. 2012. Strand-specific libraries for high  
1294 throughput RNA sequencing (RNA-Seq) prepared without poly(A) selection.  
1295 *Silence* **3**: 9. 10.1186/1758-907X-3-9

1296 Zhang, Z et al. 2014. The HP1 Homolog Rhino Anchors a Nuclear Complex that  
1297 Suppresses piRNA Precursor Splicing. *Cell* **157**: 1353–1363.  
1298 10.1016/j.cell.2014.04.030

1299 Zhang, Z, Aslam, AFM, Liu, X, Li, M, Huang, Y, Tan, A. 2015. Functional analysis of  
1300 Bombyx Wnt1 during embryogenesis using the CRISPR/Cas9 system. *Journal of*  
1301 *Insect Physiology* **79**: 73–79. 10.1016/j.jinsphys.2015.06.004

1302 Zhu, L, Mon, H, Xu, J, Lee, JM, Kusakabe, T. 2015. CRISPR/Cas9-mediated knockout  
1303 of factors in non-homologous end joining pathway enhances gene targeting in  
1304 silkworm cells. *Sci Rep* **5**: 18103. 10.1038/srep18103

1305 Zhuang, J, Wang, J, Theurkauf, W, Weng, Z. 2014. TEMP: a computational method for  
1306 analyzing transposable element polymorphism in populations. *Nucleic Acids Res*  
1307 **42**: 6826–6838. 10.1093/nar/gku323

1308 Zimyanin, VL et al. 2008. In vivo imaging of oskar mRNA transport reveals the  
1309 mechanism of posterior localization. *Cell* **134**: 843–853.  
1310 10.1016/j.cell.2008.06.053

1311

1312 **Table Legend**

1313 Table 1. Genome and gene set statistics for *T. ni* and *B. mori* (*The, ISG, 2008*).  
 1314 Cytochrome P450s, glutathione S-transferases, carboxylesterases, and ATP-binding  
 1315 cassette transporters for *B. mori* were retrieved from (*Yu, Q et al., 2008; Yu, Q-Y et al.,*  
 1316 *2009; Ai, J et al., 2011; Liu, S et al., 2011*).

	<i>T. ni</i>	<i>B. mori</i>
<b>Genome Metrics</b>		
Genome size (Mb)	368.2	431.7
Chromosome count	28	28
Scaffold N50 (Mb)	14.2	3.7
Contig N50 (kb)	621.9	15.5
Mitochondrial genome (kb)	15.8	15.7
<b>Quality Control Metrics</b>		
BUSCO complete (%)	97.5	95.5
CRP genes (%)	100%	100%
OXPPOS genes (%)	100%	100%
<b>Genomic Features</b>		
Repeat content (%)	20.5%	43.6%
GC content	35.6%	37.3%
CpG (O/E)	1.07	1.13
Coding (%)	5.58	4.11
Sex chromosomes	ZW	ZW
<b>Gene Statistics</b>		
Protein-coding genes	14,043	14,623
with Pfam matches	9,295	9,685
with GO terms	9,790	10,148
Cytochrome P450 proteins	108	83
Glutathione S-transferases	34	23
Carboxylesterases	87	76
ATP-binding cassette transporters	54	51
Universal orthologs lost	156	75
Species-specific genes	3,098	2,313

1318 **Figure Legends**

1319 **Figure 1.** Chromosomes and genes in the *T. ni* genome based on data from the Hi5 cell line. (A)  
1320 Genome assembly and annotation workflow. (B) An example of a DAPI-stained spread of Hi5  
1321 cell mitotic chromosomes used to determine the karyotype. (C) Phylogenetic tree and  
1322 orthology assignment of *T. ni* with 18 arthropod and two mammalian genomes. Colors denote  
1323 gene categories. The category 1:1:1 represents universal single-copy orthologs, allowing  
1324 absence and/or duplication in one genome. N:N:N orthologs include orthologs with variable  
1325 copy numbers across species, allowing absence in one genome or two genomes from different  
1326 orders. Lepidoptera-specific genes are present in at least three of the four lepidopteran  
1327 genomes; Hymenoptera-specific genes are present in at least one wasp or bee genome and at  
1328 least one ant genome. Coleoptera-specific genes are present in both coleopteran genomes;  
1329 Diptera-specific genes present in at least one fly genome and one mosquito genome. Insect  
1330 indicates other insect-specific genes. Mammal-specific genes are present in both mammalian  
1331 genomes. The phylogenetic tree is based on the alignment of 1:1:1 orthologs.

1332 **Figure 1—figure supplement 1.** Hi5 cell Karyotyping. Thirty images showing the  
1333 numbers of chromosomes (N) in Hi5 cells. N ranged from 103 to 122; mean  $\pm$  S.D. =  
1334  $111.7 \pm 5.45$ . Since lepidopteran cell lines are typically tetraploid, the haploid genome  
1335 likely contains 28 (mean  $\pm$  S.D. =  $27.9 \pm 1.36$ ) pairs of chromosomes.

1336 **Figure 1—figure supplement 2.** Phylogenetic tree of 21 species showing the scale, branch  
1337 lengths and bootstrap support. Strict 1:1:1 orthologs were used to compute the phylogenetic  
1338 tree using the maximum likelihood method. Black, branch length; red, bootstrap support.

1339 **Figure 1—figure supplement 3.** Opsins in insects.

1340 **Figure 2.** *T. ni* males are ZZ and females are ZW. (A) Normalized contig coverage in males and  
1341 females. (B) Relative repeat content, gene density, transcript abundance (female and male  
1342 thoraces), and piRNA density of autosomal, Z-linked, and W-linked contigs (ovary). (C)  
1343 Multiple sequence alignment of the conserved region of the sex-determining gene *masc*  
1344 among the lepidopteran species.

1345 **Figure 2—figure supplement 1.** (A) Genomic coverage comparison of Z-linked, W-  
1346 linked and autosomal contigs. Contig coverage was shuffled 1,000,000 times to  
1347 calculate the coverage ratio. Outliers are not shown. (B) Autosomal, Z-linked and W-  
1348 linked transcript abundance in Hi5 cells and *T. ni* tissues. (C) Transcript abundance  
1349 ratios of autosomal, Z-linked, and W-linked genes in Hi5 cells and *T. ni* tissues. Error  
1350 bars represent 95% confidence interval estimated from 1,000 bootstrap replicates. (D).  
1351 Sex-specific splicing of *T. ni doublesex* pre-mRNA.

1352 **Figure 2—figure supplement 2.** CpG ratios and transposons. (A) Distribution of  
1353 observed-to-expected CpG ratios in protein-coding genes (left panel) and in 500 bp  
1354 genomic windows (right panel) in *A. mellifera*, *B. mori*, *D. plexippus*, *D. melanogaster*,  
1355 *P. xylostella*, *T. castaneum*, and *T. ni*. (B) Proportion of the genome occupied by  
1356 transposons versus transposon sequence divergence. Sequence divergence was  
1357 calculated by comparing individual transposon copies with the corresponding  
1358 consensus sequence (See Materials and methods). (C) Repeat content in lepidopteran  
1359 genomes.

1360 **Figure 3.** miRNA expression in *T. ni*. (A) Comparison of miRNA abundance in male and female  
1361 *T. ni* thoraces. Solid circles, miRNAs with FDR < 0.1 and fold change > 2. Outlined circles, all  
1362 other miRNAs. (B) Comparison of the tissue distribution of the 44 most abundant miRNAs  
1363 among *T. ni* ovaries, testes, and Hi5. (C) Heat map showing the abundance of miRNAs in (B).  
1364 miRNAs are ordered according to abundance in ovary. Conservation status uses the same color  
1365 scheme in (A).

1366 **Figure 4.** siRNA. (A) Distribution of siRNAs mapping to TNCL virus in the genomic  
1367 (blue) and anti-genomic orientation (red). Inset: length distribution of TNCL virus-  
1368 mapping small RNAs. (B) Distance between the 3' and 5' ends of siRNAs on opposite  
1369 viral strands. (C) Distance between the 3' and 5' ends of siRNAs on the same viral  
1370 strand. (D) Length distribution of small RNAs from unoxidized and oxidized small RNA-  
1371 seq libraries. (E) Lepidopteran siRNAs are not 2'-O-methylated. The box plots display  
1372 the ratio of abundance (as a fraction of all small RNAs sequenced) for each siRNA in  
1373 oxidized versus unoxidized small RNA-seq libraries. The tree shows the phylogenetic  
1374 relationships of the analyzed insects. Outliers are not shown.

1375 **Figure 4—figure supplement 1.** (A) siRNA length distributions for multiple insects in  
1376 oxidized and unoxidized small RNA-seq libraries. (B) Length distribution of fully  
1377 matched and tailed TNCL virus-siRNAs.

1378 **Figure 4—figure supplement 2.** Loading asymmetry of siRNAs mapping to TNCL  
1379 RNA1 (A) and RNA2 (B). For each single-stranded siRNA species, we searched for  
1380 siRNAs on the other strand that when paired produce a typical siRNA duplex with two-  
1381 nucleotide overhanging 3' ends.

1382 **Figure 5.** piRNAs and miRNAs in the *T. ni* genome. (A) Abundance of mRNAs encoding piRNA  
1383 pathway proteins) in Hi5 cells, ovary, testis, and thorax. (B) Ideogram displaying the positions  
1384 of miRNA genes (arrowheads) and piRNA clusters in the *T. ni* genome. Color-coding reports  
1385 tissue expression for Hi5 cells, ovaries, testis, and thorax. Contigs that cannot be placed onto  
1386 chromosome-length scaffolds are arbitrarily concatenated and are marked 'Un.' (C)  
1387 Distribution of piRNAs among the autosomes, Z, and W chromosomes in Hi5 cells, ovary,  
1388 testis, and female and male thorax, compared with the fraction of the genome corresponding  
1389 to autosomes, W, and Z chromosomes.

1390 **Figure 5—figure supplement 1.** piRNA abundance (ppm) along the most productive  
1391 piRNA cluster. Top, fixed scale (some data clipped); bottom, auto-scaled.

1392 **Figure 5—figure supplement 2.** (A) piRNA clusters tend to produce piRNAs that are  
1393 antisense to transposons. The *x*-axis represents the ratio of piRNAs from the plus  
1394 strand to piRNAs from the minus strand, with the dotted lines indicating twofold  
1395 difference. The *y*-axis indicates the ratio of transposons lengths on the plus strand over  
1396 transposon length on the minus strand. The solid line indicates regression line and  
1397 shading indicates 95% confidence interval by LOWESS. Boxplot shows fractions of  
1398 antisense transposons (i.e. transposons inserted opposite to the direction of piRNAs  
1399 precursor transcription) in dual- and uni-strand clusters. Outliers are not shown.  
1400 Wilcoxon rank-sum test. (B) piRNA densities on autosomal, Z-linked and W-linked  
1401 contigs in Hi5 cells, ovary, testis, and female and male thorax. (C) Abundance of  
1402 piRNAs from putative W-linked genes.

1403 **Figure 6.** (A) Hi5-specific piRNA clusters contain younger transposon copies. RC,  
1404 rolling-circle transposons; LINE, Long interspersed nuclear elements; LTR, long  
1405 terminal repeat retrotransposon; DNA, DNA transposon. (B) Comparison of piRNA  
1406 abundance per cluster in female and male thorax. (C) piRNA precursors are rarely



1407 spliced. The number of introns supported by exon-exon junction-mapping reads is  
1408 shown for protein-coding genes and for piRNA clusters for each tissue or cell type. (D)  
1409 piRNA precursors are inefficiently spliced. Splicing efficiency is defined as the ratio of  
1410 spliced over unspliced reads. Splice sites were categorized into those inside and  
1411 outside piRNA clusters. Outliers are not shown.

1412 **Figure 6—figure supplement 1.** (A) Comparison of piRNA abundance (ppm) from  
1413 ovary and Hi5 piRNA-producing loci and from ovary and testis piRNA-producing loci. (B)  
1414 piRNA cluster lengths in *T. ni* ovary, testis, thorax, and Hi5 cells. (C) Motifs around  
1415 intron boundaries of predicted protein-coding gene models within and outside of piRNA  
1416 clusters.

1417 **Figure 7.** Genome editing in Hi5 cells. (A) Strategy for using Cas9/sgRNA RNPs to generate a  
1418 loss-of-function *TnPiwi* deletion allele. Red, protospacer-adjacent motif (PAM); blue,  
1419 protospacer sequence. Arrows indicate the diagnostic forward and reverse primers used in PCR  
1420 to detect genomic deletions ( $\Delta$ ). Sanger sequencing of the ~1700 bp PCR products validated  
1421 the *TnPiwi* deletions. (B) An example of PCR analysis of a *TnPiwi* deletion event. (C) Strategy  
1422 for using Cas9/sgRNA RNPs and a single-stranded DNA homology donor to insert EGFP and an  
1423 HA-tag in-frame with the *vasa* open reading frame. (D) An example of PCR analysis of a  
1424 successful HDR event. DNA isolated from wild type (WT) and FACS-sorted, EGFP-expressing  
1425 Hi5 cells (HDR) were used as templates.

