

# 1 The genome of the cereal pest *Sitophilus oryzae*: a 2 transposable element haven

3

4 PARISOT Nicolas<sup>1,§</sup>, VARGAS-CHAVEZ Carlos<sup>1,2,†,§</sup>, GOUBERT Clément<sup>3,4,‡,§</sup>, BAA-  
5 PUYOULET Patrice<sup>1</sup>, BALMAND Séverine<sup>1</sup>, BERANGER Louis<sup>1</sup>, BLANC Caroline<sup>1</sup>,  
6 BONNAMOUR Aymeric<sup>1</sup>, BOULESTEIX Matthieu<sup>3</sup>, BURLET Nelly<sup>3</sup>, CALEVRO Federica<sup>1</sup>,  
7 CALLAERTS Patrick<sup>5</sup>, CHANCY Théo<sup>1</sup>, CHARLES Hubert<sup>1,6</sup>, COLELLA Stefano<sup>1,§</sup>, DA  
8 SILVA BARBOSA André<sup>7</sup>, DELL'AGLIO Elisa<sup>1</sup>, DI GENOVA Alex<sup>3,6,8</sup>, FEBVAY Gérard<sup>1</sup>,  
9 GABALDON Toni<sup>9,10,11</sup>, GALVÃO FERRARINI Mariana<sup>1</sup>, GERBER Alexandra<sup>12</sup>, GILLET  
10 Benjamin<sup>13</sup>, HUBLEY Robert<sup>14</sup>, HUGHES Sandrine<sup>13</sup>, JACQUIN-JOLY Emmanuelle<sup>7</sup>, MAIRE  
11 Justin<sup>1,-</sup>, MARCET-HOUBEN Marina<sup>9</sup>, MASSON Florent<sup>1,£</sup>, MESLIN Camille<sup>7</sup>, MONTAGNE  
12 Nicolas<sup>7</sup>, MOYA Andrés<sup>2,15</sup>, RIBEIRO DE VASCONCELOS Ana Tereza<sup>12</sup>, RICHARD  
13 Gautier<sup>16</sup>, ROSEN Jeb<sup>14</sup>, SAGOT Marie-France<sup>3,6</sup>, SMIT Arian F.A.<sup>14</sup>, STORER Jessica M.<sup>14</sup>,  
14 VINCENT-MONEGAT Carole<sup>1</sup>, VALLIER Agnès<sup>1</sup>, VIGNERON Aurélien<sup>1,#</sup>, ZAIDMAN-REMY  
15 Anna<sup>1</sup>, ZAMOUM Waël<sup>1</sup>, VIEIRA Cristina<sup>3,6,\*</sup>, REBOLLO Rita<sup>1,\*</sup>, LATORRE Amparo<sup>2,15,\*</sup> and  
16 HEDDI Abdelaziz<sup>1,\*</sup>

17

18 <sup>1</sup> Univ Lyon, INSA Lyon, INRAE, BF2I, UMR 203, 69621 Villeurbanne, France.

19 <sup>2</sup> Institute for Integrative Systems Biology (I2SySBio), Universitat de València and Spanish  
20 Research Council (CSIC), València, Spain.

21 <sup>3</sup> Laboratoire de Biométrie et Biologie Evolutive, UMR5558, Université Lyon 1, Université  
22 Lyon, Villeurbanne, France.

23 <sup>4</sup> Department of Molecular Biology and Genetics, 526 Campus Rd, Cornell University, Ithaca,  
24 New York 14853, USA.

25 <sup>5</sup> KU Leuven, University of Leuven, Department of Human Genetics, Laboratory of  
26 Behavioral and Developmental Genetics, B-3000, Leuven, Belgium.

27 <sup>6</sup> ERABLE European Team, INRIA, Rhône-Alpes, France.

28 <sup>7</sup> INRAE, Sorbonne Université, CNRS, IRD, UPEC, Université de Paris, Institute of Ecology  
29 and Environmental Sciences of Paris, Versailles, France.

30 <sup>8</sup> Instituto de Ciencias de la Ingeniería, Universidad de O'Higgins, Rancagua, Chile.

31 <sup>9</sup> Life Sciences, Barcelona Supercomputing Centre (BSC-CNS), Barcelona, Spain.

32 <sup>10</sup> Mechanisms of Disease. Institute for Research in Biomedicine (IRB), Barcelona, Spain.

33 <sup>11</sup> Institut Catalán de Recerca i Estudis Avançats (ICREA), Barcelona, Spain.

34 <sup>12</sup> Laboratório de Bioinformática, Laboratório Nacional de Computação Científica, Petrópolis,  
35 Brazil.

36 <sup>13</sup> Institut de Génomique Fonctionnelle de Lyon (IGFL), Université de Lyon, Ecole Normale  
37 Supérieure de Lyon, CNRS UMR 5242, Lyon, France.

38 <sup>14</sup> Institute for Systems Biology, Seattle, WA, USA.

39 <sup>15</sup> Foundation for the Promotion of Sanitary and Biomedical Research of Valencian  
40 Community (FISABIO), València, Spain.

41 <sup>16</sup> IGEPP, INRAE, Institut Agro, Université de Rennes, Domaine de la Motte, 35653 Le Rheu,  
42 France.

43

44 † Present address: Institute of Evolutionary Biology (IBE), CSIC-Universitat Pompeu Fabra,  
45 Barcelona, Spain.

46 ‡ Present address: Human Genetics, McGill University, Montreal, QC, Canada.

47 § Present address: LSTM, Laboratoire des Symbioses Tropicales et Méditerranéennes, IRD,  
48 CIRAD, INRAE, SupAgro, Univ Montpellier, Montpellier, France.

49 £ Present address: Global Health Institute, School of Life Sciences, Ecole Polytechnique  
50 Fédérale de Lausanne (EPFL), Lausanne 1015, Switzerland.

51 - Present address: School of BioSciences, The University of Melbourne, Parkville, VIC 3010,  
52 Australia.

53 # Present address: Department of Evolutionary Ecology, Institute for Organismic and  
54 Molecular Evolution, Johannes Gutenberg University, 55128 Mainz, Germany.

55

56 § Authors contributed equally to this work.

57 \* Authors contributed equally to this work.

58

59 NP: nicolas.parisot@insa-lyon.fr

60 CV-C: carlos.vargas@uv.es

61 CG: goubert.clement@gmail.com

62 PB-P: patrice.baa-puyoulet@inrae.fr

63 SB: severine.balmand@inrae.fr

64 LB: beranger.louis.bio@gmail.com

65 CB: blancbillard.caroline@gmail.com

66 AB: aymeric.bonnamour@gmail.com

67 MB: matthieu.boulesteix@univ-lyon1.fr

68 NB: nelly.burlet@univ-lyon1.fr

69 FC: federica.calevro@insa-lyon.fr

70 PC: patrick.callaerts@kuleuven.be

71 TC: theo.chancy@hotmail.fr

72 HC: [hubert.charles@insa-lyon.fr](mailto:hubert.charles@insa-lyon.fr)  
73 SC: [stefano.colella@inrae.fr](mailto:stefano.colella@inrae.fr)  
74 ASB: [andredsb.2b@gmail.com](mailto:andredsb.2b@gmail.com)  
75 ED: [elisa.dell-aglio@insa-lyon.fr](mailto:elisa.dell-aglio@insa-lyon.fr)  
76 ADG: [digenova@gmail.com](mailto:digenova@gmail.com)  
77 GF: [gerard.febvay@inrae.fr](mailto:gerard.febvay@inrae.fr)  
78 TG: [toni.gabaldon@bsc.es](mailto:toni.gabaldon@bsc.es)  
79 MGF: [mari.ferrarini@gmail.com](mailto:mari.ferrarini@gmail.com)  
80 AG: [alegerber@Incc.br](mailto:alegerber@Incc.br)  
81 BG: [benjamin.gillet@ens-lyon.fr](mailto:benjamin.gillet@ens-lyon.fr)  
82 RH: [robert.hubley@isbscience.org](mailto:robert.hubley@isbscience.org)  
83 SH: [sandrine.hughes@ens-lyon.fr](mailto:sandrine.hughes@ens-lyon.fr)  
84 EJ-J: [emmanuelle.joly@inrae.fr](mailto:emmanuelle.joly@inrae.fr)  
85 JM: [justin.maire@animelb.edu.au](mailto:justin.maire@animelb.edu.au)  
86 MM-H: [marina.marcet@bsc.es](mailto:marina.marcet@bsc.es)  
87 FM: [masson.ff@gmail.com](mailto:masson.ff@gmail.com)  
88 CM: [camille.meslin@inrae.fr](mailto:camille.meslin@inrae.fr)  
89 NM: [nicolas.montagne@sorbonne-universite.fr](mailto:nicolas.montagne@sorbonne-universite.fr)  
90 AM: [andres.moya@uv.es](mailto:andres.moya@uv.es)  
91 ATRV: [atrv@Incc.br](mailto:atrv@Incc.br)  
92 GR: [gautier.richard@inrae.fr](mailto:gautier.richard@inrae.fr)  
93 JR: [jeb.rosen@isbscience.org](mailto:jeb.rosen@isbscience.org)  
94 M-FS: [Marie-France.Sagot@inria.fr](mailto:Marie-France.Sagot@inria.fr)  
95 AFAS: [arian.smit@isbscience.org](mailto:arian.smit@isbscience.org)  
96 JMS: [jessica.storer@isbscience.org](mailto:jessica.storer@isbscience.org)  
97 CV-M: [carole.monegat@insa-lyon.fr](mailto:carole.monegat@insa-lyon.fr)  
98 AVa: [agnes.vallier@inrae.fr](mailto:agnes.vallier@inrae.fr)  
99 AVi: [avignero@uni-mainz.de](mailto:avignero@uni-mainz.de)

100 AZ-R: [anna.zaidman@insa-lyon.fr](mailto:anna.zaidman@insa-lyon.fr)

101 WZ: [wzamoum@gmail.com](mailto:wzamoum@gmail.com)

102 CV: [cristina.vieira@univ-lyon1.fr](mailto:cristina.vieira@univ-lyon1.fr)

103 RR: [rita.rebollo@inrae.fr](mailto:rita.rebollo@inrae.fr)

104 AL: [amparo.latorre@uv.es](mailto:amparo.latorre@uv.es)

105 AH: [abdelaziz.heddi@insa-lyon.fr](mailto:abdelaziz.heddi@insa-lyon.fr)

106

## 107 Abstract

### 108 Background

109 Among beetles, the rice weevil *Sitophilus oryzae* is one of the most important pests causing  
110 extensive damage to cereal in fields and to stored grains. *S. oryzae* has an intracellular  
111 symbiotic relationship (endosymbiosis) with the Gram-negative bacterium *Sodalis*  
112 *pierantonius* and is a valuable model to decipher host-symbiont molecular interactions.

### 113 Results

114 We sequenced the *Sitophilus oryzae* genome using a combination of short and long reads to  
115 produce the best assembly for a Curculionidae species to date. We show that *S. oryzae* has  
116 undergone successive bursts of transposable element (TE) amplification, representing 72%  
117 of the genome. In addition, we show that many TE families are transcriptionally active, and  
118 changes in their expression are associated with insect endosymbiotic state. *S. oryzae* has  
119 undergone a high gene expansion rate, when compared to other beetles. Reconstruction of  
120 host-symbiont metabolic networks revealed that, despite its recent association with cereal  
121 weevils (30 Kyear), *S. pierantonius* relies on the host for several amino acids and  
122 nucleotides to survive and to produce vitamins and essential amino-acids required for insect  
123 development and cuticle biosynthesis.

## 124 Conclusions

125 In addition to being an agricultural pest and a valuable endosymbiotic system, *S. oryzae* can  
126 be a remarkable model for studying TE evolution and regulation, along with the impact of  
127 TEs on eukaryotic genomes.

128

## 129 Keywords

130 Coleoptera, weevil, *Sitophilus oryzae*, genome, transposable elements, endosymbiosis,  
131 immunity, evolution

132

## 133 Background

134 Beetles account for approximately 25% of known animals, with an estimated number of 400  
135 000 described species [1–3]. Among them, Curculionidae (true weevils) is the largest animal  
136 family described, comprising about 70 000 species [1,4,5]. Despite being often associated  
137 with ecological invasion and ecosystem degradation, only three Curculionidae genomes are  
138 publicly available to date [6–8]. Among the cereal weevils, the rice weevil *Sitophilus oryzae*  
139 is one of the most important pests of crops of high agronomic and economic importance  
140 (wheat, maize, rice, sorghum and barley), causing extensive quantitative and qualitative  
141 losses in field, stored grains and grain products throughout the world [9–11]. Moreover, this  
142 insect pest is of increasing concern due to its ability to rapidly evolve resistance to  
143 insecticides such as phosphine, a fumigant used to protect stored grains from insect pests  
144 [12–14].