1426 **Figure 8.** Hi5 cells contain nuage. (A) Schematic of single-clone selection of genome-  
1427 edited Hi5 cells using the strategy described in Figure 7C. (B) A representative field of  
1428 Hi5 cells edited to express EGFP-HA-Vasa from the endogenous locus. (C) A  
1429 representative image of a fixed, EGFP-HA-Vasa-expressing Hi5 cell stained with DAPI,

1430 anti-EGFP and anti-HA antibodies. EGFP and HA staining colocalize in a perinuclear  
1431 structure consistent with Vasa localizing to nuage.  
1432

1433 **Supplementary Files**

1434 **Supplementary file 1.** *T. ni* genome statistics. (A) BUSCO assessments of *T. ni* and six other  
1435 genomes. (B) CRP genes. (C) Genes in the OXPHOS pathway. (D) Genome comparisons.  
1436 Genomes assembled using paired-end DNA-seq data from male and female *T. ni* pupae are  
1437 compared with the Hi5 genome as the reference. The dot plots show genome alignments for  
1438 contigs  $\geq 1$  kb. (E) Numbers of genes in lepidopteran genomes. (F) Positions of telomeric  
1439 repeats: position of (TTAGG)<sub>n</sub> longer than 100 nt. (G) Transposons in *T. ni* subtelomeric  
1440 regions. (H) Repeat statistics for the *T. ni* genome. (I) Transposon family divergence rates. (J)  
1441 Manual curation of W-linked protein-coding genes and miRNAs.

1442 **Supplementary file 2.** Genes encoding small RNA pathway proteins. (A) Genes encoding  
1443 miRNA and siRNA pathway proteins. (B) Genes encoding piRNA pathway proteins (grouped by  
1444 sequence orthology).

1445 **Supplementary file 3.** *T. ni* miRNAs, siRNAs and piRNAs. (A) miRNA annotation. (B)  
1446 Mapping statistics for endogenous siRNAs in *T. ni* and *D. melanogaster*. (C) piRNA  
1447 cluster lengths. piRNA cluster coordinates in Hi5 (D), ovary (E), testis (F), female thorax  
1448 (G), and male thorax (H).

1449 **Supplementary file 4.** mirDeep2 output for *T. ni* miRNAs

1450 **Supplementary file 5.** Genomes used in this study.

1451 **Supplementary file 6.** *T. ni* detoxification-related genes. (A) P450 gene counts by  
1452 clade in *T. ni* and *B. mori*. (B) Sequences of P450 proteins. (C) Sequences of  
1453 glutathione-S-transferase proteins. (D) Carboxylesterase gene counts by clade in *T. ni*  
1454 and *B. mori*. (E) Sequences of carboxylesterase proteins. (F) ATP-binding cassette  
1455 transporter gene counts by clade in *T. ni* and *B. mori*. (G) Sequences of ATP-binding  
1456 cassette transporter proteins.

1457 **Supplementary file 7.** *T. ni* chemoreception genes. (A) Sequences of olfactory receptor  
1458 proteins. (B) Sequences of gustatory receptor proteins. (C) Sequences of ionotropic receptor  
1459 proteins.

1460 **Supplementary file 8.** Genes in the juvenile hormone biosynthesis and degradation  
1461 pathways.

1462 **Supplementary file 9.** Genome-modified sequences.

1463 **Supplementary file 10.** Single-stranded DNA donor purification  
1464

1465 **Materials and methods**

1466 **Genomic DNA libraries**

1467 Hi5 cells (ThermoFisher, Waltham, MA, USA) were cultured at 27°C in Express Five  
1468 Serum Free Medium (ThermoFisher) following the manufacturer's protocol. Thorax were  
1469 dissected from four-day-old female or male *T. ni* pupa (Benzon Research, Carlisle, PA,  
1470 USA). Cells or tissues were lysed in 2× PK buffer (200 mM Tris-HCl [pH7.5], 300 mM  
1471 NaCl, 25 mM EDTA, 2% w/v SDS) containing 200 µg/ml proteinase K at 65°C for 1 h,  
1472 extracted with phenol:chloroform:isoamyl alcohol (25:24:1; Sigma, St. Louis, MO, USA),  
1473 and genomic DNA collected by ethanol precipitation. The precipitate was dissolved in  
1474 10 mM Tris-HCl (pH 8.0), 0.1mM EDTA, treated with 20 µg/ml RNase A at 37°C for 30  
1475 min, extracted with phenol:chloroform:isoamyl alcohol (25:24:1), and collected by  
1476 ethanol precipitation. DNA concentration was determined (Qubit dsDNA HS Assay,  
1477 ThermoFisher). Genomic DNA libraries were prepared from 1 µg genomic DNA  
1478 (Illumina TruSeq LT kit, NextSeq 500, Illumina, San Diego, CA, USA).

1479 Long-read genome sequencing with a 23 kb average insert range was constructed from  
1480 16 µg genomic DNA using the SMRTbell Template Prep Kit 1.0 SPv3 (Pacific Biosciences,  
1481 Menlo Park, CA, USA) according to manufacturer's protocol. Sequence analysis was performed  
1482 using P6/C4 chemistry, 240 min data collection per SMRTcell on an RS II instrument (Pacific  
1483 Biosciences). Mate pair libraries with 2 kb and 8 kb insert sizes were constructed (Nextera Mate  
1484 Pair Library Prep Kit, Illumina) according to manufacturer's protocol from 1 µg Hi5 cell genomic  
1485 DNA. Libraries were sequenced to obtain 79 nt paired-end reads (NextSeq500, Illumina).

1486 **Hi-C**

1487 Hi-C libraries were generated from Hi5 cells as described (*Belton, J-M et al., 2012*),  
1488 except that 50 million cells were used. Hi-C Libraries were sequenced using the  
1489 NextSeq500 platform (Illumina) to obtain 79 nt, paired-end reads.

1490 **Karyotyping**

1491 Hi5 cells were first incubated in Express Five medium containing 1 µg/ml colcemid at  
1492 27°C for 8 h (*Schneider, I, 1979*), then in 4 ml 0.075 M KCl for 30 min at 37°C, and fixed  
1493 with freshly prepared methanol:acetic acid (3:1, v/v) precooled to –20°C. Mitotic  
1494 chromosomes were spread, mounted by incubation in ProLong Gold Antifade Mountant  
1495 with DAPI (4',6'-diamidino-2-phenylindole; ThermoFisher) overnight in the dark, and  
1496 imaged using a DMI8 fluorescence microscope equipped with an 63× 1.40 N.A. oil  
1497 immersion objective (HCX PL APO CS2, Leica Microsystems, Buffalo Grove, IL, USA)  
1498 as described (*Matijasevic, Z et al., 2008*).

1499 **Small RNA libraries**

1500 Ovaries, testes, and thoraces were dissected from cabbage looper adults 24–48 h after  
1501 emerging. Total RNA (30 µg) was isolated (mirVana miRNA isolation kit, Ambion,  
1502 Austin, TX, USA) and sequenced using the NextSeq500 platform (Illumina) to obtain 59  
1503 nt single-end reads as previously described (*Han, BW et al., 2015*).

1504 **RNA-seq**

1505 Adult ovaries, testes, or thoraces were dissected from cabbage looper adults 24 to 48 h after  
1506 emerging. Total RNA (3 µg) was purified (mirVana miRNA isolation kit, Ambion) and  
1507 sequenced as described (*Zhang, Z et al., 2012*) using the NextSeq500 platform (Illumina) to  
1508 obtain 79 nt, paired-end reads.

1509 **Genome assembly**

1510 Canu v1.3 (*Koren, S et al., 2017*) was used to assemble long reads into contigs,  
1511 followed by Quiver ([github.com/PacificBiosciences/GenomicConsensus](https://github.com/PacificBiosciences/GenomicConsensus)) to polish the  
1512 contigs using the same set of reads. Pilon (*Walker, BJ et al., 2014*) was used to further  
1513 polish the assembly using Illumina paired-end reads. Finally, to assemble the genome  
1514 into chromosome-length scaffolds, we joined the contigs using Hi-C reads and

1515 LACHESIS (Burton, JN et al., 2013). The mitochondrial genome was assembled  
1516 separately using MITObim (six iterations, *D. melanogaster* mitochondrial genome as  
1517 bait; [Hahn, C et al., 2013]).

1518 To evaluate the quality of the genome assembly, we ran BUSCO v3 (Simão, FA  
1519 et al., 2015) using the arthropod profile and default parameters to identify universal  
1520 single-copy orthologs. We further evaluated genome quality using conserved gene sets:  
1521 OXPHOS and CRP genes. *B. mori* and *D. melanogaster* OXPHOS and CRP protein  
1522 sequences were retrieved (Marygold, SJ et al., 2007; Porcelli, D et al., 2007) and  
1523 BLASTp was used to search for their *T. ni* homologs, which were further validated by  
1524 querying using InterPro (Jones, P et al., 2014; Mitchell, A et al., 2015). We also  
1525 assembled *T. ni* genomes from male and female animals respectively using  
1526 SOAPdenovo2 (kmer size 69; [Luo, R et al., 2012]). We then compared the animal  
1527 genomes with the *T. ni* genome assembled from Hi5 cells using QUAST (-m 500)  
1528 (Gurevich et al., 2013, #56036;) and the nucmer and mummerplot (--layout --  
1529 filter) functions from MUMmer 3.23 (Kurtz, S et al., 2004). To determine the  
1530 genomic variants, we used HaplotypeCaller from GATK (McKenna, A et al., 2010;  
1531 DePristo, MA et al., 2011; Van der Auwera, GA et al., 2013) (-ploidy 4 -  
1532 genotyping\_mode DISCOVERY').

### 1533 **Genome annotation**

1534 To annotate the *T. ni* genome, we first masked repetitive sequences and then integrated  
1535 multiple sources of evidence to predict gene models. We used RepeatModeler to define  
1536 repeat consensus sequences and RepeatMasker (-s -e ncbi) to mask repetitive  
1537 regions (Smit, AFA et al., 2017). We used RNAMmer (Lagesen, K et al., 2007) to  
1538 predict 8S, 18S, 28S rRNA genes, and Barrnap (<https://github.com/tseemann/barrnap>)  
1539 to predict 5.8S rRNA genes. We used Augustus v3.2.2 (Stanke, M et al., 2006) and  
1540 SNAP (Korf, I, 2004) to computationally predicted gene models. Predicted gene models

1541 were compiled by running six iterations of MAKER (*Campbell, MS et al., 2014*), aided  
1542 with homology evidence of well annotated genes (UniProtKB/Swiss-Prot and Ensembl)  
1543 and of transcripts from related species (*B. mori* [*Suetsugu, Y et al., 2013*] and *D.*  
1544 *melanogaster* [*Atrill, H et al., 2016*]). We used BLAST2GO (*Conesa, A et al., 2005*) to  
1545 integrate results from BLAST, and InterPro (*Mitchell, A et al., 2015*) to assign GO terms  
1546 to each gene. We used MITOS (*Bernt, M et al., 2013*) web server to predict  
1547 mitochondrial genes and WebApollo (*Lee, E et al., 2013*) for manual curation of genes  
1548 of interest. To characterize telomeres, we used (TTAGG)<sub>200</sub> (*Robertson, HM, Gordon,*  
1549 *KHJ, 2006*) as the query to search the *T. ni* genome using BLASTn with the option '-  
1550 dust no' and kept hits longer than 100 nt. The genomic coordinates of these hits were  
1551 extended by 10 kb to obtain the subtelomeric region.

## 1552 **Orthology and evolution**

1553 To place genes into ortholog groups, we compared the predicted proteomes from 21  
1554 species (Supplementary file 5). Orthology assignment was determined using OrthoMCL  
1555 (*Hirose, Y, Manley, JL, 1997*) with default parameters. MUSCLE v3.8.31 (*Edgar, RC,*  
1556 *2004*) was used for strict 1:1:1 orthologs ( $n = 381$ ) to produce sequence alignments.  
1557 Conserved blocks (66,044 amino acids in total) of these alignments were extracted  
1558 using Gblocks v0.91b (*Castresana, J, 2000*) with default parameters, and fed into  
1559 PhyML 3.0 (*Vastenhouw, NL et al., 2010*) (maximum likelihood, bootstrap value set to  
1560 1000) to calculate a phylogenetic tree. The human and mouse predicted proteomes  
1561 were used as an outgroup to root the tree. The tree was viewed using FigTree  
1562 (<http://tree.bio.ed.ac.uk/software/figtree/>) and iTOL (*Shirayama, M et al., 2012*).

## 1563 **Sex determination and sex chromosomes**

1564 To identify sex-linked contigs, we mapped genomic sequence reads from males and  
1565 females to the contigs. Reads with MAPQ scores  $\geq 20$  were used to calculate contig  
1566 coverage, which was then normalized by the median coverage. The distribution of



1567 normalized contig coverage ratios (male:female ratios, M:F ratios) was manually  
1568 checked to empirically determine the thresholds for Z-linked and W-linked contigs (M:F  
1569 ratio >1.5 for Z-linked contigs and M:F ratio < 0.5 for W-linked contigs). Lepidopteran  
1570 *misc* genes were obtained from Lepbase (*Challis, RJ et al., 2016*). Z/AA ratio was  
1571 calculated according to (*Gu, L et al., 2017*).

## 1572 **Gene families for detoxification and chemoreception**

1573 To curate genes related to detoxification and chemoreception, we obtained seed  
1574 alignments from Pfam (*Finn, RD et al., 2016*) and ran hmmbuild to build HMM profiles of  
1575 cytochrome P450 (P450), amino- and carboxy-termini of glutathione-S-transferase  
1576 (GST), carboxylesterase (COE), ATP-binding cassette transporter (ABCs), olfactory  
1577 receptor (OR), gustatory receptor (GR), ionotropic receptor (IR), and odorant binding  
1578 (OBP) proteins, (Supplementary file 6, 7 and 8). We then used these HMM profiles to  
1579 search for gene models in the predicted *T. ni* proteome (hmmsearch, e-value cutoff:  $1 \times$   
1580  $10^{-5}$ ). We also retrieved reference sequences of P450, GST, COE, ABC, OR, GR, IR,  
1581 OBP, and juvenile hormone pathway genes from the literature (*Hekmat-Safe, DS et al.,*  
1582 *2002*; ; *Xavier, B et al., 2005*; *Wanner, KW, Robertson, HM, 2008*; *Yu, Q et al., 2008*;  
1583 *Benton, R et al., 2009*; *Gong, D-P et al., 2009*; *Yu, Q-Y et al., 2009*; *Croset, V et al.,*  
1584 *2010*; *Ai, J et al., 2011*; *Liu, S et al., 2011*; *Dermauw, W, Van Leeuwen, T, 2014*; *van*  
1585 *Schooten, B et al., 2016*). These were aligned to the *T. ni* genome using tBLASTx  
1586 (*Altschul, SF et al., 1990*) and Exonerate (*Slater, GSC, Birney, E, 2005*) to search for  
1587 homologs. Hits were manually inspected to ensure compatibility with RNA-seq data,  
1588 predicted gene models, known protein domains (using CDD [*Marchler-Bauer, A et al.,*  
1589 *2015*]) and homologs from other species. P450 genes were submitted to David Nelson's  
1590 Cytochrome P450 Homepage (*Nelson, DR, 2009*) for nomenclature and classification.  
1591 Sequences and statistics of these genes are in Supplementary files 6, 7 and 8.

1592 To determine the phylogeny of these gene families, we aligned the putative  
1593 protein sequences from *T. ni* and *B. mori* genomes using MUSCLE (Edgar, RC, 2004),  
1594 trimmed the multiple sequence alignments using TrimAl (Capella-Gutiérrez, S et al.,  
1595 2009) (with the option -automated1), and performed phylogenetic analysis (PhyML  
1596 3.0 [Vastenhouw, NL et al., 2010], with parameters: -q --datatype aa --run\_id  
1597 0 --no\_memory\_check -b -2). Phylogenetic trees were visualized using FigTree  
1598 (<http://tree.bio.ed.ac.uk/software/figtree/>).

1599 To curate opsin genes, we used opsin mRNA and peptide sequences from other  
1600 species (Zimyanin, VL et al., 2008; Futahashi, R et al., 2015) to search for homologs in  
1601 *T. ni*. To discriminate opsin genes from other G-protein-coupled receptors, we required  
1602 that the top hit in the NCBI non-redundant database and UniProt were opsins.

### 1603 **Transposon analysis**

1604 To determine transposon age, we calculated the average percent divergence for each  
1605 transposon family: the percent divergence (RepeatMasker) of each transposon copy  
1606 was multiplied by its length, and the sum of all copies were divided by the sum of  
1607 lengths of all copies in the family (Pace, JK, Feschotte, C, 2007). We used TEMP  
1608 (Zhuang, J et al., 2014) to identify transposon insertions in the Hi5 genome.

### 1609 **miRNA and siRNA analysis**

1610 mirDeep2 (Friedländer, MR et al., 2008; Friedlander, MR et al., 2012) with default  
1611 parameters predicted miRNA genes. Predicted miRNA hairpins were required to have  
1612 homology (exact seed matches and BLASTn e-value  $< 1 \times 10^{-5}$ ) to known miRNAs  
1613 and/or miRDeep2 scores  $\geq 10$ . miRNAs were named according to exact seed matches  
1614 and high sequence identities (BLASTn e-value  $< 1 \times 10^{-5}$ ) with known miRNA hairpins.  
1615 To determine the conservation status of *T. ni* miRNAs, putative *T. ni* miRNAs were  
1616 compared with annotated miRNAs from *A. aegypti*, *A. mellifera*, *B. mori*, *D.*

1617 *melanogaster*, *H. sapiens*, *M. musculus*, *M. sexta*, *P. xylostella*, and *T. castaneum*:  
1618 conserved miRNAs were required to have homologous miRNAs beyond Lepidoptera.

1619 To compare siRNA abundance in oxidized and unoxidized small RNA-seq  
1620 libraries, we normalized siRNA read counts to piRNA cluster-mapping reads (piRNA  
1621 cluster read counts had >0.98 Pearson correlation coefficients between oxidized and  
1622 unoxidized libraries in all cases). Because piRNA degradation products can be 20–22 nt  
1623 long, we excluded potential siRNA species that were prefixes of piRNAs (23–35 nt).

1624 To search for viral transcripts in *T. ni*, we downloaded viral protein sequences  
1625 from NCBI (<http://www.ncbi.nlm.nih.gov/genome/viruses/>) and used using tBLASTn to  
1626 map them to the *T. ni* genome and to the transcriptomes of Hi5 cells and five *T. ni*  
1627 tissues. We filtered hits (percent identity  $\geq 0.80$ , e-val  $\leq 1 \times 10^{-20}$ , and alignment length  
1628  $\geq 100$ ) and mapped small RNA-seq reads to the identified viral transcripts.

1629 Candidate genomic hairpins were defined according to (*Okamura, K et al.,*  
1630 *2008b*). Candidate *cis*-NATs were defined according to (*Ghildiyal, M et al., 2008*).

### 1631 piRNA analysis

1632 To determine the genomic coordinates of piRNA-producing loci, we mapped small  
1633 RNAs to the genome as described (*Han, BW et al., 2014*). We then calculated the  
1634 abundance of piRNAs in 5 kb genomic windows. For each window, we counted the  
1635 number of uniquely mapped reads and the number of reads mapped to multiple loci  
1636 (multimappers) by assigning reads using an expectation-maximization algorithm. Briefly,  
1637 each window had the same initial weight. The weight was used to linearly apportion  
1638 multimappers. During the expectation (E) step, uniquely mapped reads were  
1639 unambiguously assigned to genomic windows; multimappers were apportioned to the  
1640 genomic windows they mapped to, according to the weights of these windows. At the  
1641 maximization (M) step, window weights were updated to reflect the number of reads  
1642 each window contained from the E step. The E and M steps were run iteratively until the

1643 Manhattan distance between two consecutive iterations was smaller than 0.1% of the  
1644 total number of reads.

1645 To identify differentially expressed piRNA loci, we used the ppm and rpkm  
1646 values, normalized to the total number of uniquely mapped reads, to measure piRNA  
1647 abundance. For analyses including all mapped reads (uniquely mapped reads and  
1648 multimappers), reads were apportioned by the number of times that they were mapped  
1649 to the genome. To make piRNA loci comparable across tissues, we merged piRNA loci  
1650 from ovary, testis, female and male thorax, and Hi5 cells. For the comparison between  
1651 female and male thoraces, the cluster on tig00001980 was removed as this cluster likely  
1652 corresponds to a mis-assembly. We used Spearman correlations to calculate the  
1653 pairwise correlations of piRNA abundances. As for defining sex-linked contigs, we  
1654 calculated M:F ratios and used the same thresholds to determine whether a piRNA  
1655 cluster was sex-linked. A piRNA locus was considered to be differentially expressed if  
1656 the ratio between the two tissues was  $>2$  or  $<0.5$  and FDR  $<0.1$  (after t-test).

1657 Splice sites were deemed to be supported by RNA-seq data when supported by  
1658 at least one data set. We used AUGUSTUS (*Stanke, M et al., 2006*), with the model  
1659 trained for *T. ni* genome-wide gene prediction, to predict gene models and their splice  
1660 sites in *T. ni* piRNA clusters.

### 1661 **$\beta$ -elimination**

1662 Total RNAs were extracted from Hi5 cells using mirVana kit as described previously.  
1663 We then incubated 100  $\mu$ g total RNA with 25 mM NaIO<sub>4</sub> in borate buffer (148 mM  
1664 Borax, 148 mM Boric acid, pH 8.6) for 30 min at room temperature, beta-elimination was  
1665 performed in 50 mM NaOH at 45°C for 90 min (*Horwich, MD et al., 2007*). The resultant  
1666 RNA was collected by ethanol precipitation.

1667 **sgRNA design**

1668 sgRNAs for the target loci (5'-end of *TnPiwi* and 5'-end of *vasa*) were designed using  
1669 crispr.mit.edu (Hsu, PD et al., 2013) to retrieve all possible guide sequences, and guide  
1670 sequences adjacent to deletion or insertion targets were chosen. Supplementary file 9  
1671 lists guide sequences.

1672 **ssDNA donor purification**

1673 Donor template sequence was produced as a gBlock (Integrated DNA Technologies,  
1674 San Diego, CA, USA). A biotinylated forward primer and a standard reverse primer were  
1675 used in PCR to generate a double-stranded, biotinylated DNA donor. The biotinylated  
1676 DNA was captured on M-280 streptavidin Dynabeads (ThermoFisher), and the  
1677 biotinylated strand was separated from the non-biotinylated strand essentially as  
1678 described in the manufacturer's protocol. Supplemental file 10 provides a detailed  
1679 protocol.

1680 **Transfection of Hi5 cells**

1681 sgRNAs were transcribed using T7 RNA polymerase, gel purified, then incubated with  
1682 Cas9 in serum-free Hi5 culture medium supplemented with 18 mM L-glutamine. The  
1683 resulting sgRNA/Cas9 RNPs were incubated with Trans-IT insect reagent (Mirus Bio,  
1684 Madison, WI, USA) for 15 min at room temperature, then evenly distributed onto 90%  
1685 confluent Hi5 cells. Culture medium was replaced with fresh medium 12 h later.  
1686 Genomic DNA was isolated and analyzed by PCR 48 h later.

1687 **PCR to validate genomic editing in transfected cells**

1688 Forty eight hours after transfection, Hi5 cells from one 90% confluent well of a six-well  
1689 plate (Corning, Corning, NY, USA) were collected, washed once with PBS  
1690 (ThermoFisher) and lysed in 2× PK buffer containing 200 µg/ml proteinase K, extracted  
1691 with phenol:chloroform:isoamyl alcohol (25:24:1), and then genomic DNA collected by

1692 ethanol precipitation. Deletions in *TnPiwi* were detected by PCR using primers flanking  
1693 the deleted region (Supplementary file 9). To confirm deletions by sequencing PCR  
1694 products were resolved by agarose gel electrophoresis, purified (QIAquick Gel  
1695 Extraction Kit, QIAGEN, Germantown, MD, USA), and cloned into pCR-Blunt II-Topo  
1696 vector (ThermoFisher). The recombinant plasmid was transformed into Top10  
1697 competent *E.coli* (ThermoFisher) following supplier's protocol. PCR products amplified  
1698 using M13 (-20) forward and M13 reverse primers from a sample of a single bacterial  
1699 colony were sequenced by GENEWIZ (South Plainfield, NJ, USA).

#### 1700 **Single clone selection**

1701 Wild-type Hi5 cells were seeded into a 96-well Transwell permeable support receiver  
1702 plate (Corning, Corning, NY, USA) at 30% confluence and incubated overnight in serum  
1703 free medium with 100 U/ml penicillin and 100 µg/ml streptomycin. A Transwell  
1704 permeable support insert plate with media in each well was inserted into the receiver  
1705 plate, and a single EGFP-positive cell was sorted into each insert well by FACS. After  
1706 14 days incubation at 27°C, wells were examined for EGFP-positive cell clones using a  
1707 DMI8 fluorescent microscope (Leica).