145 Like other holometabolous insects, the life cycle of *S. oryzae* can be divided into four stages:  
146 egg, larva, pupa and adult (Figure 1). Females drill a small hole in the grain, deposit a single  
147 egg and seal it with secretions from their ovipositor. Up to six eggs can be laid daily by each

148 female, totaling around 400 eggs over its entire lifespan [15]. Larvae develop and pupate  
149 within the grain kernel, metamorphose, and exit the grain as adults. The whole process  
150 takes on average 30 days [10]. Like many insects living on nutritionally poor diets, cereal  
151 weevils permanently associate with nutritional intracellular bacteria (endosymbionts) that  
152 supply them with nutrients that are not readily available in the grains, thereby increasing their  
153 fitness and invasive power. The endosymbiont of *S. oryzae*, the gamma-proteobacterium  
154 *Sodalis pierantonius* [16,17], is housed within specialized host cells, named bacteriocytes,  
155 that group together into an organ, the bacteriome [18]. Contrasting with most studied  
156 symbiotic insects, the association between *Sitophilus* spp. and *S. pierantonius* was  
157 established recently (less than 30 000 years ago), probably following the replacement of the  
158 ancestor endosymbiont, *Candidatus Nardonella*, in the Dryophthorinae subfamily [19,20]. As  
159 a result, contrary to long-lasting endosymbiotic associations, the genome of *S. pierantonius*  
160 is GC rich (56.06%), and its size is similar to that of free-living bacteria (4.5 Mbp) [16].  
161 Moreover, it encodes genes involved in bacterial infection, including Type Three Secretion  
162 Systems (TTSS), as well as genes encoding Microbial Associated Molecular Patterns  
163 (MAMPs) that trigger Pattern Recognition Receptors (PRR), and are usually absent or  
164 reduced in bacteria involved in long-lasting associations [16,21,22]. Nevertheless, many  
165 features indicate that the genome of *S. pierantonius* is in a process of degradation, as it  
166 contains many pseudogenes (43% of the predicted protein-coding sequences) and a large  
167 number of mobile elements (18% of the genome size) [16,23]. Finally, it is important to note  
168 that no other symbionts, with the exception of the familiar *Wolbachia* endosymbiont in some  
169 strains, have been described in *S. oryzae*.

170 In order to help unravel potential adaptive functions and features that could become the  
171 basis for identifying novel control strategies for weevils and other major insect pests, we  
172 have undertaken the sequencing, assembly and annotation of the genome of *S. oryzae*.  
173 Strikingly, the repeated fraction of *S. oryzae*'s genome (repeatome), composed mostly of  
174 transposable elements (TEs), is among the largest found to date in insects. TEs, the most  
175 versatile DNA units described to date, are sequences present in multiple copies and capable

176 of relocating or replicating within a genome. While most observed TE insertions evolve  
177 neutrally or are slightly deleterious, there are a number of documented cases where TEs  
178 may facilitate host adaptation (for reviews see [24–26]). For instance, gene families involved  
179 in xenobiotic detoxification are enriched in TEs in *Drosophila melanogaster* [27], *Plutella*  
180 *xylostella* [28], a major crop pest, and *Myzus persicae*, another phytophagous insect causing  
181 significant agronomic losses [29]. TEs have also been frequently associated with insecticide-  
182 resistance in *Drosophila* species [30–32]. In addition, population genetics studies suggested  
183 that more than 84 TE copies in *D. melanogaster* may play a positive role in fitness-related  
184 traits [33], including xenobiotic resistance [32] and immune response to Gram-negative  
185 bacteria [34].

186 In eukaryotes, TE content varies drastically and contributes significantly to the size and  
187 organization of the genome. From TE-rich genomes as maize (85% [35]), humans ( $\approx$ 45%  
188 [36]), and the recently sequenced lungfish ( $\approx$ 90% [37]) for instance, to TE-poor genomes, as  
189 *D. melanogaster* (12-15% [38]), or *Arabidopsis thaliana* ( $\approx$ 10% [39]), repeatomes thrive on a  
190 high level of diversity. These drastic variations are also observed within animal clades, such  
191 as insects, where the proportion of TE ranges from 2% in the Antarctic midge (*Belgica*  
192 *antarctica*) to 65% in the migratory locust (*Locusta migratoria*) [40–42] and up to 75% in  
193 morabine grasshoppers (*Vandiemenella viatica* species) [43]. In addition to the overall TE  
194 content, the number of different TE families (homogeneous groups of phylogenetically  
195 related TE sequences), their size (number of copies per family) and sequence diversity are  
196 also very high among insect species [44]. For instance, SINEs (Short INterspersed  
197 Elements) are almost absent from most insect genomes, but many lepidopterans harbor  
198 these elements [44]. In flies, Long Terminal Repeats retrotransposons (LTRs) are a staple of  
199 the *Drosophila* genus, but such TEs are nearly absent from other dipteran genomes (e.g.  
200 *Glossina brevipalpis* and *Megaselia scalaris*) [44]. Recent advances in sequencing have  
201 dramatically increased the level to which TEs can be studied across species and reveal that  
202 such variations can persist even within recently diverged groups, as observed within  
203 *Drosophila* species [45] or among *Heliconius* butterflies [46]. An increasing number of insect



204 genomes are reported with large repeatomes (e.g. *Aedes aegypti* and *Ae. albopictus* 40-  
205 50% [47,48], *L. migratoria* 60-65% [40,41], *Dendrolimus punctatus* 56% [49], *Vandiemena*  
206 *viatica* species 66-75% [43]).

207 Here we present the genome of *S. oryzae*, with a strong focus on the repeatome, its  
208 largest genomic compartment, spanning over  $\approx$ 74% of the assembly. *S. oryzae* represents a  
209 model system for stored grain pests, host-TE evolutionary biology, and the study of the  
210 molecular mechanisms acting at the early steps of symbiogenesis. Moreover, the features  
211 uncovered suggest that *S. oryzae* and its relatives have the potential to become a platform  
212 to study the interplay between TEs, host genomes and endosymbionts.

213






## 214 Results and discussion

### 215 Genome assembly and annotation

216 We have sequenced and assembled the genome of the rice weevil *S. oryzae* at a base  
217 coverage depth of 142X using a combination of short and long read strategies (see  
218 Methods). The karyotype of *S. oryzae* comprises 22 chromosomes [50], and the genome  
219 assembly consists of 2 025 scaffolds spanning 770 Mbp with a N50 of 2.86 Mbp,  
220 demonstrating a high contiguity compared to other Coleopteran genomes (Table 1). The  
221 assembly size is consistent with the genome size measured through flow cytometry (769  
222 Mbp in females and 768 Mbp in males [50]). We assessed the completeness of the genome  
223 assembly using BUSCO [51] (97.9% complete and 0.7% fragmented), and along with the  
224 aforementioned statistics, *S. oryzae* is the best assembled Curculionidae genome to date  
225 [7,52,53] (Table 1). The complete analysis of gene content and function can be found in  
226 Additional file 1.

227

228 Table 1. Assembly statistics of *S. oryzae*'s genome in comparison to Curculionidae genomes  
 229 and *T. castaneum*

Statistics	<i>Sitophilus oryzae</i> 	<i>Rhynchophorus ferrugineus</i> [8] 	<i>Hypothenemus hampei</i> [6] 	<i>Dendroctonus ponderosae</i> [7] 	<i>Tribolium castaneum</i> [53] 
Order, Family	Coleoptera, Curculionidae	Coleoptera, Curculionidae	Coleoptera, Curculionidae	Coleoptera, Curculionidae	Coleoptera, Tenebrionidae
No. chromosomes	2n=22 [50]	2n=22 [54]	2n=14 [55]	2n=24 [56]	2n=20 [57]
No. scaffolds	2,025	4,807	15,896	8,188	2,149
Total length (Mb)	770	782	151	253	166
Scaffold N50 (Kb)	2,861	64,117	39	629	4,456
GC%	32.9	30.5	27.8	38.4	35.2
Gap length (Mb)	12.6	40.6	20.9	51.0	13.5
Median coverage	142x	108x	100x	443x	-
BUSCO (% complete/partial)	98/99	92/94	97/98	96/97	99/100
No. protein-coding genes	15,057	25,567*	19,222*	13,021	12,862

230 \*All genes, no NCBI RefSeq annotation report available.

231

## 232 Annotation of the *Sitophilus oryzae* genome

233 Among the different pathways we were able to decipher in the genome of *S. oryzae*, we  
 234 present here highlights of the main annotation efforts, followed by a detailed analysis of the  
 235 TE content and impact on the host genome. A comprehensive analysis for each highlight is  
 236 presented as Supplemental Notes in Additional file 1.

## 237 Phylome and horizontal gene transfer

238 *Sitophilus oryzae* has a high gene expansion rate when compared to other beetles. Some of  
239 the families with the largest expansions include genes coding for proteins with DNA binding  
240 motifs, potentially regulating functions specific to this clade. Olfactory receptors,  
241 antimicrobial peptides (AMPs) and P450 cytochromes were expanded as well, probably in  
242 response to their ecological niche and lifestyle. Additionally, we noticed an expansion of  
243 plant cell wall degrading enzymes that originated from horizontal gene transfer (HGT) events  
244 from both bacteria and fungi. Given the intimate relationship between *S. oryzae* and its  
245 endosymbiont, including the permanent infection of the female germline, we searched for  
246 evidence for HGT in the weevil genome possibly coming from *S. pierantonius*. Contrary to  
247 the genome of the tsetse fly *Glossina*, where at least three HGT events from *Wolbachia*  
248 have been reported [58], we were unable to pinpoint any HGT event from either the ancient  
249 endosymbiont *Nardonella*, *Wolbachia*, or the recently acquired *S. pierantonius*. A detailed  
250 description is reported in Additional file 1: Supplemental Note 1 and Note 3 for digestive  
251 enzymes.

## 252 Global analysis of metabolic pathways

253 Using the CycADS [59] pipeline and Pathway Tools [60] we have generated BioCyc  
254 metabolism reconstruction databases for *S. oryzae* and its endosymbiont *S. pierantonius*.  
255 We compared *S. oryzae* metabolism to that of other arthropods available in the  
256 ArthropodaCyc collection and we explored the metabolic exchanges between weevils and  
257 their endosymbionts (see Additional file 1: Supplemental Note 2). The metabolic  
258 reconstruction reveals that, despite its large genome for an endosymbiotic bacterium, *S.*  
259 *pierantonius* relies on its host for several central compounds, including alanine and proline,  
260 but also isocitrate, Inosine MonoPhosphate (IMP) and Uridine MonoPhosphate (UMP), to  
261 produce essential molecules to weevils, including the essential amino acids tryptophan,  
262 phenylalanine, lysine and arginine, the vitamins pantothenate, riboflavin and dihydropterolate  
263 as a folate precursor, as well as nicotinamide adenine dinucleotide (NAD) (Additional file 1:

264 Supplemental Note 2). Among the amino acids listed above, phenylalanine, in particular is  
265 an essential precursor for the cuticle synthesis in emerging adults [61]. In addition, several  
266 studies have shown that *S. pierantonius* improves host fitness, including fertility,  
267 developmental time and flight capacity, in part by supplying the host with vitamins and  
268 improving its mitochondrial energy metabolism [62–64].

## 269 Development

270 The annotation of developmental genes uncovered a high level of conservation in  
271 comparison to the red flour beetle *Tribolium castaneum*, a model coleopteran. When  
272 compared to *D. melanogaster*, several key coordinate group genes are absent in *T.*  
273 *castaneum* and *S. oryzae*. Moreover, a number of genes with two homologs in the  
274 *Drosophila* genome are represented by a single ortholog in *T. castaneum* and *S. oryzae*. We  
275 also observed that homologs for several signaling pathway ligands could not be identified,  
276 which, given the presence of conserved receptors, is probably due to divergent primary  
277 sequence of the ligands. A detailed description is reported in Additional file 1: Supplemental  
278 Note 4.

## 279 Cuticle protein genes

280 Among the distinctive biological features of coleopterans is the ability to generate a hard and  
281 thick cuticle that protects them against dehydration and represents the first physical barrier  
282 from infections and topical insecticide penetration. The analysis of cuticle proteins (CPs)  
283 showed that *S. oryzae* has an average number of CPs, but with an enrichment of members  
284 of the CPAP1 family. While some members of this family are known to be involved in molting  
285 and maintaining the integrity of the cuticle in *T. castaneum*, most are still uncharacterized  
286 [65,66]. Thus, these proteins might be involved in the development of specific cuticular  
287 tissues in *S. oryzae* or other weevils. The total number of CPs did not follow the taxonomy of  
288 beetles, suggesting instead that it might be an adaptation to their diverse lifestyles. For  
289 details see Additional file 1: Supplemental Note 5.