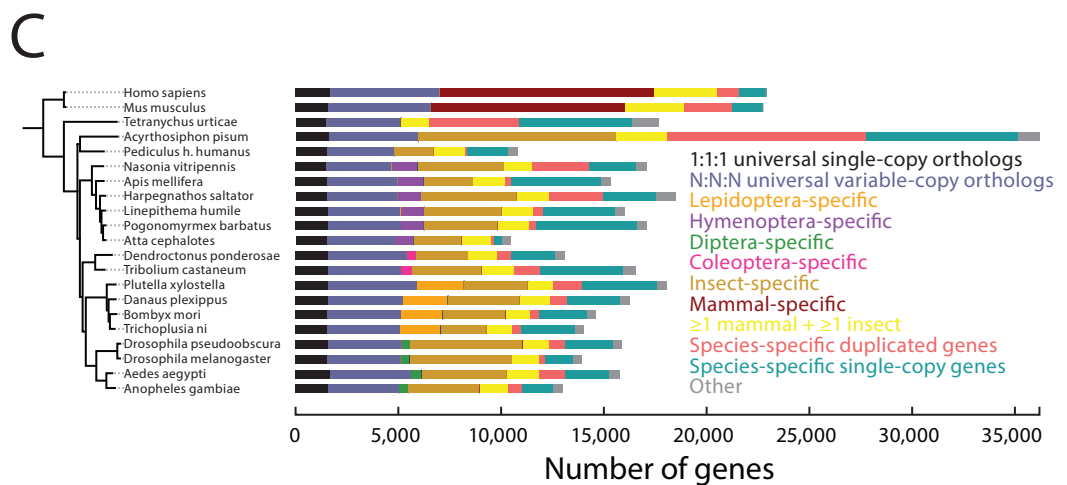
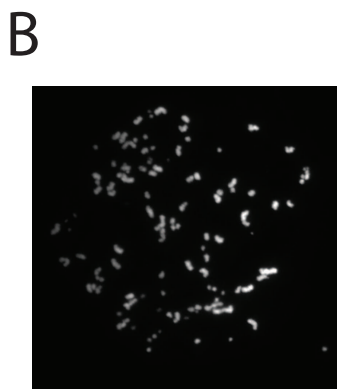
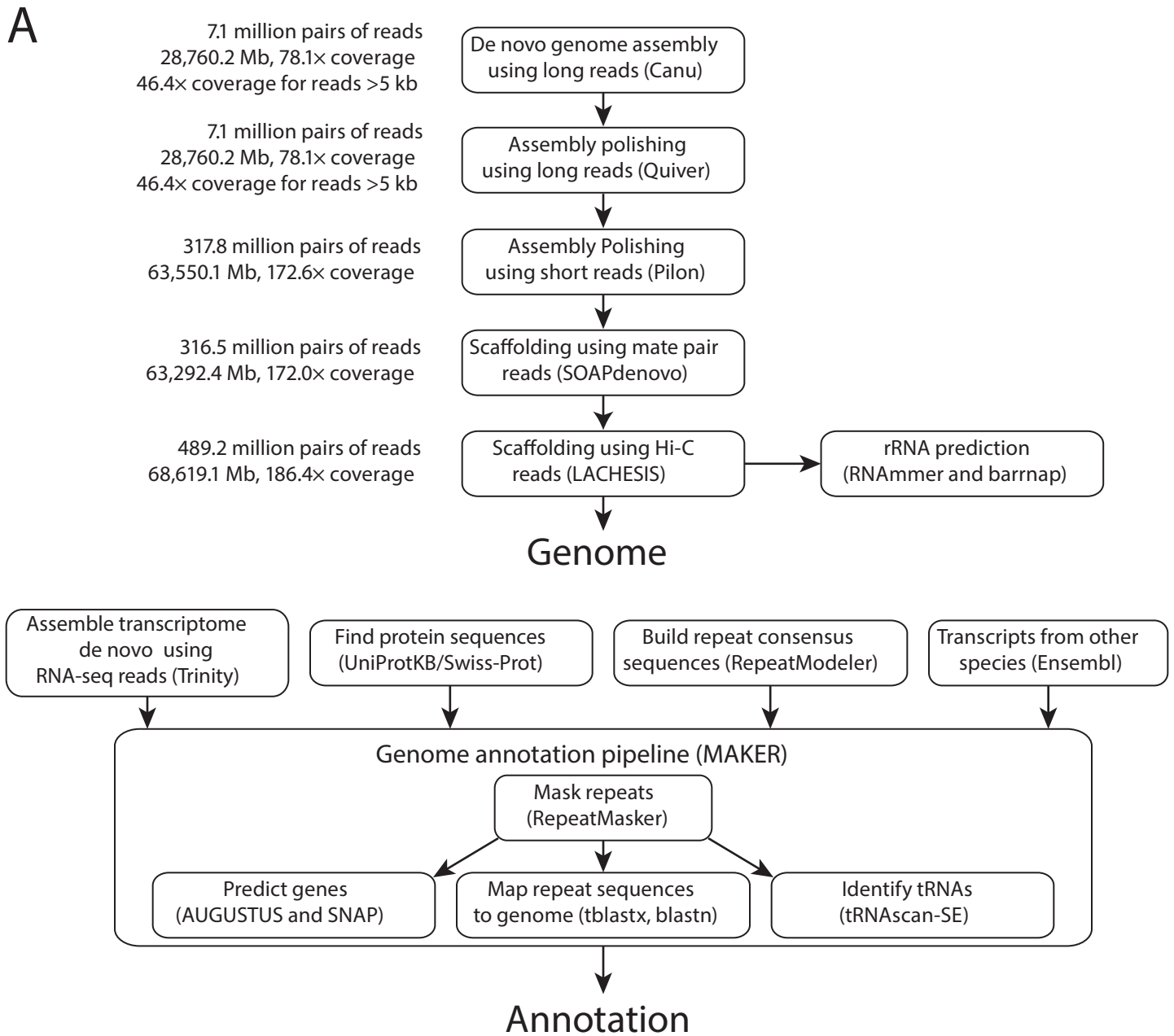
#### 1708 **Immunostaining**

1709 EGFP-HA-Vasa-expressing Hi5 cells were seeded on 22 × 22 mm cover slips (Fisher  
1710 Scientific, Pittsburgh, PA, USA) in a well of a six-well plate (Corning). After cells had  
1711 attached to the coverslip, the medium was removed and cells were washed three times  
1712 with PBS (Gibco). Cells were fixed in 4% (w/v) methanol-free formaldehyde  
1713 (ThermoFisher) in PBS at room temperature for 15 min, washed three times with PBS,  
1714 permeabilized with 0.1% (w/v) Triton X-100 in PBS for 15 min at room temperature, and  
1715 then washed three times with PBS. For antibody labeling, cells were incubated in 0.4%  
1716 (v/v) Photo-Flo in 1× PBS for 10 min at room temperature, then 10 min in 0.1% (w/v)  
1717 Triton X-100 in PBS and 10 min in 1× ADB-PBS (3 mg/ml bovine serum albumen, 1%

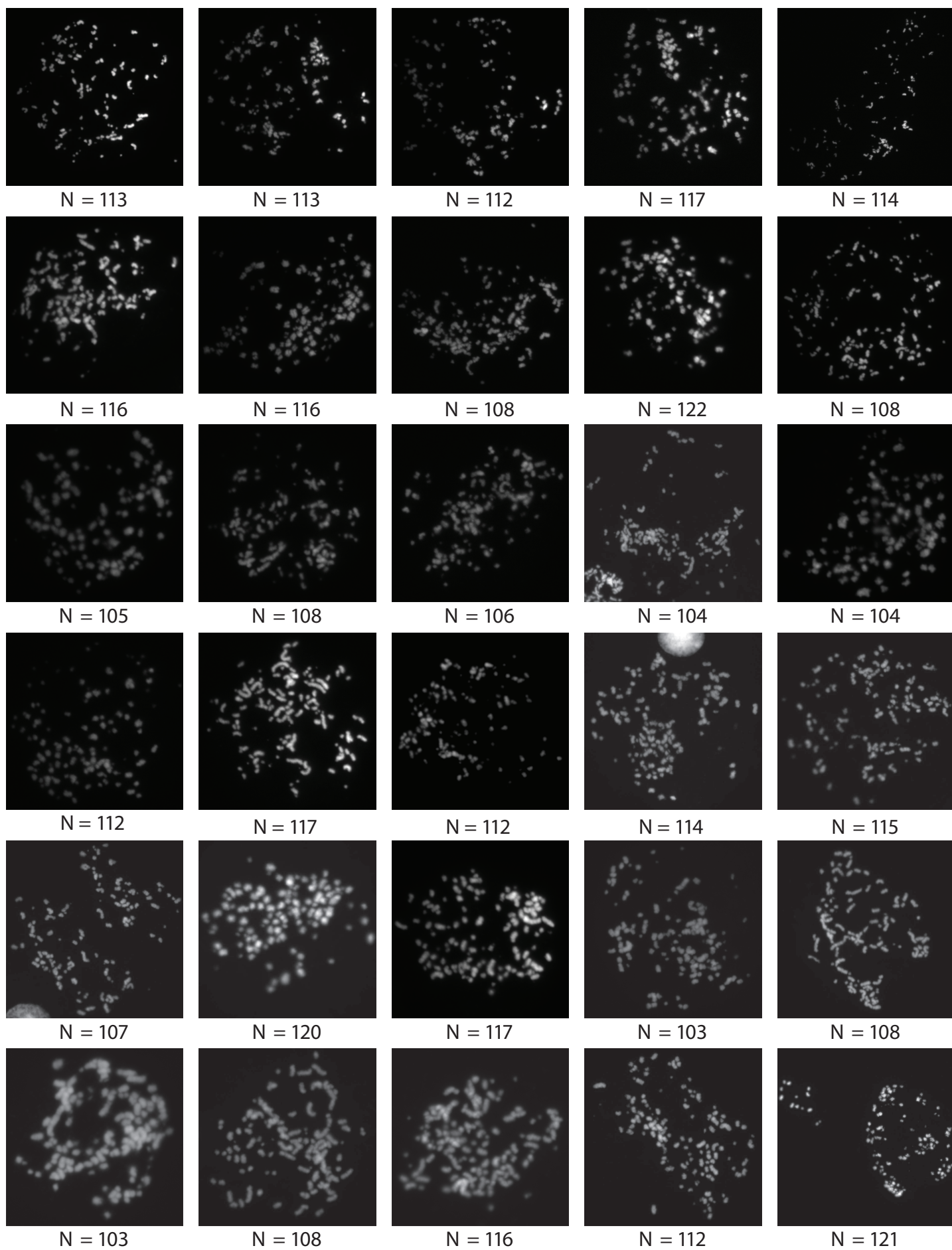
1718 (v/v) donkey serum, 0.005% (w/v) Triton X-100 in 1× PBS). Next, cells were incubated  
1719 with primary antibodies (mouse anti-GFP antibody (GFP-1D2, Developmental Studies  
1720 Hybridoma Bank, Iowa City, IA, USA) and rabbit anti-HA Tag antibody (C29F4, Cell  
1721 Signaling, Danvers, MA, USA), diluted 1:200 in ADB (30 mg/ml BSA, 10% (v/v) donkey  
1722 serum, 0.05% (w/v) Triton X-100 in 1× PBS) at 4°C overnight. After three washes in  
1723 PBS, cells were incubated sequentially in 0.4% (v/v) Photo-Flo in 1× PBS, 0.1% (w/v)  
1724 Triton X-100 in PBS, and 1× ADB-PBS, each for 10 min at room temperature. Cells  
1725 were then incubated with secondary Alexa Fluor 488-labeled donkey anti-mouse  
1726 (ThermoFisher) and Alexa Fluor 680-labeled donkey anti-rabbit (ThermoFisher)  
1727 antibodies, diluted 1:500 in ADB at room temperature for one hour. After washing three  
1728 times with 0.4% (v/v) Photo-Flo in 1× PBS and once with 0.4% (v/v) Photo-Flo in water,  
1729 coverslips were air dried in the dark at room temperature. Slides were mounted in  
1730 ProLong Gold Antifade Mountant with DAPI and examined by confocal microscopy  
1731 (TCS SP5 II Laser Scanning Confocal, Leica).

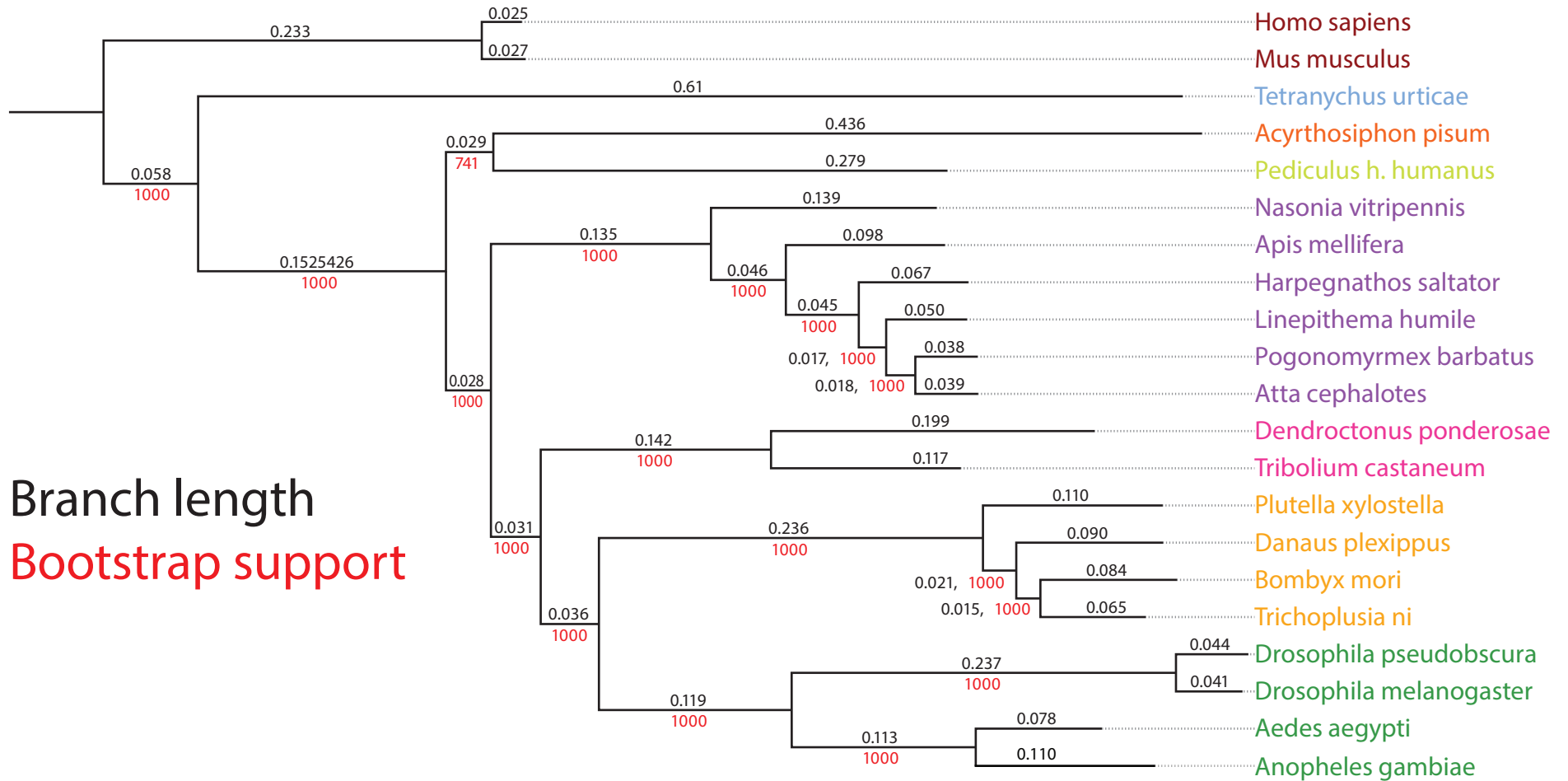
### 1732 **Data deposition**

1733 The *T. ni* Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank  
1734 under the accession NKQN00000000. The version described here is version  
1735 NKQN01000000. All sequencing data are available through the NCBI Sequence Read  
1736 Archive under the accession number PRJNA336361. Further details are available at the  
1737 Cabbage Looper Database (<http://cabbagelooper.org/>).



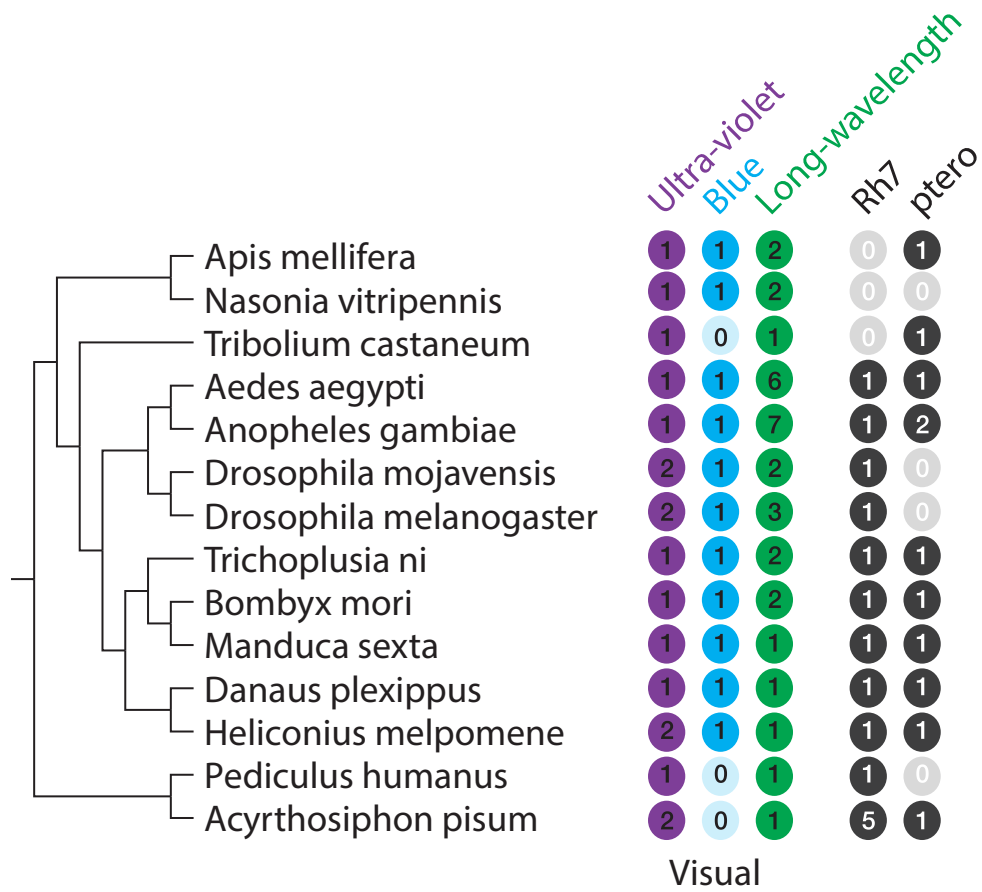




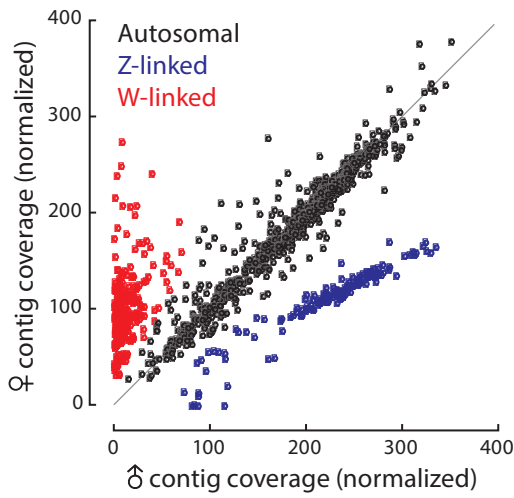


Branch length  
Bootstrap support

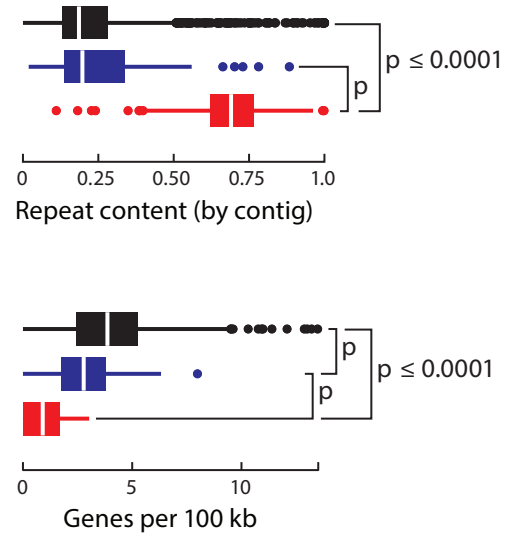
Insects: **Diptera**, **Lepidoptera**, **Coleoptera**, **Hymenoptera**, **Phthiraptera**, **Homoptera**  
**Arachnida**  
**Mammals**



A



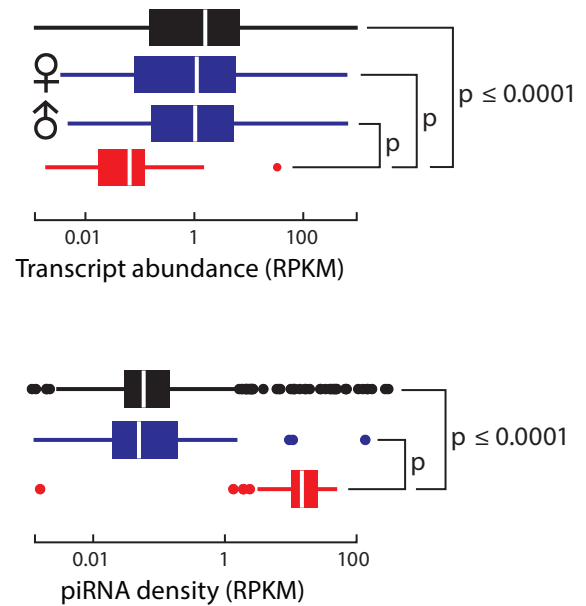
B



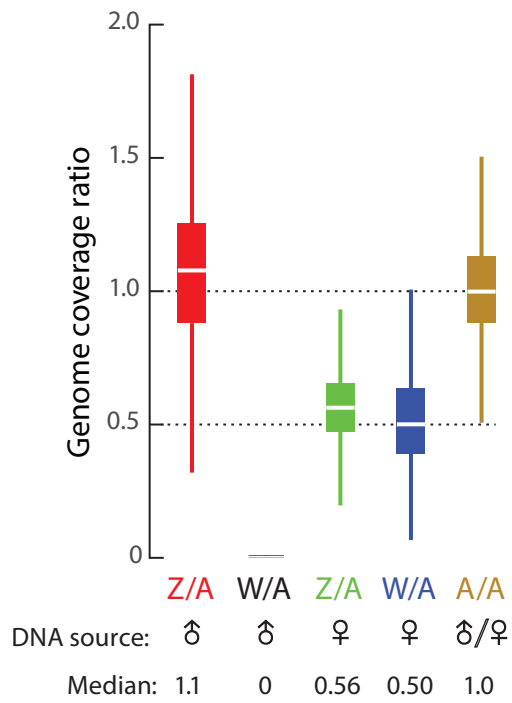
C

	294	297	300	
Trichoplusia ni	L	C	E	S
Bombyx mori	D	C	N	N
Trilocho varians	D	C	N	N
Ostrinia furnacalis	M	C	E	Q
Papilio polytes	H	C	E	N
Bicyclus anynana	L	C	E	Q
Danaus plexippus	Q	C	D	L
Phoebis sennae	L	C	D	S
Heliconius melpomene	H	C	E	K
	C	C	Q	R