## 290 Innate immune system

291 The analysis of immunity-related genes revealed that the genome of *S. oryzae* encodes the  
292 canonical genes involved in the three main antimicrobial pathways Toll, Imd and JAK-STAT,  
293 suggesting functional conservation of these pathways in cereal weevils. The conservation of  
294 the Imd pathway in the *S. oryzae* genome is of particular interest as its degradation in other  
295 symbiotic insects (*Acyrtosiphon pisum* [67], *B. tabaci* [68], or *Rhodnius prolixus* [69]) was  
296 initially correlated to their symbiotic status. The Imd pathway is not only present in *S. oryzae*,  
297 but it is also functional [70,71], and has evolved molecular features necessary for  
298 endosymbiont control [70] and host immune homeostasis [71]. Thus, not only is the Imd  
299 pathway conserved in cereal weevils, contrary to aphids and some other hemimetabolous  
300 insects, but it seems to have been evolutionary “rewired” toward additional functions in  
301 symbiotic homeostasis [70]. A detailed description can be seen in Additional file 1:  
302 Supplemental Note 6.

## 303 Detoxification and insecticide resistance

304 Fumigation using phosphine, hydrogen phosphide gas (PH<sub>3</sub>), is by far the most widely used  
305 treatment for the protection of stored grains against insect pests due to its ease of use, low  
306 cost, and universal acceptance as a residue-free treatment [72,73]. However, high-level  
307 resistance to this fumigant has been reported in *S. oryzae* from different countries [13,74–  
308 81]. Hence, we searched for genes associated with detoxification and resistance to  
309 insecticide and more generally to toxins, including plant allelochemicals. The *S. oryzae*  
310 repertoire of detoxification and insecticide resistance genes includes more than 300  
311 candidates, similar to what is seen in other coleopteran genomes. For more details see  
312 Additional file 1: Supplemental Note 7.

## 313 Odorant receptors

314 One promising pest management strategy relies on modifying insect behavior through the  
315 use of volatile organic compounds that act on odorant receptors (ORs) [82,83]. ORs play a

316 significant role in many crucial behaviors in insects by mediating host-seeking behavior,  
317 mating, oviposition, and predator avoidance [84]. Interfering with the behavior of pest insects  
318 and modulating their ability to find suitable hosts and mates has been shown to reduce  
319 population numbers, notably using plants that are capable of producing attractants and  
320 repellents [85,86]. *Sitophilus* spp. are known to use kairomones for host detection [87,88], as  
321 well as aggregation pheromones [89,90]. We annotated 100 candidate OR genes in *S.*  
322 *oryzae* (named SoryORs), including the gene encoding the co-receptor Orco. Of these  
323 genes, 46 were predicted to encode a full-length sequence. The global size of the SoryOR  
324 gene repertoire is in the range of what has been described in other species of the  
325 coleopteran suborder Polyphaga (between 46 in *Agrilus planipennis* and more than 300 in *T.*  
326 *castaneum*) and close to the number of OR genes annotated in the closely related species  
327 *Dendroctonus ponderosae* (85 genes, [91]) (Additional file 1: Supplemental Note 8).

328

## 329 Massive expansion of TE copies in the genome of *S. oryzae*

### 330 Detection and annotation of the repeatome

331 The repeatome represents the fraction of the genome categorized as repetitive. It  
332 encompasses TEs, satellites, tandem, and simple repeats. Eukaryotic TEs can be separated  
333 into two classes, depending on their replication mode [92]. DNA (Class II) based elements  
334 are able to directly move within a genome, and include terminal inverted repeat (TIR),  
335 Crypton, rolling-circle (RC/Helitron), and large composite elements (Maverick). Conversely,  
336 retrotransposons (Class I) have an RNA intermediate and replicate through RNA  
337 retrotranscription. Retrotransposons can be further divided into long terminal repeat (LTR),  
338 and non-LTR elements, including long and short interspersed nuclear repeat elements  
339 (LINEs and SINEs). Other retrotransposons include Penelope-like (PLEs) and DIRS-like  
340 elements. Each one of these TE orders can be further classified into specific superfamilies

341 (as for instance Copia or Gypsy LTR elements, and hAT or Tc1/Mariner TIR elements), that  
342 may encompass hundreds of TE families, each containing thousands of copies. The intrinsic  
343 diversity of TEs complicates their identification and annotation, especially in understudied  
344 species genera.

345 We used multiple state-of-the-art TE detection tools, including RepeatModeler2 and EDTA  
346 [93,94], to generate consensus sequences of the TE families in *S. oryzae*. After an initial  
347 discovery step, more than 10 000 likely redundant TE families were identified by the  
348 dedicated programs; we combined their results using multiple sequence alignments and  
349 clustering (see Methods and Additional file 1: Figure S1) to reduce this number to 3 399.  
350 Due to the evolutionary distance between *S. oryzae* and other known coleopterans, the  
351 consensus sequences obtained were further classified using a thorough combination of  
352 sequence homology and structure (see Methods).

353 The *S. oryzae* genome is among the most TE-rich insect genomes to date  
354 We uncovered 570 Mbp of repeat sequences, corresponding to  $\approx 74\%$  of the *S. oryzae*  
355 genome:  $\approx 2\%$  of satellite sequences, simple or low-complexity repeats, and  $\approx 72\%$  of other  
356 mobile elements, including TEs, (Figure 2A, Additional file 2). Given the limitation of the  
357 sequencing technologies, the proportion of satellites and TEs usually abundant in the  
358 heterochromatin is likely underestimated. We took advantage of a recent comparative  
359 analysis of TE content in 62 insect species [40] to contrast with the *S. oryzae* TE  
360 compartment. The *S. oryzae* genome ranks among those with the highest TE fraction  
361 observed in insects (Figure 2B and 2C). Within the largest insect order, Coleoptera, very  
362 little is known regarding TE distribution and evolution. *T. castaneum* harbors only 6% of TEs  
363 [53], *Hypothenemus hampei* contains 8.2% of TEs [6,95], while *Dichotomius schiffleri*  
364 harbors 21% [96]. The species closest to *S. oryzae*, *Rhynchophorus ferrugineus*, has a TE  
365 content of 45% [8]. Therefore, while TE content has been described to follow phylogenetic  
366 relationships in most insects [44,45] there is a large variation among the few Coleoptera  
367 species with available genomes. It is important to note that the pipeline we used to detect

368 and annotate TEs in *S. oryzae* differs from the method implemented by Petersen and  
369 colleagues [40], as we incorporated 31 manually curated TE references for *S. oryzae*, and  
370 specifically annotated DNA/TIR elements based on their sequence structure (see Methods),  
371 increasing the annotation sensitivity.

## 372 Class II (DNA) elements dominate *S. oryzae*'s genome

373 The most striking feature of the genome of *S. oryzae* is the high abundance of Class II  
374 (DNA) elements ( $\approx 32\%$  of the genome,  $\approx 43\%$  of the TE content) (Figure 2A), which is the  
375 highest observed among all 62 insect species included in this analysis [40–42]. The most  
376 DNA transposon-rich genomes include mosquito *Culex quinquefasciatus*, and *Ae. aegypti*,  
377 harboring 25% and 20% of DNA transposon content in their genome, amounting to 54% and  
378 36% of the total TE compartment, respectively [6]. The TE-rich grasshopper *L. migratoria*  
379 repeatome comprises only 14% of DNA transposons, while LINE retroelements (Class I)  
380 amount to 25%. Morabine grasshoppers, with up to 75% of TE content, show equivalent  
381 amounts of DNA, LINE and Helitrons [43]. Finally, among Coleoptera, a large diversity of  
382 repeatomes is observed (Figure 2C) with *A. planipennis*, *Leptinotarsa decemlineata* and  
383 *Onthophagus taurus* carrying an abundant LINE content, while *S. oryzae*, *T. castaneum* and  
384 *Anoplophora glabripennis* show larger DNA transposon content.

385 Among the Class II elements present in *S. oryzae*, the majority belongs to the TIR  
386 subclass but has not been assigned a known superfamily (Figure 2D), while Tc Mariner  
387 make up  $\approx 6\%$  of DNA elements. Among the consensus sequences we were able to  
388 assemble from 5'TIR to 3'TIR (highest confidence, see Methods), the length distribution  
389 shows a continuum starting at a couple of hundred bases to a maximum of  $\sim 5$  Kbp (see  
390 Figure 2E). We hypothesize that most of the smaller TIR families observed are miniature  
391 inverted repeat elements (MITEs). MITEs are non-autonomous elements, deriving from  
392 autonomous ClassII/TIR copies, comprising two TIRs flanking a unique, non-coding, region  
393 (sometimes absent) of variable length. While the TE detection pipeline used was able to  
394 detect and annotate most Class II/TIR elements based on transposase homologies, we also



395 specifically searched for non-autonomous TIR sequences, allowing the detection of putative  
396 MITEs that lack protein coding regions (Additional file 1: Figure S1). Among all Class II/TIR  
397 superfamilies, TIR length varies between tens of base pairs to  $\approx 1$  Kbp (Figure 2E). We  
398 identified short elements, composed mostly of their TIR sequences (Figure 2E), typical of  
399 MITEs. Interestingly, the unknown TIR families show an average size smaller than 1 Kbp,  
400 while TIRs with an annotated superfamily, show larger sizes (Additional file 1: Figure S3),  
401 suggesting that most unknown families could be indeed non-autonomous MITEs. MITE size  
402 ranges were previously described from around 100 bp to copies reaching more than 1 Kbp  
403 [97]. Finally, the distribution of the proportions of TIR relative to the consensus length  
404 appears superfamily-specific (Figure 2E and Additional file 1: Figure S3), and unknown  
405 families recapitulate these patterns. In conclusion, while most unknown TIR families seem to  
406 be composed of MITEs, we cannot exclude that our homology database is limited, likely  
407 missing some unknown protein domains. The most abundant TE family detected in the *S.*  
408 *oryzae* genome is indeed a MITE element (TE2641\_SO2\_FAM0704), with 10 486 genomic  
409 hits (or the equivalent of  $\approx 4$  117 copies based on the consensus size), corresponding to  
410 1.3% of the genome. Large fractions of MITEs were also reported in Class II-rich genomes,  
411 such as the aforementioned mosquitoes [48,98] and the invasive *Ae. albopictus* [47], but  
412 also in many plant species such as the rice *Oryza sativa* [99–101]. Among Class II elements,  
413 we have also detected Crypton (0.9% of the genome), RC/Helitrons (0.4% of the genome)  
414 and Mavericks (0.3% of the genome).

415 LINE elements are the second most abundant TE subclass, representing  $\approx 11\%$  of the *S.*  
416 *oryzae* genome, among which  $\approx 35\%$  are assigned to RTE elements and  $\approx 22\%$  to I elements  
417 (Figure 2D). No SINE families have been detected. LTRs are rather scarce, representing  
418 only  $\approx 3\%$  of the genome (Figure 2D), and the vast majority belongs to the Gypsy superfamily  
419 ( $\approx 30\%$ ). Another retrotransposon order detected are Penelope (PLEs), reaching nearly 2%  
420 of *S. oryzae*'s genome, and DIRS (Tyrosine recombinase retrotransposons, 0.14% of the  
421 genome).

422 Finally, around 22% of the genome is composed of repeats for which our pipeline could not  
423 assign a known TE class (Figure 2D). These unknown families highlight the wealth and  
424 diversity of TEs among insects and Coleopteran genomes in particular, and could represent  
425 an overlooked reservoir of genomic innovations.

426 TE copies make up most of non-coding sequences of *S. oryzae*'s genome

427 TE copies are interspersed around the *S. oryzae* genome. TEs are less frequently found  
428 close to gene transcription start sites (TSS), 5' and 3' untranslated regions (5' and 3' UTRs)  
429 and exons (Figure 3A), as expected. On the contrary, introns and intergenic sequences  
430 harbor the highest TE content (Figure 3A), amounting to around 50% of TE density, close to  
431 the general TE proportion in the genome (72%), suggesting that most non-coding DNA  
432 sequences in the *S. oryzae* genome are virtually made of TEs. To grasp the impact of TEs  
433 on intron size, we compared intron length in *S. oryzae* with two very well assembled  
434 genomes: *D. melanogaster* with a very compact and small genome, and the large, TE-rich  
435 human genome (Figure 3B). In *D. melanogaster*, introns are small and harbor few TEs, while  
436 in humans, introns are much larger potentially due to high TE accumulation [102]. *S. oryzae*  
437 intron sizes also seem to be due, at least partly, to TE accumulation. Interestingly, the *S.*  
438 *oryzae* genome presents a bimodal distribution, with a large proportion of small introns, as  
439 found in *D. melanogaster*, but also a noticeable amount of larger, TE packed and more  
440 human-like introns. This could suggest that specific regions of the genome could be more  
441 prone to TE elimination, and be associated with high rates of recombination and/or signature  
442 of purifying selection.

443 TE activity inferred by evolutionary history

444 Within reconstructed TE families, nucleotide substitution levels (Kimura 2 Parameters, K2P)  
445 between copies and their consensus sequences allowed estimation of their relative ages and  
446 identified potentially active ones (Figure 4A). Such "TE landscapes" are extremely helpful to  
447 pinpoint potential TE amplifications (modes in the distribution) and extinctions (valleys) within

448 the 0-30% K2P range (beyond, the increased divergence between copies affects negatively  
449 the sensitivity of the alignments, such that TE-derived sequences are no longer  
450 recognisable). The landscape analysis revealed a heterogeneous distribution of the TE copy  
451 divergence to their consensus within and between the main TE subclasses (Figure 4A). Most  
452 identified TE copies have a K2P divergence under 10, which is often observed in insects,  
453 and strikingly distinguishes itself from TE-rich mammalian genomes (RepeatMasker.org,  
454 [40]). While *S. oryzae*'s TE density and distribution evokes the architecture of mammalian  
455 genomes, this relatively younger TE landscape suggests higher deletion rates, and possibly  
456 a higher TE turnover rate, as observed in *Drosophila* [103,104]. LINEs and DNA transposons  
457 have the wider spectrum of divergence levels, suggesting an aggregation of distinct  
458 dynamics for the TE families present in *S. oryzae*. By contrast, the rare LTR copies identified  
459 appear to be the most homogeneous within families, with only a few substitutions between  
460 copies and their consensuses, suggesting a very recent amplification in this subclass.  
461 Finally, unknown TEs share a large part of their K2P distribution with TIR elements, though  
462 relatively less divergent from their consensus sequences as a whole. A breakdown of the  
463 K2P distributions at the superfamily level reveals specific evolutionary dynamics (Figure 4B).  
464 Diverse superfamilies, such as Tc-Mar and hAT (TIR) or RTE (LINE), show more uniform  
465 distributions, suggesting sustained activity of some of its members throughout *S. oryzae*'s  
466 genome evolution, though this could also indicate that these subfamilies could be subdivided  
467 further. As observed at the class level, all three identified LTR superfamilies (Pao, Gypsy  
468 and Copia) show families within the lowest K2P range.

469 TEs are transcriptionally active in somatic and germline tissues

470 The TE K2P landscape suggests that LTR elements as well as some LINE families and  
471 several Class II subclasses are among the youngest, and thus potentially active. In order to  
472 estimate the transcriptional activity of *S. oryzae*'s TE families, we have produced somatic  
473 (midgut) and germline (ovary) transcriptomic data. While germline tissues allow identification  
474 of potential TE families capable of producing vertically transmitted new copies, TE

475 derepression in somatic tissues represents the potential mutational burden due to TEs. The  
476 expression of TE families varied extensively within a class and the proportion of  
477 transcriptionally active/inactive TE families between classes was also distinct (Figure 5A). In  
478 total, 1 594 TE families were differentially expressed between ovary and midgut tissues  
479 (Figure 5B, Additional file 3); of which 329 have an absolute Log<sub>2</sub> fold change higher than 2  
480 (71 downregulated and 258 upregulated in midgut). In total, we detected 360 TE families  
481 downregulated in midgut when compared to ovaries: A much larger set of upregulated TE  
482 families was detected in midgut when compared to ovaries (1 236), illustrating the tighter  
483 regulation of TE copies in germline tissues. Moreover, the distribution of Log<sub>2</sub> fold changes  
484 were similar between TE subclasses but different for LTRs, which had a higher proportion of  
485 upregulated TE families in ovaries compared to other classes (Figure 5C. Kruskal and  
486 Wallis rank-sum test:  $H = 36.18$ ,  $P < 0.01$ ; LTR vs. LINE, Class II or Unknown: Dunn's test:  
487  $P\text{-adj} < 0.01$ ). In conclusion, the large TE compartment in *S. oryzae* shows abundantly  
488 expressed TE families, and tissue-specific expression patterns.

489 To estimate the TE transcriptional load imposed on *S. oryzae*, we computed the  
490 percentage of total RNA-seq poly-A enriched reads mapping to TE consensus sequences  
491 (Additional file 1: Figure S4). Around 5% of the midgut transcriptome corresponds to TE  
492 sequences. We compared such transcriptional burden to a TE-poor (*D. melanogaster*,  
493  $\approx 12\%$ ) and a TE-rich (*Ae. albopictus*  $\approx 50\%$ ) genome, using similar technology in equivalent  
494 tissues (adult midgut, see Methods). It is important to note that, despite being a TE-poor  
495 genome, *D. melanogaster* harbors many young LTR elements that have been recurrently  
496 shown to transpose [105]. We did not detect a direct correlation between genomic TE  
497 content and TE expression (Additional file 1: Figure S4). *S. oryzae* bears the highest  
498 proportion of RNAseq reads mapped against TE consensus sequences ( $\approx 5\%$ ), followed by  
499 *D. melanogaster* ( $\approx 1\%$ ) and *Ae. albopictus* ( $\approx 0.01\%$ ). Henceforth, not only is *S. oryzae* a TE-  
500 rich genome, but the transcriptional load from TEs is higher than in other TE-rich genomes  
501 (*Ae. albopictus*), and in genomes harboring young and active TE copies (*D. melanogaster*,  
502 [38,106]).

503 Finally, it is important to note that while transcriptional activation of TE copies may have an  
504 impact on the host genome, it does not indicate high transposition and therefore higher  
505 mutation rates. The high transcriptional load of *S. oryzae* compared to other species might  
506 stem from differences in TE regulation. In insects, TEs are mainly silenced by small RNAs  
507 and repressive chromatin marks [107]. More specifically, piwi-interacting RNAs (piRNAs) are  
508 able to target post-transcriptional repression of TEs, and guide chromatin silencing  
509 complexes to TE copies [107–109]. Therefore, we have annotated genes implicated in small  
510 RNA biogenesis and found that all three pathways (piRNAs but also microRNAs and small  
511 interfering RNAs biogenesis pathways) are complete (Additional file 1: Supplemental Note  
512 9). Genes involved in piRNA biosynthesis are expressed mainly in ovaries and testes, while  
513 somatic tissues (midgut) show smaller steady-state levels (Additional file 1: Supplemental  
514 Note 9), suggesting the piRNA pathway is potentially functional in *S. oryzae* ovaries, and  
515 could efficiently reduce transposition.

#### 516 TE content is variable among *Sitophilus* species

517 Cereal weevils are part of the Dryophthoridae family that includes more than 500 species.  
518 Very little is known about genome dynamics in this massive phylogenetic group. Because of  
519 the unusual high TE copy number and landscape observed in *S. oryzae*, we analyzed three  
520 other closely related species namely *Sitophilus zeamais*, *Sitophilus granarius* and *Sitophilus*  
521 *linearis*. We produced low coverage sequencing and estimated the TE content from raw  
522 reads using our annotated *S. oryzae* TE library with dnaPipeTE [47]. Remarkably, among  
523 *Sitophilus* species, repeat content is variable (Figure 6A), with *S. linearis* harboring the  
524 smaller repeat load ( $\approx 54\%$ ) compared to *S. oryzae* ( $\approx 80\%$ ), *S. zeamais* ( $\approx 79\%$ ), and *S.*  
525 *granarius* ( $\approx 65\%$ ). Most importantly, Class II (DNA) elements of *S. oryzae* are nearly absent  
526 from *S. linearis*, and no recent burst of LTR elements is observed, contrary to the other  
527 *Sitophilus* species, suggesting alternative TE evolutionary histories (Figure 6B). It is  
528 important to note that our analysis is biased towards *S. oryzae*, as the library used to  
529 annotate the TEs in the other *Sitophilus* species stems from automatic and manual

530 annotation of the *S. oryzae* genome. Finally, the relatively higher dnaPipeTE estimations of  
531 the LTR content in *S. oryzae* compared to the assembled genome supports the hypothesis  
532 that LTR elements have seen a recent burst of transposition, as young elements tend to  
533 collapse in genome assemblies and eventually diminish their estimated copy number.

534 Overall, the comparison of TE content in closely related species highlights the influence of  
535 phylogenetic inertia, but reveals a possible TE turnover in the *S. linearis* lineage. In addition  
536 to the regulation mechanisms that strongly contribute to TE amount and variation, TE  
537 accumulation is conditioned by the drift/selection balance in populations. Indeed, effective  
538 population size has been suggested to be a major variable influencing TE content, as small,  
539 inbred or expanding populations suffer drift, allowing detrimental insertions to stay in the  
540 gene pool and thus favor TE fixation [110]. Such hypotheses should be addressed in the  
541 future, especially on recently sequenced TE-rich but rather small (<1 Gbp) genomes such as  
542 *S. oryzae*.

#### 543 Endosymbionts impact TE transcriptional regulation

544 The four *Sitophilus* species studied have different ecologies. *S. oryzae* and *S. zeamais* infest  
545 field cereals and silos, while *S. granarius* is mainly observed in cereal-containing silos. *S.*  
546 *linearis*, however, lives in a richer environment, *i.e.* tamarind seeds. In association with their  
547 diets, the interaction of *Sitophilus* species with endosymbiotic bacteria differs: the cereal  
548 weevils (*S. oryzae*, *S. zeamais* and *S. granarius*) harbor the intracellular gram-negative  
549 bacteria *S. pierantonius*, albeit at very different loads. While *S. oryzae* and *S. zeamais* show  
550 high bacterial load, *S. granarius* has a smaller bacterial population [61]. In contrast, *S.*  
551 *linearis* has no nutritional endosymbionts, in correlation with its richer diet. We wondered  
552 whether the presence of intracellular bacteria impacts TE regulation, and took advantage of  
553 artificially obtained aposymbiotic *S. oryzae* animals to search for TE families differentially  
554 expressed in symbiotic versus aposymbiotic ovaries. There were 50 TE families upregulated  
555 in symbiotic ovaries compared to artificially obtained aposymbiotic ones, while 15 families  
556 were downregulated (Figure 7 and Additional file 3). Only three families presented an

557 absolute Log<sub>2</sub> fold change higher than 2: one LINE and two LTR/Gypsy elements. The three  
558 of them were upregulated both in symbiotic *versus* aposymbiotic ovaries, and in ovaries  
559 *versus* midgut (Additional file 5), suggesting that such elements have tissue specificity, and  
560 their expression is modulated by the presence of intracellular bacteria. Such TE families  
561 would be ideal candidates to further study the crosstalk between host genes, intracellular  
562 bacteria and TE transcriptional regulation.

563

## 564 Conclusion

565 The success of obtaining a TE-rich genome assembly complete enough to understand  
566 genome architecture and regulatory networks relies on the use of multiple sequencing  
567 platforms [111]. Here, we describe the first assembly of the repeat-rich (74%) *S. oryzae*  
568 genome, based on a combination of long and short read sequencing, and a new assembly  
569 method, WENGAN [112]. While this first assembly reaches quality standards similar to other  
570 coleopteran species (Table 1), it is important to stress that new sequencing methods have  
571 emerged in order to improve genome assemblies, including linked-reads and optical  
572 mapping [111].

573 We uncovered around 74% of repeated sequences in the *S. oryzae* genome, mostly TE  
574 families. While the TE landscape is marked by a wealth of Class II elements, especially non-  
575 autonomous MITE elements, 22% of the genome is composed of unknown repeats. Large  
576 duplicated gene families can be present in such a category, but it is tempting to speculate  
577 that the majority is composed of novel Class II elements. Indeed, Unknown and TIR  
578 elements share the same K2P landscapes, and many Class II elements have only been  
579 detected through an inverted repeat search for TIRs, and not proteins, excluding therefore  
580 TE copies old enough that TIRs are too divergent to be recognized. Moreover, we have  
581 shown that many TE families in *S. oryzae* are present in the transcriptome, suggesting that  
582 several families can be transcriptionally active. How such TE families are able to escape

583 host silencing remains unknown. It seems obvious today that insect models such as *D.*  
584 *melanogaster* only represent a small window on the complex biology and evolution of TEs,  
585 and the sequencing and annotation of species with high TE content -- while challenging  
586 [113] -- is key to understanding how genomes, their size, their structure and their function  
587 evolve. In conclusion, *S. oryzae* constitutes an excellent model to understand TE dynamics  
588 and regulation and the impact on genome function.