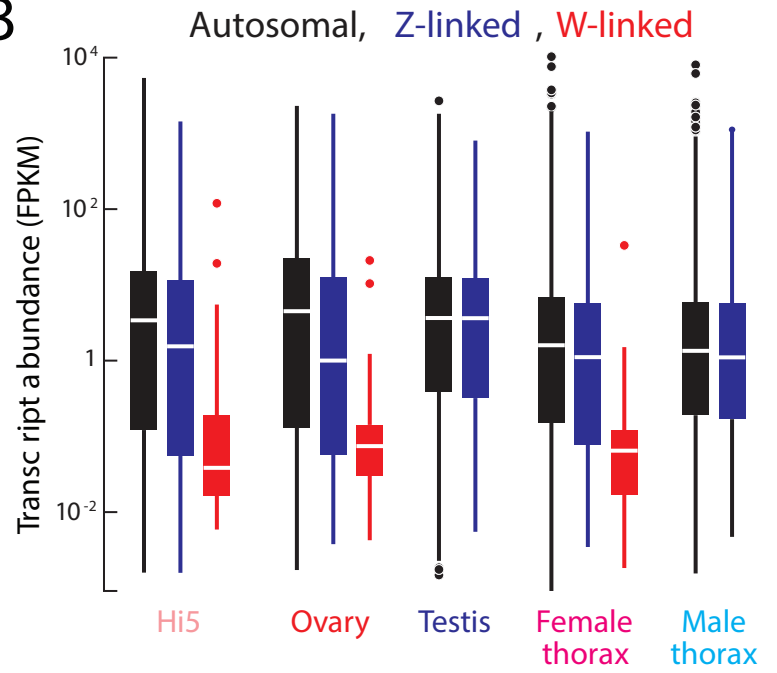
Moths  
Butterflies



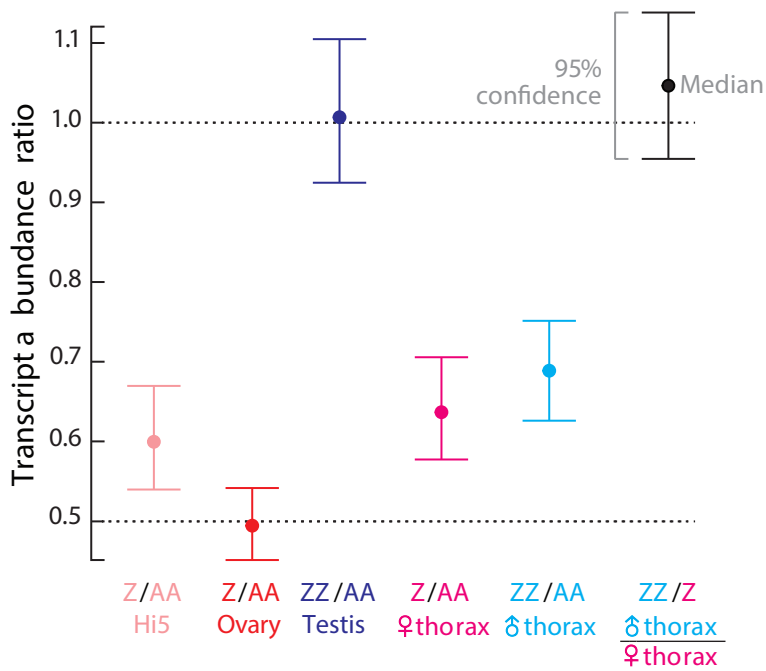
**A**



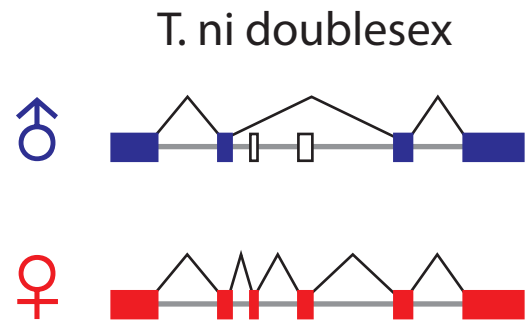
**B**



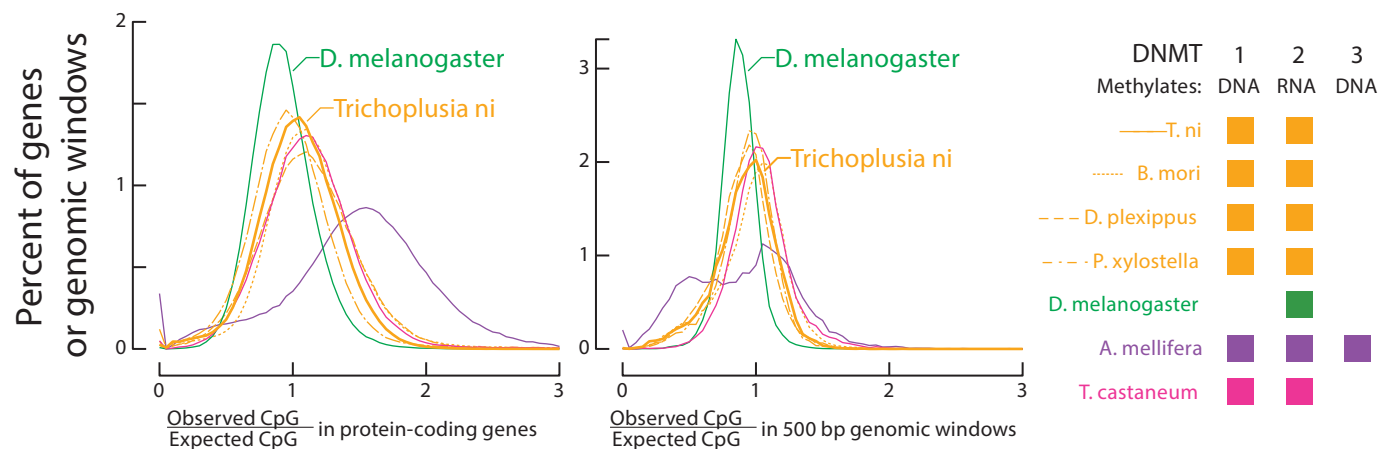
**C**



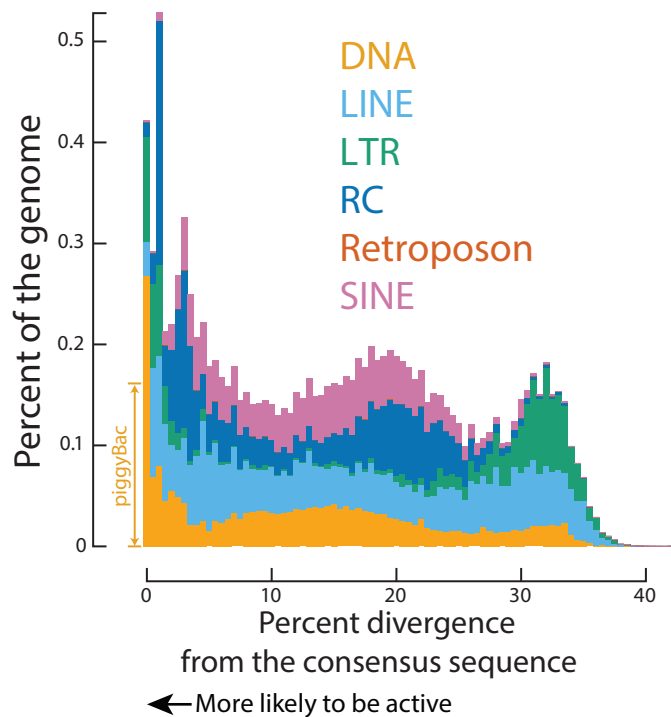
**D**



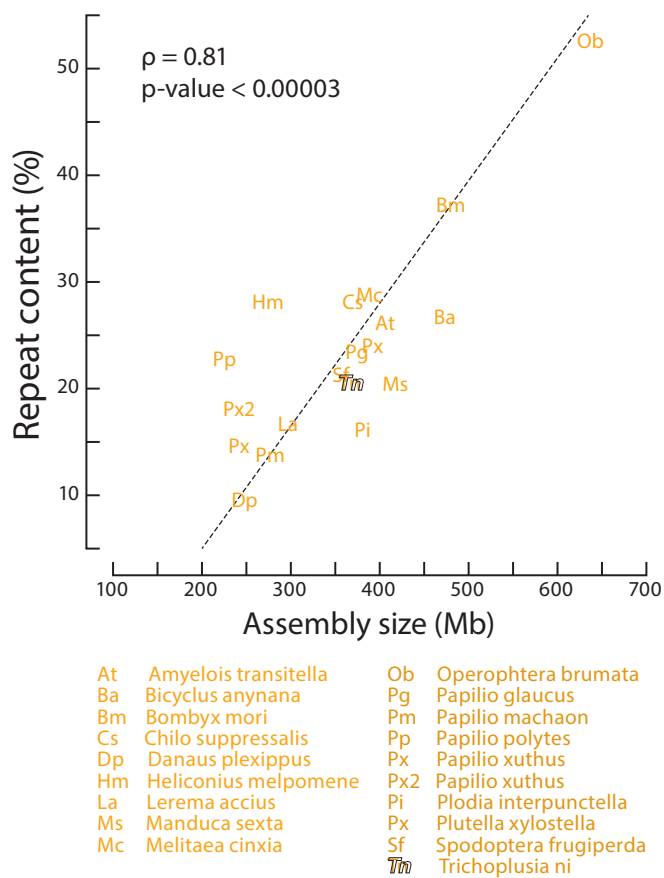
**A**

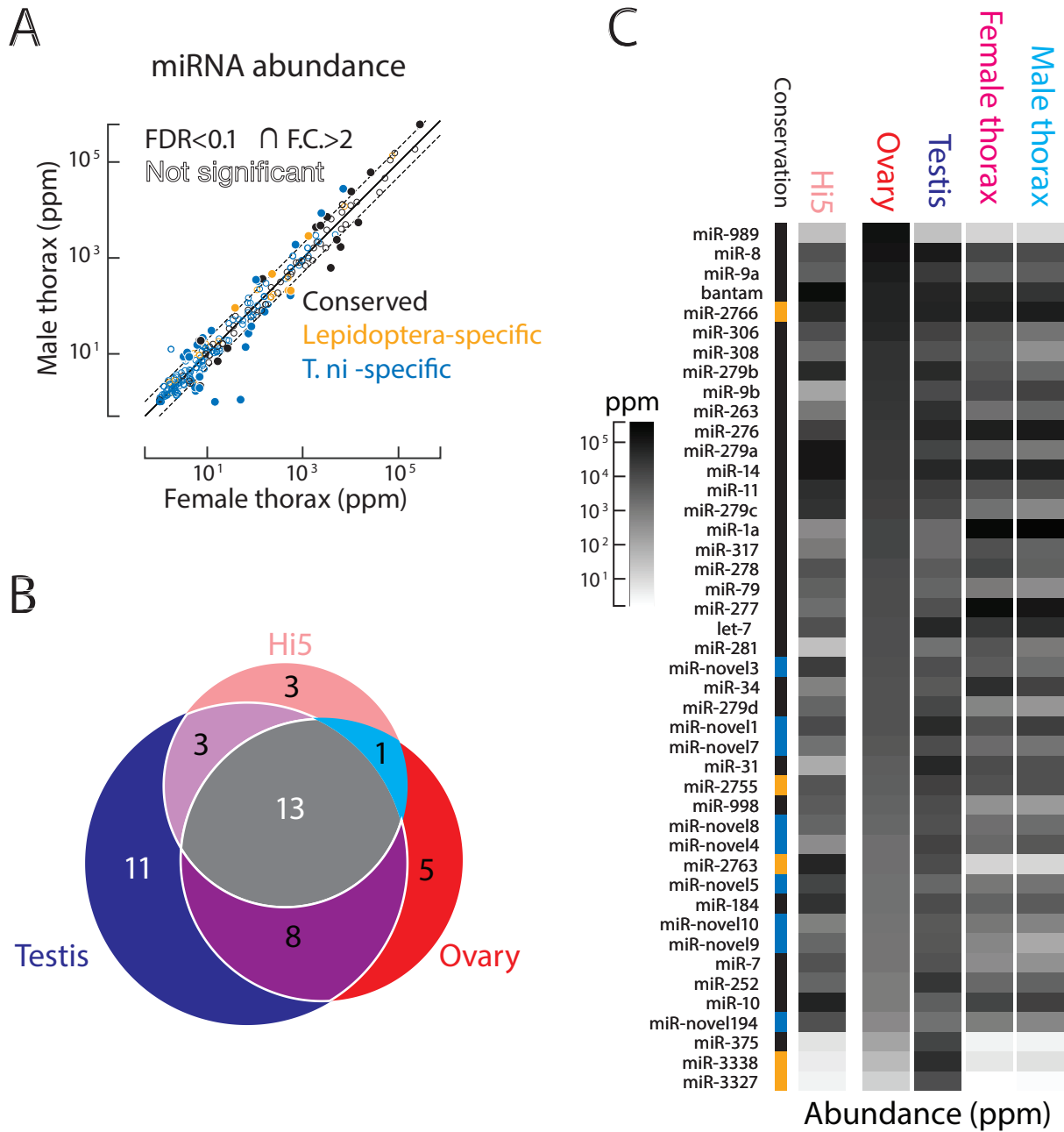


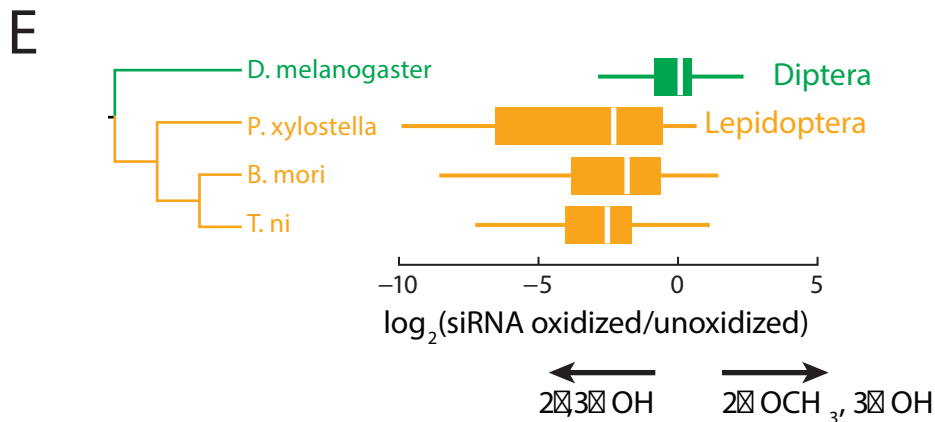
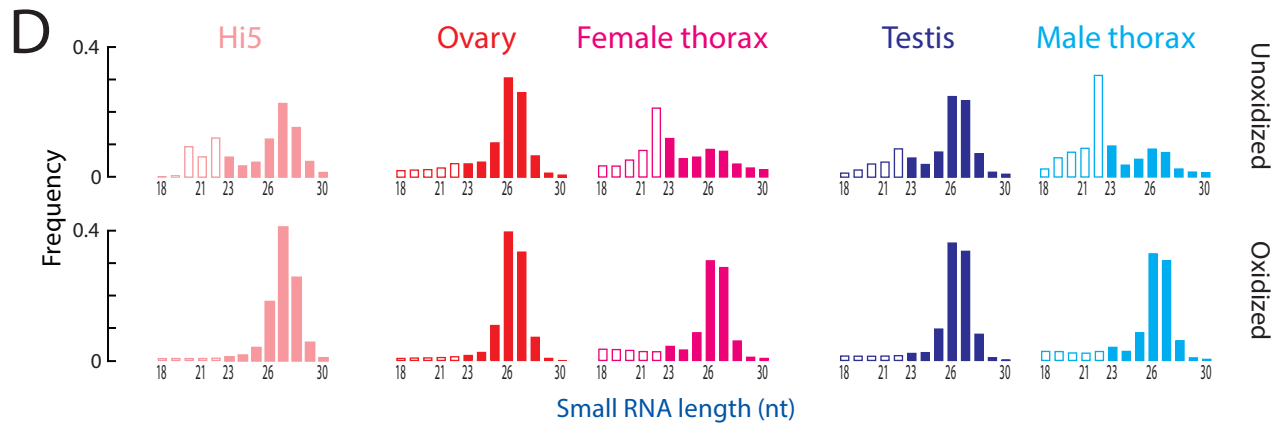
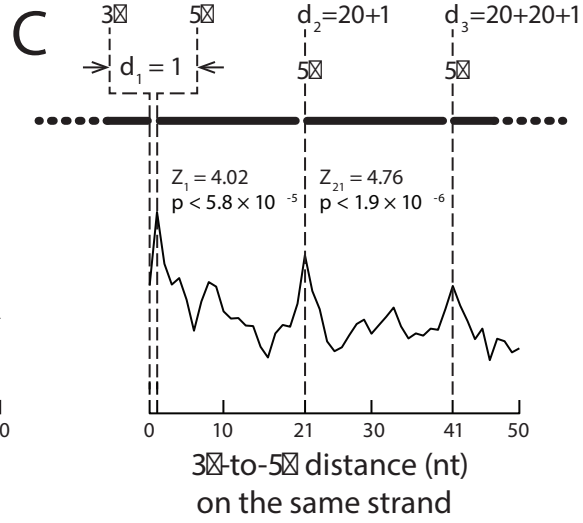
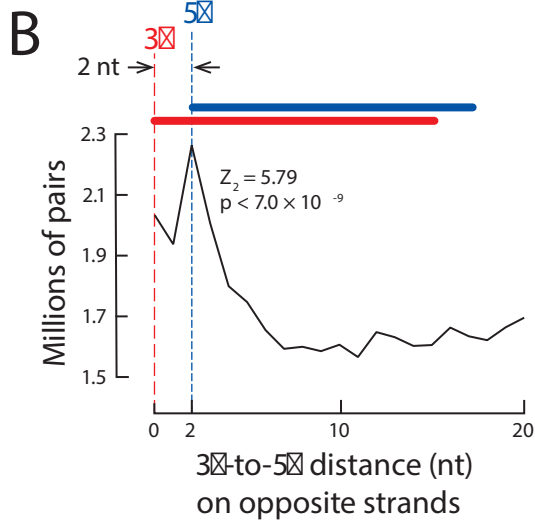
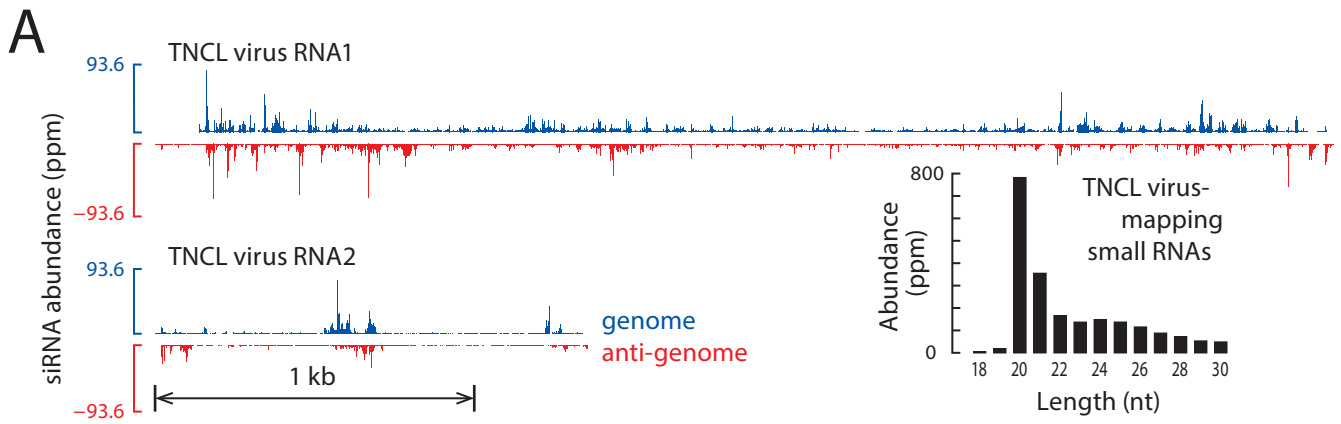
**B**