589 *Sitophilus* species not only differ in their TE landscape, but also in their ecology and as a  
590 consequence, their association with intracellular bacteria. Comparison of TE content within  
591 the *Sitophilus* genus shows variable TE amount and diversity. In addition, intracellular  
592 bacterium impacts transcription of specific TE families in ovaries. The molecular  
593 mechanisms behind the co-evolution between an insect, its endosymbiotic bacterium and  
594 TEs remains unexplored. The impact of intracellular bacteria on host genomes is poorly  
595 studied, and the *Sitophilus* genus offers a simpler experimental setting, with a single  
596 intracellular bacterium present within specific host cells [19,62], and a well established  
597 knowledge of host-bacteria interaction [61,70,71,114,115].

598

## 599 Methods

### 600 DNA extraction and high-throughput sequencing

601 Individuals of both sexes of *S. oryzae* were reared on wheat grains at 27.5 °C with 70%  
602 relative humidity. The aposymbiotic strain was obtained by treating the symbiotic strain  
603 during one month at 35 °C and 90% relative humidity as previously described [116]. This  
604 strain is viable, fertile and was raised in the same conditions as the symbiotic strain. The  
605 aposymbiotic status was confirmed by PCR and histology. Male and female adults of *S.*  
606 *oryzae* were used for DNA extraction. Only the gonads were used to minimize DNA



607 contamination from its diet, which could be still present in the gut. The reproductive organs  
608 were obtained from aposymbiotic adults and a DNA extraction protocol specific for *Sitophilus*  
609 weevils was performed. DNA extractions were performed using a STE buffer (100 mM NaCl,  
610 1 mM Na<sub>2</sub>EDTA pH 8, 10 mM Tris HCl pH 8). Tissues were homogenized in STE buffer, then  
611 treated successively by SDS 10%, proteinase K and RNase. Briefly, genomic DNA was  
612 purified by two successive extractions with phenol:chloroform:isoamyl alcohol (25/24/1)  
613 followed by extraction with 1 vol of chloroform:isoamyl alcohol (24/1). Genomic DNA was  
614 then precipitated by 0.7 vol isopropanol. After washing the pellet with 70% ethanol, genomic  
615 DNA was recovered in TE (1 mM EDTA, 10 mM Tris HCl pH8) buffer. Using this protocol, we  
616 obtained six different DNA samples: four from males and two from females. Each sample  
617 corresponds to the genomic DNA from 20 individuals. Five additional DNA samples were  
618 obtained using a high molecular weight DNA extraction protocol consisting of a single  
619 phenol:chloroform:isoamyl alcohol (25/24/1) extraction step from the genomic DNA of 100  
620 males. The DNA concentration in each of these samples was quantified using a NanoDrop  
621 spectrophotometer (ThermoFisher Scientific, Waltham, MA, USA).

622 Sequencing was performed using a combination of Illumina, PacBio and Nanopore  
623 technologies (Additional file 4). For each sex, two Illumina libraries were generated: one  
624 paired-end library with an average fragment size of 500 bp and one mate pair library with an  
625 average fragment size of 5 Kbp. The libraries were sequenced using an Illumina HiSeq 2000  
626 platform with the V3 chemistry and a read size of 101 bp; the paired-end (PE) libraries were  
627 sequenced at the "Génomique & Microgénomique" service from ProfileXpert (Lyon, France)  
628 while the mate paired (MP) were sequenced at Macrogen (Seoul, South Korea). Two male  
629 samples were used to build (i) an Illumina library with an average fragment size of 200 bp  
630 which was sequenced on a HiSeq 2500 instrument using the V4 chemistry and a read size  
631 of 125 bp, and (ii) a PacBio library sequenced on seven SMRT cells using the P6-C4  
632 chemistry. These two libraries were sequenced at KeyGene (Wageningen, The  
633 Netherlands). Finally, five male samples were used to build Nanopore libraries with the SQK-

634 LSK109 kit and without DNA fragmentation step. The libraries were independently  
635 sequenced on five MinION R9.4 flow cells. These libraries were built and sequenced at the  
636 sequencing platform of the IGFL (Institut de Génomique Fonctionnelle de Lyon, Ecole  
637 Normale Supérieure de Lyon, France). Statistics and accession numbers from all the  
638 sequencing runs are listed in the Additional file 3.

639

## 640 Genome assembly and annotation

641 First, the Illumina reads were error-corrected using BFC release 181 [117]. The PacBio and  
642 Nanopore reads were error-corrected using LORDEC v0.9 [118] with the error-corrected  
643 Illumina overlapping PE reads, a k-mer size of 19 and solidity threshold of 3. Overlapping  
644 reads were then merged using FLASH2 v2.2 [119]. Based on the merged Illumina reads, a  
645 first short-read assembly was produced using a modified version of MINIA v3.2.1 [120] with  
646 a k-mer length of 211. A hybrid assembly was then performed using WENGAN v0.1 [112] on  
647 the MINIA short-read assembly and the raw Nanopore reads. The resulting assembly was  
648 polished using two rounds of PILON v1.23 [121] using the error-corrected Illumina  
649 overlapping PE reads and the --diploid option. A first scaffolding was then performed with  
650 two rounds of FAST-SG v06/2019 [122] and SCAFFMATCH v0.9 [123] with the error-  
651 corrected Illumina MP, Illumina PE, PacBio and Nanopore libraries. The LR\_GAPCLOSER  
652 algorithm v06/2019 [124] was used for the gap-filling step using the error-corrected PacBio  
653 and Nanopore libraries. An additional scaffolding step was performed using RASCAF v1.0.2  
654 [125] with the available RNA-seq libraries from the Sequence Read Archive (SRX1034967-  
655 SRX1034972 and SRX3721133-SRX3721138). The resulting scaffolds were then gap-filled  
656 using a new round of LR\_GAPCLOSER as previously described followed by two rounds of  
657 SEALER v2.1.5 [126] using the error-corrected Illumina overlapping PE reads and k-mer  
658 sizes of 64 and 96. Two rounds of PILON, as previously described, were performed to  
659 produce the final assembly. Quality of the assembly was assessed by computing several

660 metrics using i) QUAST v5.0.2 [127] with a minimal contig size of 100 bp and the --large and  
661 -k options, ii) BUSCO v4.0.5 [51] using the Insecta ODB10 database and the -geno option,  
662 and iii) KMC v3.0.0 [128] to evaluate the percentage of shared 100-mers between the  
663 assembly and the merged Illumina reads.

664 Three contaminant scaffolds corresponding to the mitochondrial genome and an artefact  
665 were removed from the assembly prior to the annotation step. The 'NCBI *Sitophilus oryzae*  
666 Annotation Release 100' was produced using the NCBI Eukaryotic Genome Annotation  
667 Pipeline v8.2.

668

### 669 Low-coverage genome sequencing of other *Sitophilus* species

670 Twenty pairs of ovaries were dissected from *S. oryzae*, *S. zeamais*, *S. granarius* and *S.*  
671 *linearis* females. Ovaries were homogenized in 100 mM NaCl, 1 mM EDTA pH 8, 10 mM  
672 Tris-HCl pH 8 using a small piston. Proteinase K digestion followed in the presence of SDS  
673 for 2 h at 55 °C with shaking and for 1 h at 37 °C with RNase A. A typical phenol chloroform  
674 extraction was then performed and genomic DNA was isopropanol precipitated. Eight whole  
675 genome sequencing libraries with a median insert size of 550 bp were constructed using the  
676 Illumina TruSeq DNA PCR-free sample preparation kit (Illumina, San Diego, CA, USA),  
677 according to manufacturer's protocols. Briefly, 2 µg of each gDNA were sheared using a  
678 Covaris M220 Focused-ultrasonicator (Covaris, Inc. Woburn, MA, USA), end-repaired, A-  
679 tailed, and adapter ligated. Library quality control was performed using the 2100 Bioanalyzer  
680 System with the Agilent High Sensitivity DNA Kit (Agilent Technologies, Santa Clara, CA,  
681 USA). The libraries were individually quantified via qPCR using a KAPA Library  
682 Quantification Kits (Kapa Biosystems, Wilmington, MA, USA) for Illumina platforms, then  
683 they were pooled together in equimolar quantities and sequenced in a MiSeq sequencing  
684 system. 2x300 paired-end reads were obtained using a MiSeq Reagent Kits (600-cycles).

685

## 686 TE library construction

687 In order to annotate the *S. oryzae* repeatome, we collected and combined cutting-edge  
688 bioinformatic tools to (i) create and (ii) classify a non-redundant library of repeated elements  
689 (Additional file 1: Figure S1). First, we separately ran RepeatModeler2 [93] and EDTA [94]  
690 on the assembled genome. Together, these programs include most of the recent and long-  
691 trusted tools used to detect generic repeats, but also include specific modules, such as for  
692 LTR and TIR elements. Preliminary analyses of the *S. oryzae* genome with RepeatModeler1  
693 [129] and dnaPipeTE [47] suggested a rather large fraction of Class II DNA elements with  
694 terminal inverted repeats (TIRs). Thus, MITE-Tracker [130] was incorporated in our pipeline  
695 and ran independently on the genome assembly using 1- and 2-Kbp size cutoffs to detect  
696 Class II elements harboring TIRs with high sensitivity. Following this initial step, 15 510  
697 consensus sequences obtained from RM2, EDTA and the two runs of MITE-tracker were  
698 successively clustered using MAFFT [131], Mothur [132], and Refiner [129] to reduce  
699 redundancy in the repeat library to a total of 2 754 consensus sequences (Additional file 1:  
700 Figure S1A, <https://github.com/clemgoub/So2>). Then, we inspected the quality of the raw  
701 library by calculating the genomic coverage of each consensus. We ran the library against  
702 the genome using RepeatMasker (52) and implemented a simple algorithm “TE-trimmer.sh”  
703 to trim or split a consensus sequence wherever the genomic support drops below 5% of the  
704 average consensus coverage (Additional file 1: Figure S1A,  
705 <https://github.com/clemgoub/So2>). To mitigate any redundancy generated by the splitting,  
706 the newly trimmed library was clustered before being re-quantified using RepeatMasker  
707 [129]. At this step, we removed any consensus under 200 bp and represented by less than  
708 the equivalent of two full-length copies (in total bp). In addition, TAREAN [133] was used to  
709 detect and quantify candidate satellite repeats. We obtained an *ab-initio* repeat library of 3  
710 950 consensus sequences automatically generated (Additional file 1: Figure S1A).

711 To refine and improve the quality of the TE consensus sequences, we then turned it over  
712 to DFAM [134] who processed the *ab initio* library. First, any sequences mostly composed of

713 tandem repeats were removed using a custom script to remove any sequences that were  
714 greater than 80% masked and/or had a sequence less than 100 bp. To generate seed  
715 alignments for each consensus, the consensus sequences were used as a search library for  
716 RepeatMasker to collect interspersed repeats. Seed alignments in the form of stockholm  
717 files were generated using the RepeatMasker output. To extend potentially truncated  
718 elements, the instances in the stockholm file for each model were extended into neighboring  
719 flanking sequences until the alignment was below a threshold equivalent to ~3 sequences in  
720 agreement. More specifically, all sequences are extended using full dynamic programming  
721 matrices using an improved affine gap penalty (default: -28 open, -6 extension) and a full  
722 substitution matrix (default: 20 percent divergence, 43% GC background). The termination of  
723 extension occurs when the improvement by adding a further column to the multiple  
724 alignment does not exceed 27 (with default scoring system). This is equivalent to a net gain  
725 of ~3 sequences in agreement. Following extension, the new consensus were collected and  
726 consensus sequences greater than 80% similar for 80% of their length were considered  
727 duplicates and only one consensus was kept.

728 Upon completion, we used RepeatMasker to quantify the improved library. We selected  
729 the top 50 elements (by abundance in the genome) represented in each of the “LTR”,  
730 “LINE”, “Class II” and “Unknown” classes for manual inspection (these categories represent  
731 the 4 most abundant classes of repeats in the *S. oryzae* genome). While most consensus  
732 sequences were correctly extended and annotated (200) we noticed some cases of over-  
733 extension with LTR (consensus doubled in size) and flagged others with non-supported  
734 fragments for further trimming (Additional file 2 | tab 1). Once our quality check completed  
735 and the sequences curated, we removed fragments with 100% identity against a previously  
736 established consensus (Additional file 2 | tab 2). The final TE library contains 3 399  
737 sequences to classify.