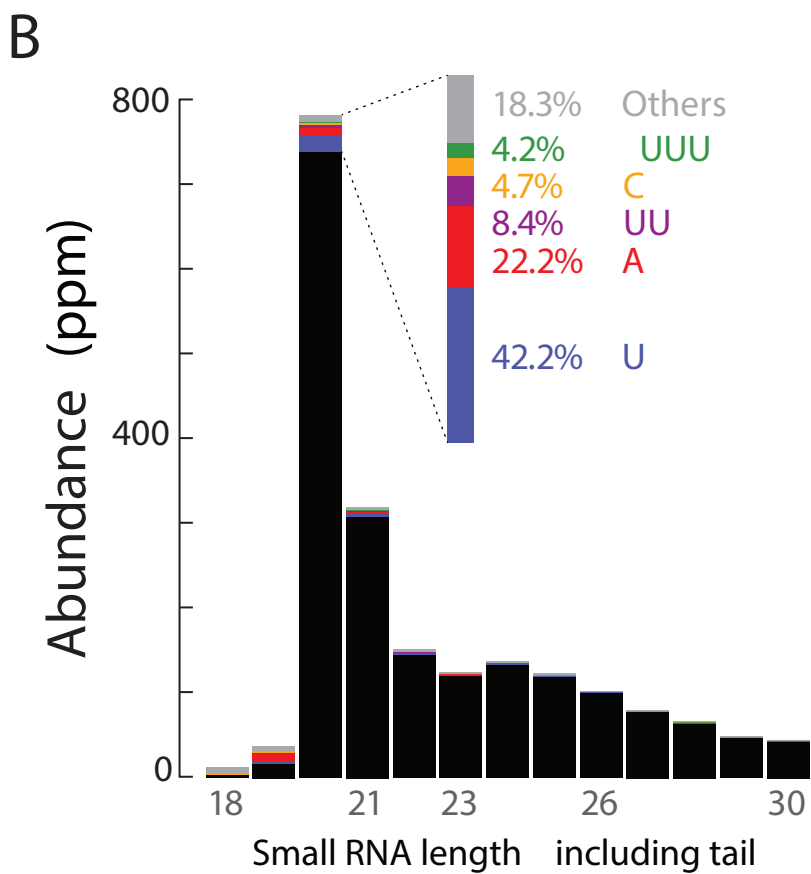
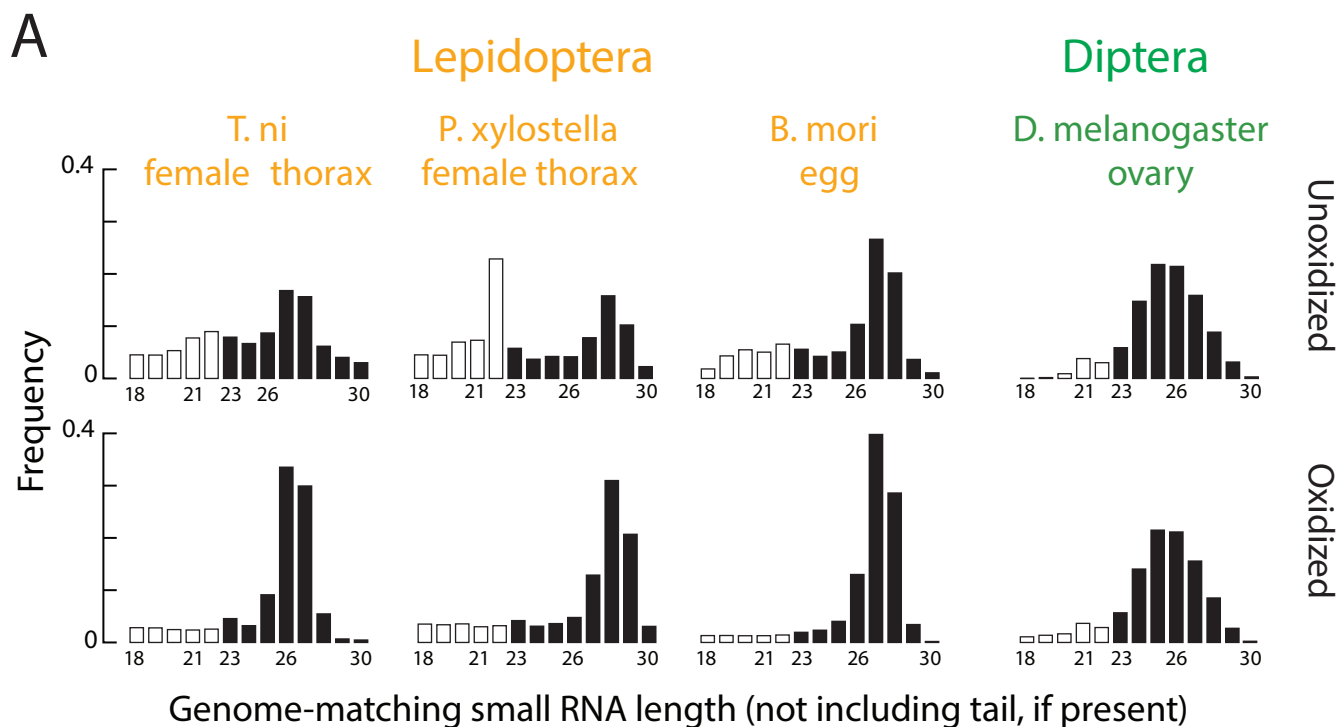
**C**



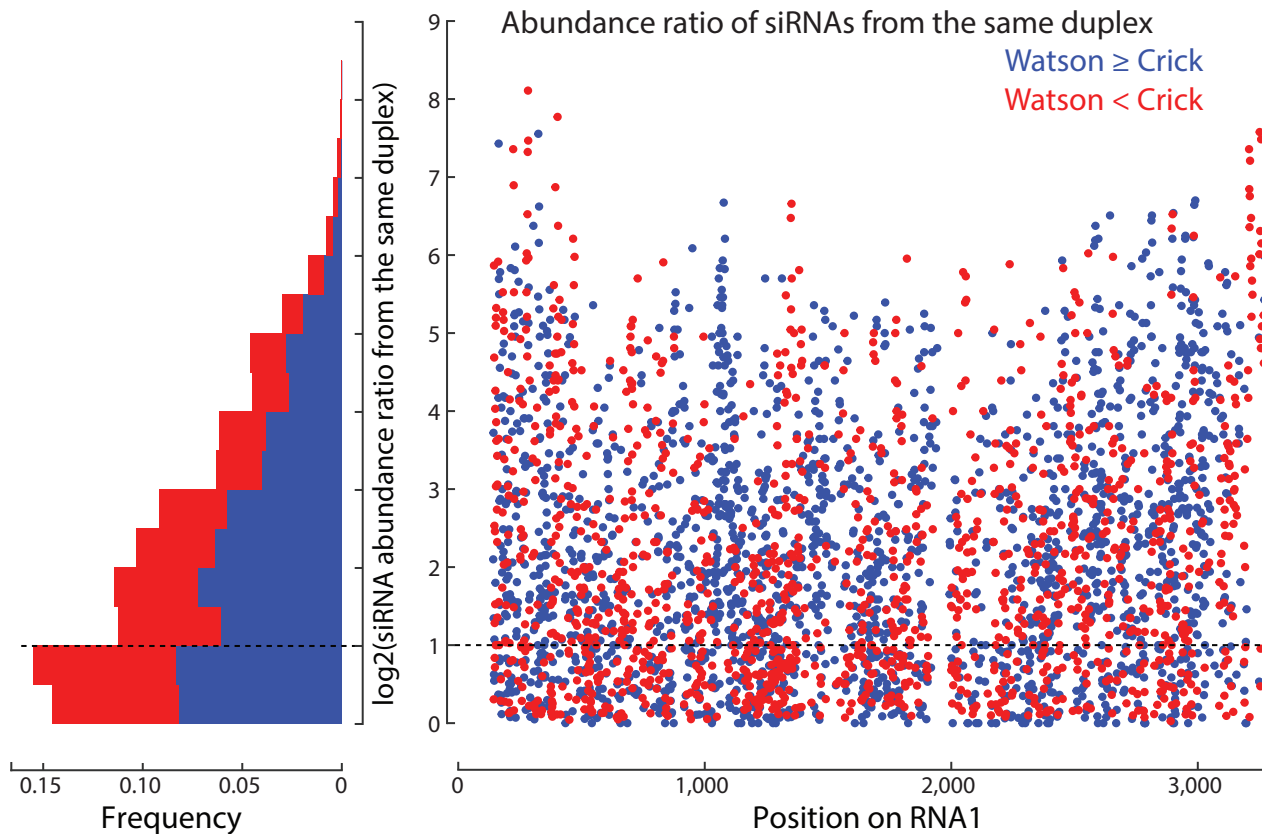




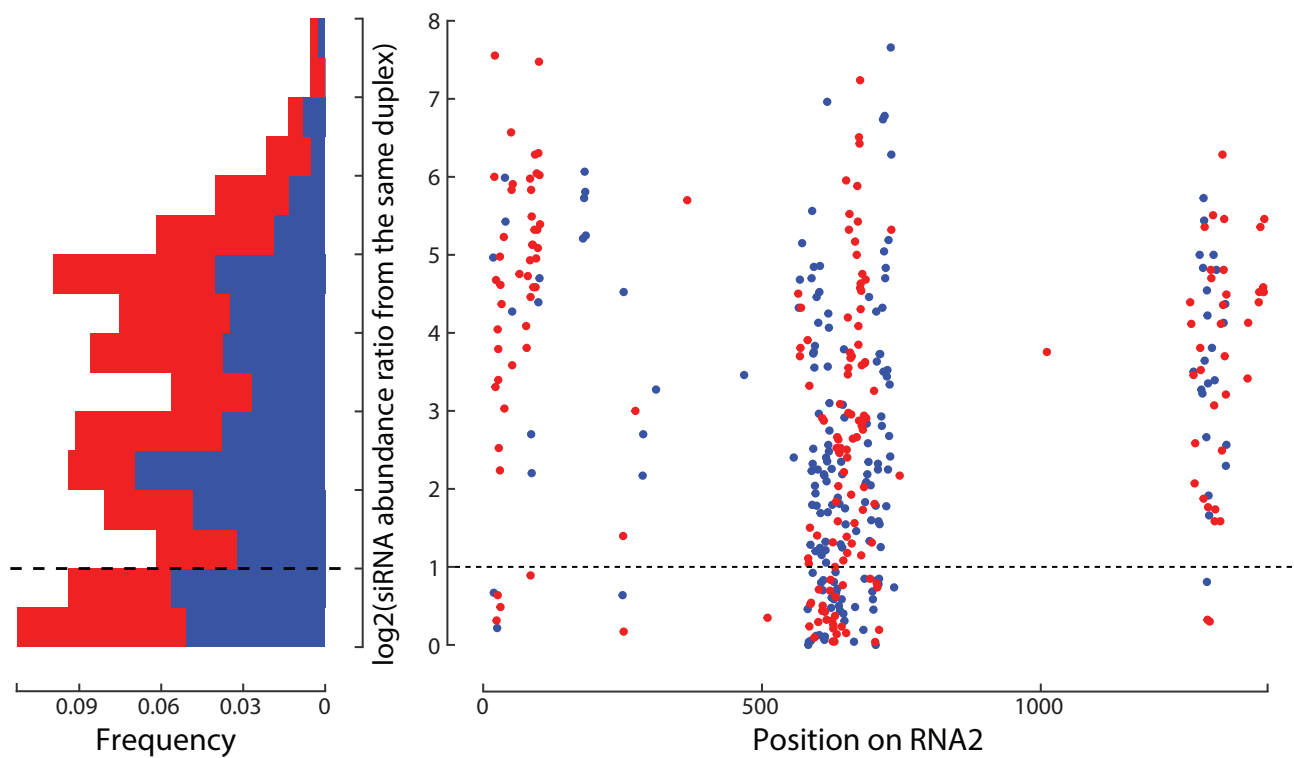




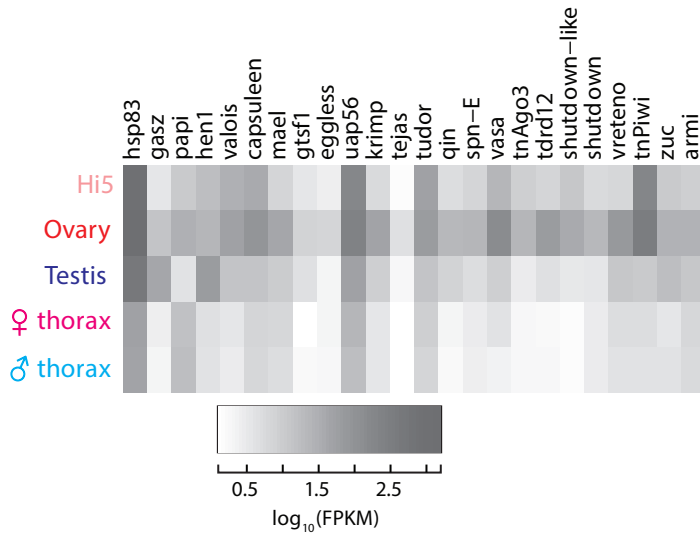
A



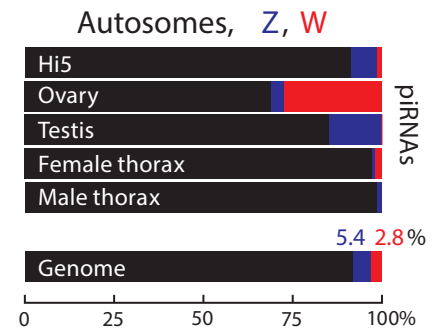
B



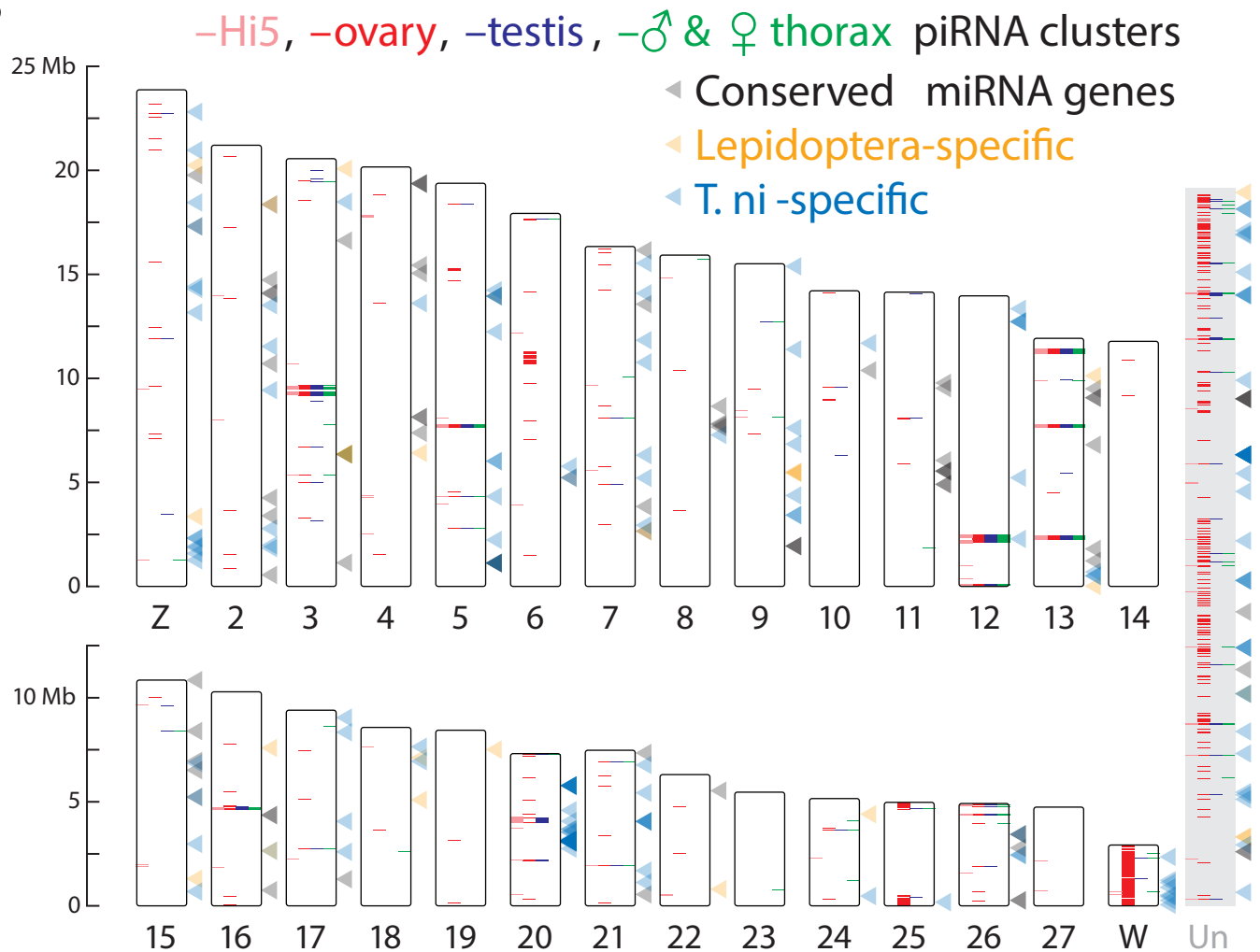
A

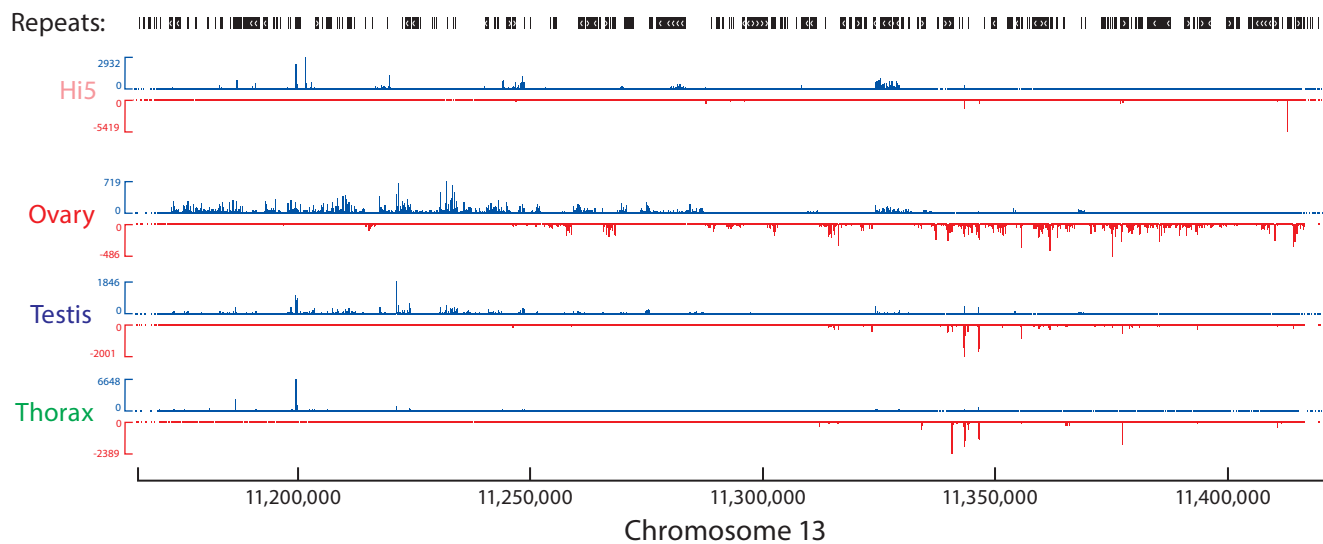
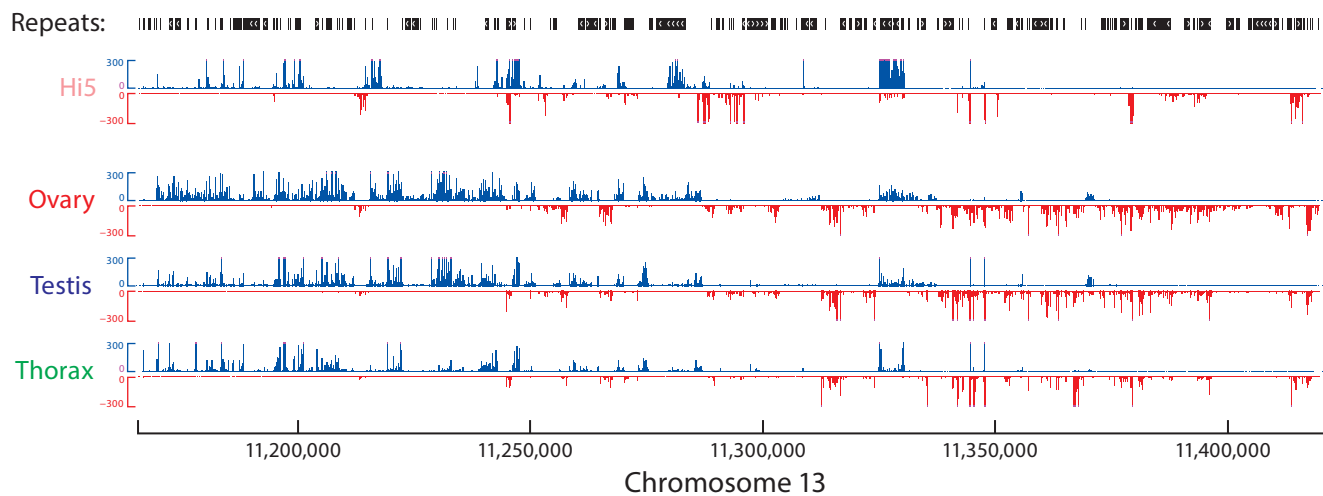


C

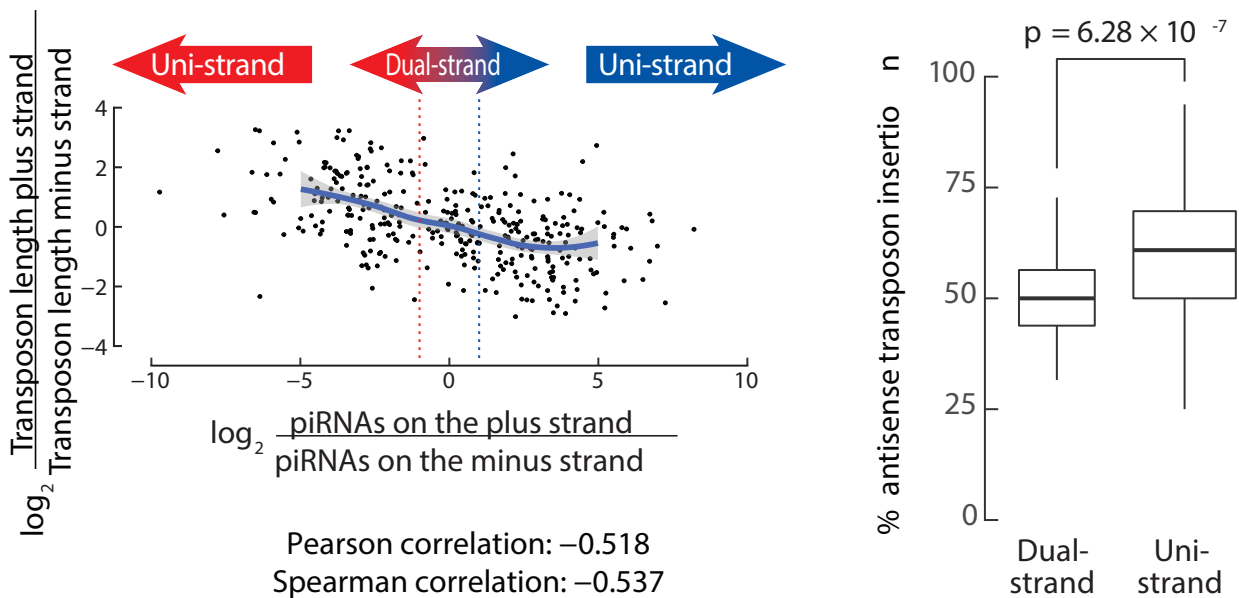


B

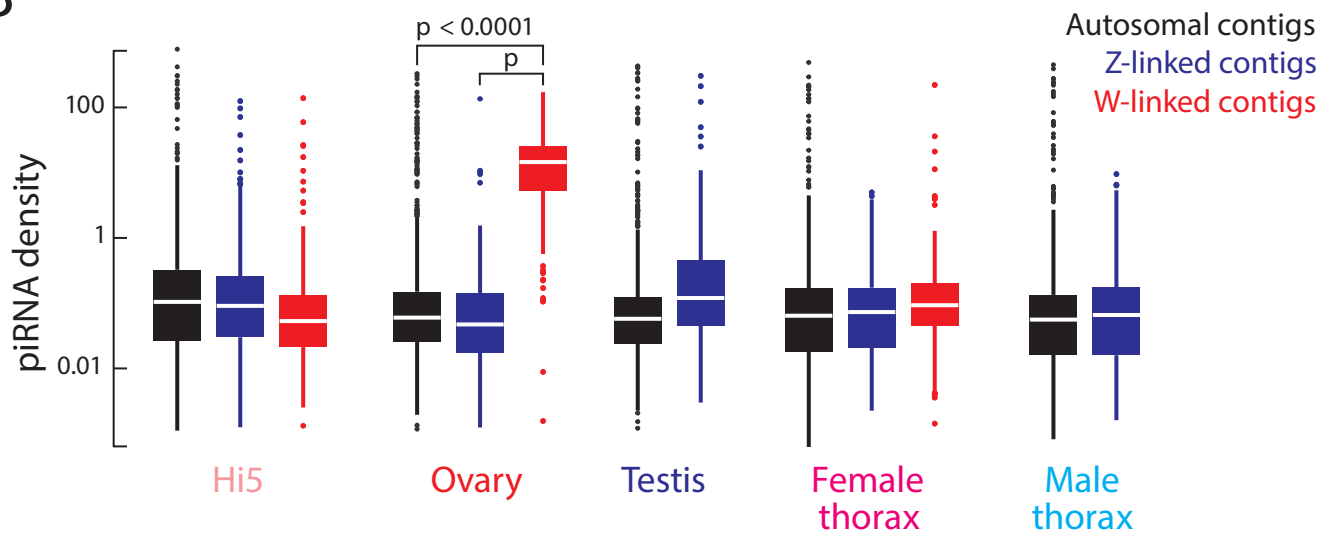




**A**



**B**



**C**