738 The classification of the final repeat library was done in successive rounds combining  
739 homology and structure methods (Additional file 1: Figure S1B). Before the final TE library  
740 was completed, we manually curated and annotated the sequences of 31 transposable

741 elements and satellites among the most represented in *S. oryzae*. These high-confidence  
742 references were added to the default libraries used by the following programs and Repeatbase  
743 v.2017 [135]. We searched for nucleotide homology using RepeatMasker (V.4.1.1 [129]) with  
744 -s “-slow” search settings. Best hits were chosen based on the highest score at the  
745 superfamily level allowing non-overlapping hits of related families to contribute to the same  
746 hit. In addition we used blastx [136] to query each consensus against a curated collection of  
747 TE proteins (available with RepeatMasker), as well as those identified in the 31 manual  
748 consensus sequences. We kept the best protein hit based on the blastx score. Based on the  
749 200 consensus sequences manually inspected (see above), we set a hit length / consensus  
750 size threshold of 0.08 (RepeatMasker) and 0.03 (blastx) to keep a hit. In our hands, these  
751 thresholds were conservative to automate the classification. As an alternate homology-  
752 based method, we also ran RepeatClassifier (RepeatModeler2). Finally, because DNA  
753 elements are often represented by non-autonomous copies (unidentifiable or absent  
754 transposase) we further used einverted to flag terminal inverted repeats located less than  
755 100 bp of the ends of each sequence. The complete library of 3 399 consensus sequences  
756 was first annotated at the subclass level (see DFAM taxonomy:  
757 <https://dfam.org/classification/tree>) if two out of RepeatMasker, RepeatClassifier and blastx  
758 annotations agreed. Further, the same rule was applied for the superfamilies if possible. At  
759 this stage, consensus sequences without annotation by homology but with TIRs as flagged  
760 by einverted, were classified as TIR and all other sequences classified as Unknown. We  
761 further divided the subclass “DNA” into “MAV” (Mavericks), “RC” (Rolling circle/Helitron),  
762 “CRY” (Crypton) and “TIR” (terminal inverted repeats). Finally, the classifications  
763 automatically given as “Unknown” to 16/274 manually inspected consensus sequences were  
764 replaced to match the manually reported classification.  
765

## 766 Estimation of the repeat content

767 The total repeat content of the *S. oryzae* genome was analyzed using RepeatMasker  
768 (v.4.1.1) and our classified library of 3 399 consensus sequences and the following  
769 parameters: -s -gccalc -no\_is  
770 -cutoff 200. The subsequent alignments were parsed with the script “parseRM.pl” [137]  
771 <https://github.com/4ureliek/Parsing-RepeatMasker-Outputs>) to remove hits overlap and  
772 statistically analyzed with R version 4.0.2.

773

## 774 Genomic distribution of TE copies

775 The distribution of TE copies across the *S. oryzae* genome was assessed using two different  
776 approaches over six different genomic regions namely TSS  $\pm$  3 Kbp, 5' UTRs, exons,  
777 introns, 3' UTRs and intergenic regions. Briefly, the coverage of all TE copies was computed  
778 over a sliding window of 100 bp across the whole genome sequence using the  
779 makewindows and coverage tools from the bedtools package [138] and the  
780 bedGraphToBigWig UCSC gtfToGenePred tool. Then the different genomic regions were  
781 retrieved from the *S. oryzae* annotation file (GFF format) using the gencode\_regions script  
782 ([https://github.com/saketkc/gencode\\_regions](https://github.com/saketkc/gencode_regions)) and the UCSC gtfToGenePred tool  
783 (<https://github.com/ENCODE-DCC/kentUtils>). A matrix containing the TE coverage per  
784 genomic region was generated using the computeMatrix tool from deepTools [139] and used  
785 to generate metaplots using the plotProfile tool.

786

## 787 TE landscapes

788 The relative age of the different TE families identified in the genome assembly was drawn  
789 performing a “TE-landscape” analysis on the RepeatMasker outputs. Briefly, the different  
790 copies of one TE family identified by RepeatMasker are compared to their consensus

791 sequence and the divergence (Kimura substitution level, CpG adjusted, see RepeatMasker  
792 webpage: <http://repeatmasker.org/webrepeatmaskerhelp.html>) is calculated. The TE  
793 landscape consists of the distribution of these divergence levels. In the end, the relative age  
794 of a TE family can be seen as its distribution within the landscape graph: “older” TE families  
795 tend to have wider and flatter distribution spreading to the right (higher substitution levels)  
796 than the “recent” TE families, which are found on the left of the graph and have a narrower  
797 distribution. TE landscapes were drawn from the RepeatMasker output parsed with the  
798 options -l of “parseRM.pl”. We report here the TE landscape at the level of the TE subclass  
799 (LINE, LTR, TIR, CRY, MAV, DIRS, PLE, RC and Unknown).

800

## 801 dnaPipeTE comparative analysis in *Sitophilus* species

802 To compare the TE content of *S. oryzae* to four related species of *Sitophilus* (*S. granarius*,  
803 *S. zeamais*, *S. linearis*) we used dnaPipeTE v.1.3 [47]. dnaPipeTE allows unbiased  
804 estimation and comparison of the total repeat content across different species by assembling  
805 and quantifying TE from unassembled reads instead of a linear genome assembly. Reads  
806 for *Sitophilus* species were produced as described above. Using our new classified library (3  
807 390 consensus) as TE database in dnaPipeTE, we were further able to identify the  
808 phylogenetic depth of the repeat identified in *S. oryzae*.

809

## 810 RNA sequencing and TE expression analysis

811 Adapter sequences and low quality reads were filtered out with Trimmomatic (v0.36) [140]  
812 and clean reads were aligned to the *S. oryzae* genome with STAR aligner (v2.5.4b, [141])  
813 and featureCounts from subread package [142] to obtain gene counts. We also used the  
814 STAR aligner to map the clean reads against all TE copies extracted from the genome with  
815 the following options: --outFilterMultimapNmax 100 --winAnchorMultimapNmax 100 --



816 outMultimapperOrder Random --outSAMmultNmax 1. The mapped bam files were used as  
817 input to TEtools software [143] to determine TE family expression. Genes and TE family  
818 counts were used as input for DESeq2 package [144] to determine differential TE  
819 expression between Ovary vs Gut tissues as well as Ovaries from symbiotic and  
820 aposymbiotic weevils. Differentially expressed TEs were defined whenever the adjusted p-  
821 value was smaller than 0.05 and Log2 fold change was higher than 1 or smaller than -1. We  
822 used the aforementioned STAR alignment parameters to map transcriptomic sequencing  
823 reads from midgut of *S. oryzae* (Accession: SRX1034971, and SRX1034972), *D.*  
824 *melanogaster* (Accession: SRX029389, and SRX045361), and *Ae. albopictus* (Accession:  
825 SRX1512976, SRX1898481, SRX1898483, SRX1898487, SRX3939061, and SRX3939054)  
826 against the TE consensus sequences for each species.

827

## 828 Abbreviations

829 AMPs: AntiMicrobial Peptides

830 CPs: cuticle proteins

831 HGT: horizontal gene transfer

832 IMP: Inosine MonoPhosphate

833 K2P: Kimura 2 Parameters

834 LINE: long INterspersed Element

835 LTR: Long Terminal Repeat

836 MAMPs: Microbial Associated Molecular Patterns

837 MITEs: miniature inverted repeat elements

- 838 ORs: odorant receptors
- 839 PRR: Pattern Recognition Receptors
- 840 PLE: penelope-like
- 841 RC: rolling circle
- 842 SINE: Short INterspersed Element
- 843 TIR: terminal inverted repeat
- 844 TSS: transcription start sites
- 845 TEs: transposable elements
- 846 TTSS: Type Three Secretion Systems
- 847 UTR: untranslated regions
- 848 UMP: Uridine MonoPhosphate

849

## 850 **Declarations**

851 Ethics approval and consent to participate

852 Not applicable

853

854 Consent for publication

855 Not applicable

856

## 857 Availability of data and materials

858 This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank under the  
859 accession PPTJ00000000. The version described in this paper is version PPTJ02000000.  
860 The assembly can be visualised, along with gene models and supporting data, on a  
861 dedicated genome browser ([https://bipaa.genouest.org/sp/sitophilus\\_oryzae/](https://bipaa.genouest.org/sp/sitophilus_oryzae/)). Raw reads  
862 from low coverage genome sequencing of *S. zeamais*, *S. granarius* and *S. linearis* have  
863 been deposited at NCBI Sequence Read Archive (SRA) under the BioProject accessions  
864 PRJNA647530, PRJNA647520 and PRJNA647347 respectively. TE annotation (GFF) and  
865 consensus sequences can be found at <https://dx.doi.org/10.5281/zenodo.4570415>. Bisulfite-  
866 seq reads have been deposited at NCBI SRA, under the BioProject accession  
867 PRJNA681724.

868

## 869 Competing interests

870 The authors declare that they have no competing interests.

871

## 872 Funding

873 Funding for this project was provided by the Fondation de l'Institut National des Sciences  
874 Appliquées-Lyon (INSA-Lyon), the Institut National de Recherche pour l'Agriculture,  
875 l'Alimentation et l'Environnement (INRAE), the French ANR-10-BLAN-1701  
876 (ImmunSymbArt), the French ANR-13-BSV7-0016-01 (IMetSym), the French ANR-  
877 17\_CE20\_0031\_01 (GREEN), the French ANR-17-CE20-0015 (UNLEASH), the Santé des  
878 Plantes et Environnement (SPE) department at INRAE , the IDEX-Lyon PALSE IMPULSION  
879 initiative and a grant from la Région Rhône-Alpes. The project was also funded by European  
880 Regional Development Fund (ERDF) and Ministerio de Ciencia, Innovación y Universidades

881 (Spain) PGC2018-099344-B-I00 to AL, and PID2019-105969GB-I00 to AM and Conselleria  
882 d'Educació, Generalitat Valenciana (Spain), grant number PROMETEO/2018/133 to AM.  
883 CV-C was a recipient of a fellowship from the Ministerio de Economía y Competitividad  
884 (Spain).

885

## 886 Authors' contributions

887 AH and AL conceived the original sequencing project and were joined by NP, RR, CG, CV-  
888 C, AM and CV who participated in the coordination of the project. AVa, CV-M, ED, JM, FM  
889 and AVi reared the inbred lines and AVa extracted genomic DNA and RNA that was used for  
890 library construction and sequencing. BG and SH produced and sequenced the Nanopore  
891 libraries. CG, AVa, MB, NB, CV, AG and ATRV produced and sequenced the low-coverage  
892 Illumina libraries. NP, CV-C, ADG and M-FS performed the genome assembly and  
893 automated gene prediction. CV-C, MM-H and TG analyzed and wrote the *phylome and*  
894 *horizontal gene transfer* note. PB-P, GF, SC, HC and FC analyzed and wrote the *global*  
895 *analysis of metabolic pathways* note. NP analyzed and wrote the *digestive enzymes* and the  
896 *detoxification and insecticide resistance* notes. PC analyzed and wrote the *development*  
897 note. CV-C analyzed and wrote the *cuticle protein genes* note. CV-M, CV-C, NP, JM, LB,  
898 AB, WZ, FM, AVi and AZ-R analyzed and wrote the *innate immune system* note. NM, CM,  
899 ASB and EJ-J analyzed and wrote the *odorant receptors* note. TC, CB, AVa and RR  
900 produced the data for the *epigenetic pathways* note. TC, CB, AVa, GR, CV-C, CV and RR,  
901 analyzed and wrote the *epigenetic pathways* note. MGF, CG, ED, RR, SB, GF, NM, CV-M  
902 and NP produced the figures. CG, RR, JMS, JR, RH and AFAS annotated and analyzed the  
903 TE content while MGF analyzed the TE RNAseq data. NP, CV-C, CG, CV, RR, AL and AH  
904 wrote the manuscript. All authors read and approved the final manuscript.

905

## 906 Acknowledgments

907 The authors acknowledge supercomputing resources made available by the Rhône-Alpes  
908 Bioinformatics center (PRABI-AMSB, <http://www.prabi.fr>) to perform the NGS data analyses.  
909 The authors would like to thank Stéphanie Robin, Fabrice Legeai and Anthony Bretaudeau  
910 at the INRAE Bioinformatics Platform for Agro-ecosystems Arthropods (BIPAA)  
911 (<https://bipaa.genouest.org>) for the *S. oryzae* genome integration. The authors would like to  
912 thank Fabienne Barbet, Séverine Croze and Nicolas Nazaret from profileXpert for Illumina  
913 sequencing of *Sitophilus oryzae* libraries. The authors thank the network REaCTION and its  
914 coordinators (N. Ponts, G. Le Trionnaire, I. Fudal, M. Jubault), funded by INRAE-SPE, for  
915 organizing a Bisulfite-seq workshop that allowed the production of the data presented here.  
916 We also thank J.-Y. Rasplus for collecting *Sitophilus linearis* individuals from Niger. The  
917 authors would like to thank Cédric Feschotte at Cornell University for the support, insights  
918 and the occasional use of bioinformatic resources.

## 919 Figure legends

920 Figure 1. *Sitophilus oryzae* overview. A. Life cycle of cereal weevil *Sitophilus oryzae*. The  
921 embryo develops into a larva and pupa, and metamorphoses into a young adult, exiting the  
922 grain around 3 days after metamorphosis completion. The developmental times indicated  
923 are from a rearing condition at 27 °C and 70% relative humidity. B. Photos of adult *S.*  
924 *oryzae*. Lower panel shows an adult exiting the grain.

925

926 Figure 2. A. Proportion of repeat content in *S. oryzae*'s genome. The majority of repeats  
927 detected in *S. oryzae* are represented by Class II (TIR) elements, LINEs (Class I), and  
928 unclassified repeats (unknown). NR: non repetitive. B. Variation of genome size and TE  
929 content in 62 insect species from [40] and *S. oryzae*. Coleopteran species are depicted in  
930 dark blue, and *S. oryzae* in light blue. *S. oryzae* is clearly a TE-rich genome. C. TE  
931 proportion across 11 insect species, including six coleoptera. In agreement with the data  
932 used for comparison [40], PLEs are included in the LINE superfamilies, DIRS in LTRs, and  
933 RC, CRY, MAV and TIR in the DNA superfamilies. NR: non repetitive. *S. oryzae* harbors the  
934 largest TE content among Coleopterans and most insect species studied to date. Within  
935 Coleoptera, there is a large variation in TE content and type, with *A. planipennis*, *L.*  
936 *decemlineata* and *O. taurus* carrying an abundant LINE content, while *S. oryzae*, *T.*  
937 *castaneum* and *A. glabripennis* show larger DNA content. Cladogram based on [145]. D.  
938 Classification of the 570 Mbs of TEs present in the *S. oryzae* genome. Most TIR families  
939 detected were not classified into known superfamilies. RTE LINE and Gypsy LTR elements  
940 are the most abundant superfamilies among retrotransposons. Around 22% of repeats in *S.*  
941 *oryzae*'s genome were not classified by our pipeline, and remain unknown (grey). E.  
942 Distribution of TIR length sequences (right) detected by einverted, and the internal region  
943 present between both TIRs (left) for complete consensus of TIR superfamilies (color) and  
944 unknown TIR families (grey).

945

946 Figure 3. TE distribution in *S. oryzae*'s genome. A. Density of TE copies within gene regions.  
947 TE copies are the least abundant within TSSs, 5' and 3' UTRs and exons, while introns and  
948 intergenic regions are riddled with TEs. TSS: transcription start site, UTR: untranslated  
949 regions. B. Relationship between intron length and TE per intron in *D. melanogaster* (red), *H.*  
950 *sapiens* (blue) and *S. oryzae* (yellow). *S. oryzae* shares characteristics of both *Drosophila*  
951 with short and TE poor introns and Humans with a significant number of large, TE-packed  
952 introns.

953

954 Figure 4. A. TE divergence landscape. Distribution of the divergence (Kimura two  
955 parameters, K2P) between TE copies and their consensus, aggregated by TE class reported  
956 in percent of the genome. The less divergent superfamilies are distributed to the left and  
957 suggest recent activity. Strikingly, most of the TE copies have less than 10% divergence to  
958 their consensus, with a large number of copies under 5% (dotted line). The distribution of the  
959 "unknown" class overlaps with the leftmost mode of the TIR distribution, suggesting that  
960 many more TIR families are yet to be described in *S. oryzae*. Strikingly, LTR elements are  
961 the least diverged altogether with the mode of the distribution on the 0-1% divergence bin. B.  
962 Mean K2P distributions within TE superfamilies. Left panel depicts Class II families, and all  
963 Class I (retrotransposons) and unknown families are on the right panel. LTR superfamilies  
964 harbor some of the least divergent TE families, suggesting that this class may host some of  
965 the youngest TE.

966

967 Figure 5. TE family expression in midguts and ovaries from *S. oryzae*. A. Log10 normalized  
968 counts in midguts and ovaries triplicates. Normalized counts show different proportions of  
969 transcriptionally active TE families in different TE classes. B. Log10 of base mean average  
970 expression of TE families in ovaries and midguts from three biological replicates. Depicted in  
971 color only TE families which had differential expression between ovary and gut tissues  
972 ( $\text{padj} < 0.05$ ,  $|\log_2\text{FC}| > 2$ ). Most TE families are upregulated in midguts compared to ovaries.  
973 C. Distribution of all significant ( $\text{padj} < 0.05$ ). Log2FC depicts specifically deregulated TE  
974 classes in each tissue. LTR elements are predominantly upregulated in ovaries.

975

976 Figure 6. TE landscape across *Sitophilus* species. A. Proportion of TE per species estimated  
977 from short reads with dnaPipeTE and a custom TE library including Repbase (release 2017)  
978 and annotated TE consensus discovered in *S. oryzae*, *S. zeamais* and *S.*  
979 *granarius* harbor similar TE content, while *S. granarius* presents a smaller TE load, and *S.*  
980 *linearis* harbors the smallest TE content and the higher proportion of unknown repeats. The  
981 proportion of unknown repeats only found by dnaPipeTE (black) increases from *S. oryzae* to  
982 *S. linearis* with the phylogenetic distance. B. Distribution of divergence values between raw  
983 reads and repeats contig assembled with dnaPipeTE (blastn) across four *Sitophilus* species.  
984 *S. oryzae* appears to share its TE landscape with *S. zeamais* and *S. granarius*, but the three  
985 species display a distinct repeatome than *S. linearis*, in spite of their phylogenetic proximity.  
986 SO2: *S. oryzae*'s TE library produced in this analysis, DPTE: DNAPipeTE TE annotation  
987 (repeats only found by dnaPipeTE).

988

989 Figure 7. Differentially expressed TE families between symbiotic and aposymbiotic *S. oryzae*  
990 ovaries. Log10 of base mean average expression of TE families in symbiotic vs  
991 aposymbiotic ovaries from two biological replicates. Depicted in color only TE families which  
992 had differential expression between both ovary types ( $\text{padj} < 0.05$ ,  $|\log_2\text{FC}| > 2$ ). Two LTR  
993 elements and one LINE element are upregulated ( $\log_2\text{FC} > 2$ ) in symbiotic ovaries.



994

## 995 Additional files

996 Additional file 1: Supplementary notes, supplementary figures, and small tables. (PDF)

997 Additional file 2: Transposable elements annotation tables. (XLSX)

998 Additional file 3: STAR and TTools mapping statistics. (XLSX)

999 Additional file 4: Summary of sequencing libraries produced for *S. oryzae*. (XLSX)

1000 Additional file 5: Large supporting tables and datasets. (XLSX)

1001

## 1002 References

- 1003 1. Hunt T, Bergsten J, Levkanicova Z, Papadopoulou A, John OS, Wild R, et al. A  
1004 comprehensive phylogeny of beetles reveals the evolutionary origins of a superradiation.  
1005 Science. 2007;318:1913–6.
- 1006 2. Stork NE, McBroom J, Gely C, Hamilton AJ. New approaches narrow global species  
1007 estimates for beetles, insects, and terrestrial arthropods. Proc Natl Acad Sci U S A.  
1008 2015;112:7519–23.
- 1009 3. Hammond P. Species Inventory. In: Groombridge B, editor. Global biodiversity: Status of  
1010 the Earth's living resources. 1992. Chapman and Hall, London. p. 17–39.
- 1011 4. McKenna DD, Sequeira AS, Marvaldi AE, Farrell BD. Temporal lags and overlap in the  
1012 diversification of weevils and flowering plants. Proc Natl Acad Sci U S A. 2009;106:7083–8.
- 1013 5. Oberprieler RG, Marvaldi AE, Anderson RS. Weevils, weevils, weevils everywhere\*.  
1014 Zootaxa. 2007;1668:491–520.
- 1015 6. Vega FE, Brown SM, Chen H, Shen E, Nair MB, Ceja-Navarro JA, et al. Draft genome of  
1016 the most devastating insect pest of coffee worldwide: the coffee berry borer, *Hypothenemus*  
1017 *hampei*. Sci Rep. 2015;5:12525.

- 1018 7. Keeling CI, Yuen MM, Liao NY, Roderick Docking T, Chan SK, Taylor GA, et al. Draft  
1019 genome of the mountain pine beetle, *Dendroctonus ponderosae* Hopkins, a major forest  
1020 pest. *Genome Biol.* 2013;14:R27.
- 1021 8. Hazzouri KM, Sudalaimuthasari N, Kundu B, Nelson D, Al-Deeb MA, Le Mansour A, et  
1022 al. The genome of pest *Rhynchophorus ferrugineus* reveals gene families important at the  
1023 plant-beetle interface. *Commun Biol.* 2020;3:1–14.
- 1024 9. Zunjare R, Hossain F, Muthusamy V, Jha SK, Kumar P, Sekhar JC, et al. Genetic  
1025 variability among exotic and indigenous maize inbreds for resistance to stored grain weevil  
1026 (*Sitophilus oryzae* L.) infestation. *Cogent Food Agric.* 2016;2:1137156.
- 1027 10. Longstaff BC. Biology of the grain pest species of the genus *Sitophilus* (Coleoptera:  
1028 Curculionidae): a critical review. *Prot Ecol.* 1981;3:83–130.
- 1029 11. Grenier A-M, Mbaiguinam M, Delobel B. Genetical analysis of the ability of the rice  
1030 weevil *Sitophilus oryzae* (Coleoptera, Curculionidae) to breed on split peas. *Heredity.*  
1031 1997;79:15–23.
- 1032 12. Champ BR, Dyte CE. FAO global survey of pesticide susceptibility of stored grain pests.  
1033 FAO Plant Protec Bull. 1977;25(2):49-67.
- 1034 13. Nguyen TT, Collins PJ, Ebert PR. Inheritance and characterization of strong resistance  
1035 to phosphine in *Sitophilus oryzae* (L.). *PLoS One.* 2015;10:e0124335.
- 1036 14. Mills KA. Phosphine resistance: Where to now? In: Donahaye, EJ, Navarro, S and  
1037 Leesch JG, editors. *Proceeding International Conference on Controlled Atmosphere and*  
1038 *Fumigation in Stored Products*; 2000 Oct 29-Nov 3; Fresno, USA. 2000:583–91.
- 1039 15. Campbell JF. Fitness Consequences of Multiple Mating on Female *Sitophilus oryzae* L.  
1040 (Coleoptera: Curculionidae). *Environ Entomol.* 2005;34:833–43.
- 1041 16. Oakeson KF, Gil R, Clayton AL, Dunn DM, von Niederhausern AC, Hamil C, et al.  
1042 Genome degeneration and adaptation in a nascent stage of symbiosis. *Genome Biol Evol.*  
1043 2014;6:76–93.
- 1044 17. Heddi A, Charles H, Khatchadourian C, Bonnot G, Nardon P. Molecular characterization  
1045 of the principal symbiotic bacteria of the weevil *Sitophilus oryzae*: a peculiar G + C content of

- 1046 an endocytobiotic DNA. *J Mol Evol.* 1998;47:52–61.
- 1047 18. Heddi A, Charles H, Khatchadourian C. Intracellular bacterial symbiosis in the genus  
1048 *Sitophilus*: the 'biological individual' concept revisited. *Res Microbiol.* 2001;152:431–7.
- 1049 19. Lefèvre C, Charles H, Vallier A, Delobel B, Farrell B, Heddi A. Endosymbiont  
1050 phylogenesis in the Dryophthoridae weevils: evidence for bacterial replacement. *Mol Biol*  
1051 *Evol.* 2004;21:965–73.
- 1052 20. Clayton AL, Oakeson KF, Gutin M, Pontes A, Dunn DM, Niederhausern AC von, et al. A  
1053 Novel human-infection-derived bacterium provides insights into the evolutionary origins of  
1054 mutualistic insect–bacterial symbioses. *PLoS Genet.* 2012;8:e1002990.
- 1055 21. Akman L, Yamashita A, Watanabe H, Oshima K, Shiba T, Hattori M, et al. Genome  
1056 sequence of the endocellular obligate symbiont of tsetse flies, *Wigglesworthia glossinidia*.  
1057 *Nat Genet.* 2002;32:402–7.
- 1058 22. Shigenobu S, Watanabe H, Hattori M, Sakaki Y, Ishikawa H. Genome sequence of the  
1059 endocellular bacterial symbiont of aphids *Buchnera sp.* APS. *Nature.* 2000;407:81–6.
- 1060 23. Gil R, Belda E, Gosalbes MJ, Delaye L, Vallier A, Vincent-Monégat C, et al. Massive  
1061 presence of insertion sequences in the genome of SOPE, the primary endosymbiont of the  
1062 rice weevil *Sitophilus oryzae*. *Int Microbiol Off J Span Soc Microbiol.* 2008;11:41–8.
- 1063 24. Rebollo R, Romanish MT, Mager DL. Transposable elements: an abundant and natural  
1064 source of regulatory sequences for host genes. *Annu Rev Genet.* 2012;46:21–42.
- 1065 25. Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, et al. Ten  
1066 things you should know about transposable elements. *Genome Biol.* 2018;19:199.
- 1067 26. Chuong EB, Elde NC, Feschotte C. Regulatory activities of transposable elements: from  
1068 conflicts to benefits. *Nat Rev Genet.* 2017;18:71–86.
- 1069 27. Chen S, Li X. Transposable elements are enriched within or in close proximity to  
1070 xenobiotic-metabolizing cytochrome P450 genes. *BMC Evol Biol.* 2007;7:46.
- 1071 28. You M, Yue Z, He W, Yang X, Yang G, Xie M, et al. A heterozygous moth genome  
1072 provides insights into herbivory and detoxification. *Nat Genet.* 2013;45:220–5.
- 1073 29. Singh KS, Troczka BJ, Duarte A, Balabanidou V, Trissi N, Paladino LZC, et al. The

- 1074 genetic architecture of a host shift: An adaptive walk protected an aphid and its  
1075 endosymbiont from plant chemical defenses. *Sci Adv.* 2020;6:eaba1070.
- 1076 30. Carareto CMA, Hernandez EH, Vieira C. Genomic regions harboring insecticide  
1077 resistance-associated Cyp genes are enriched by transposable element fragments carrying  
1078 putative transcription factor binding sites in two sibling *Drosophila* species. *Gene.*  
1079 2014;537:93–9.
- 1080 31. Rostant WG, Wedell N, Hosken DJ. Chapter 2 - Transposable elements and insecticide  
1081 resistance. In: Goodwin SF, Friedmann T, Dunlap JC, editors. *Adv Genet.* Academic Press;  
1082 2012. p. 169–201.
- 1083 32. Mateo L, Ullastres A, González J. A transposable element insertion confers xenobiotic  
1084 resistance in *Drosophila*. *PLoS Genet.* 2014;10:e1004560.
- 1085 33. Rech GE, Bogaerts-Márquez M, Barrón MG, Merenciano M, Villanueva-Cañas JL,  
1086 Horváth V, et al. Stress response, behavior, and development are shaped by transposable  
1087 element-induced mutations in *Drosophila*. *PLoS Genet.* 2019;15:e1007900.
- 1088 34. Ullastres A, Merenciano M, González J. Natural transposable element insertions drive  
1089 expression changes in genes underlying *Drosophila* immune response. *bioRxiv.*  
1090 2019;655225.
- 1091 35. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, et al. The B73 maize  
1092 genome: complexity, diversity, and dynamics. *Science.* 2009;326:1112–5.
- 1093 36. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial  
1094 sequencing and analysis of the human genome. *Nature.* 2001;409:860–921.
- 1095 37. Meyer A, Schloissnig S, Franchini P, Du K, Woltering JM, Irisarri I, et al. Giant lungfish  
1096 genome elucidates the conquest of land by vertebrates. *Nature.* 2021;1–6.
- 1097 38. Adams MD, Celniker SE, Holt RA, Evans CA, Gocayne JD, Amanatides PG, et al. The  
1098 genome sequence of *Drosophila melanogaster*. *Science.* 2000;287:2185–95.
- 1099 39. The Arabidopsis Genome Initiative. Analysis of the genome sequence of the flowering  
1100 plant *Arabidopsis thaliana*. *Nature.* 2000;408:796–815.
- 1101 40. Petersen M, Armisén D, Gibbs RA, Hering L, Khila A, Mayer G, et al. Diversity and

- 1102 evolution of the transposable element repertoire in arthropods with particular reference to  
1103 insects. *BMC Evol Biol.* 2019;19:11.
- 1104 41. Wang X, Fang X, Yang P, Jiang X, Jiang F, Zhao D, et al. The locust genome provides  
1105 insight into swarm formation and long-distance flight. *Nat Commun.* 2014;5:2957.
- 1106 42. Kelley JL, Peyton JT, Fiston-Lavier A-S, Teets NM, Yee M-C, Johnston JS, et al.  
1107 Compact genome of the Antarctic midge is likely an adaptation to an extreme environment.  
1108 *Nat Commun.* 2014;5:4611.
- 1109 43. Palacios-Gimenez OM, Koelman J, Palmada-Flores M, Bradford TM, Jones KK, Cooper  
1110 SJB, et al. Comparative analysis of morabine grasshopper genomes reveals highly abundant  
1111 transposable elements and rapidly proliferating satellite DNA repeats. *BMC Biol.*  
1112 2020;18:199.
- 1113 44. Gilbert C, Peccoud J, Cordaux R. Transposable elements and the evolution of insects.  
1114 *Annu Rev Entomol.* 2021;66:355-372.
- 1115 45. Sessegolo C, Burllet N, Haudry A. Strong phylogenetic inertia on genome size and  
1116 transposable element content among 26 species of flies. *Biol Lett.* 2016;12:20160407.
- 1117 46. Ray DA, Grimshaw JR, Halsey MK, Korstian JM, Osmanski AB, Sullivan KAM, et al.  
1118 Simultaneous TE Analysis of 19 Heliconiine butterflies yields novel insights into rapid TE-  
1119 based genome diversification and multiple SINE births and deaths. *Genome Biol Evol.*  
1120 2019;11:2162–77.
- 1121 47. Goubert C, Modolo L, Vieira C, Valiente-Moro C, Mavingui P, Boulesteix M. *De novo*  
1122 assembly and annotation of the Asian tiger mosquito (*Aedes albopictus*) repeatome with  
1123 dnaPipeTE from raw genomic reads and comparative analysis with the Yellow fever  
1124 mosquito (*Aedes aegypti*). *Genome Biol Evol.* 2015;7:1192–205.
- 1125 48. Nene V, Wortman JR, Lawson D, Haas B, Kodira C, Tu Z (Jake), et al. Genome  
1126 sequence of *Aedes aegypti*, a major arbovirus vector. *Science.* 2007;316:1718–23.
- 1127 49. Zhang S, Shen S, Peng J, Zhou X, Kong X, Ren P, et al. Chromosome-level genome  
1128 assembly of an important pine defoliator, *Dendrolimus punctatus* (Lepidoptera;  
1129 Lasiocampidae). *Mol Ecol Resour.* 2020;20:1023–37.

- 1130 50. Silva AA, Braga LS, Corrêa AS, Holmes VR, Johnston JS, Oppert B, et al. Comparative  
1131 cytogenetics and derived phylogenetic relationship among *Sitophilus* grain weevils  
1132 (Coleoptera, Curculionidae, Dryophthorinae). *Comp Cytogenet.* 2018;12:223–45.
- 1133 51. Seppey M, Manni M, Zdobnov EM. BUSCO: Assessing genome assembly and  
1134 annotation completeness. In: Kollmar M, editor. *Gene Prediction. Methods Mol Biol.*  
1135 2019;1962. p. 227–45.
- 1136 52. McKenna DD, Scully ED, Pauchet Y, Hoover K, Kirsch R, Geib SM, et al. Genome of the  
1137 Asian longhorned beetle (*Anoplophora glabripennis*), a globally significant invasive species,  
1138 reveals key functional and evolutionary innovations at the beetle-plant interface. *Genome*  
1139 *Biol.* 2016;17:227.
- 1140 53. Tribolium Genome Sequencing Consortium, Richards S, Gibbs RA, Weinstock GM,  
1141 Brown SJ, Denell R, et al. The genome of the model beetle and pest *Tribolium castaneum*.  
1142 *Nature.* 2008;452:949–55.
- 1143 54. Al-Qahtani AH, Al-Khalifa MS, Al-Saleh AA. Karyotype, meiosis and sperm formation in  
1144 the red palm weevil *Rhynchophorus ferrugineus*. *Cytologia.* 2014;79:235–42.
- 1145 55. Brun LO, Stuart J, Gaudichon V, Aronstein K, French-Constant RH. Functional  
1146 haplodiploidy: a mechanism for the spread of insecticide resistance in an important  
1147 international insect pest. *Proc Natl Acad Sci U S A.* 1995;92:9861–5.
- 1148 56. Lanier GN, Wood DL. Controlled mating, karyology, morphology, and sex-ratio in the  
1149 *Dendroctonus ponderosae* complex. *Ann Entomol Soc Am.* 1968;61:517–26.
- 1150 57. Stuart JJ, Mocelin G. Cytogenetics of chromosome rearrangements in *Tribolium*  
1151 *castaneum*. *Genome.* 1995;38(4):673-80.
- 1152 58. Initiative IGG. Genome sequence of the Tsetse fly (*Glossina morsitans*): Vector of  
1153 African trypanosomiasis. *Science.* 2014;344:380–6.
- 1154 59. Vellozo AF, Véron AS, Baa-Puyoulet P, Huerta-Cepas J, Cottret L, Febvay G, et al.  
1155 CycADS: an annotation database system to ease the development and update of BioCyc  
1156 databases. *Database.* 2011;2011:bar008.
- 1157 60. Karp PD, Midford PE, Billington R, Kothari A, Krummenacker M, Latendresse M, et al.

- 1158 Pathway Tools version 23.0 update: software for pathway/genome informatics and systems  
1159 biology. *Brief Bioinform.* 2019;
- 1160 61. Vigneron A, Masson F, Vallier A, Balmand S, Rey M, Vincent-Monégat C, et al. Insects  
1161 recycle endosymbionts when the benefit is over. *Curr Biol.* 2014;24:2267–73.
- 1162 62. Heddi A, Grenier A-M, Khatchadourian C, Charles H, Nardon P. Four intracellular  
1163 genomes direct weevil biology: Nuclear, mitochondrial, principal endosymbiont, and  
1164 *Wolbachia*. *Proc Natl Acad Sci U S A.* 1999;96:6814–9.
- 1165 63. Grenier AM, Nardon C, Nardon P. The role of symbiotes in flight activity of *Sitophilus*  
1166 weevils. *Entomol Exp Appl.* 1994;70:201–8.
- 1167 64. Rio RVM, Lefevre C, Heddi A, Aksoy S. Comparative genomics of insect-symbiotic  
1168 bacteria: influence of host environment on microbial genome composition. *Appl Environ*  
1169 *Microbiol.* 2003;69:6825–32.
- 1170 65. Jasarapura S, Arakane Y, Osman G, Kramer KJ, Beeman RW, Muthukrishnan S. Genes  
1171 encoding proteins with peritrophin A-type chitin-binding domains in *Tribolium castaneum* are  
1172 grouped into three distinct families based on phylogeny, expression and function. *Insect*  
1173 *Biochem Mol Biol.* 2010;40:214–27.
- 1174 66. Jasarapura S, Specht CA, Kramer KJ, Beeman RW, Muthukrishnan S. Gene families of  
1175 cuticular proteins analogous to peritrophins (CPAPs) in *Tribolium castaneum* have diverse  
1176 functions. *PLoS One.* 2012;7:e49844.
- 1177 67. Gerardo NM, Altincicek B, Anselme C, Atamian H, Barribeau SM, de Vos M, et al.  
1178 Immunity and other defenses in pea aphids, *Acyrtosiphon pisum*. *Genome Biol.*  
1179 2010;11:R21.
- 1180 68. Zhang C-R, Zhang S, Xia J, Li F-F, Xia W-Q, Liu S-S, et al. The immune strategy and  
1181 stress response of the mediterranean species of the *Bemisia tabaci* complex to an orally  
1182 delivered bacterial pathogen. *PLoS ONE.* 2014;9:e94477.
- 1183 69. Salcedo-Porras N, Guarneri A, Oliveira PL, Lowenberger C. *Rhodnius prolixus*:  
1184 Identification of missing components of the IMD immune signaling pathway and functional  
1185 characterization of its role in eliminating bacteria. *PLoS ONE.* 2019;14:e0214794.

- 1186 70. Maire J, Vincent-Monégat C, Masson F, Zaidman-Rémy A, Heddi A. An IMD-like  
1187 pathway mediates both endosymbiont control and host immunity in the cereal weevil  
1188 *Sitophilus* spp. *Microbiome*. 2018;6:6.
- 1189 71. Maire J, Vincent-Monégat C, Balmand S, Vallier A, Hervé M, Masson F, et al. Weevil  
1190 *pgrp-Ib* prevents endosymbiont TCT dissemination and chronic host systemic immune  
1191 activation. *Proc Natl Acad Sci U S A*. 2019;116:5623–32.
- 1192 72. Chaudhry MQ. Phosphine resistance. *Pestic Outlook*. 2000;11:88–91.
- 1193 73. Chaudhry MQ. A review of the mechanisms involved in the action of phosphine as an  
1194 insecticide and phosphine resistance in stored-product insects. *Pestic Sci*. 1997;49:213–28.
- 1195 74. Athié I, Gomes RAR, Bolonhezi S, Valentini SRT, De Castro MFPM. Effects of carbon  
1196 dioxide and phosphine mixtures on resistant populations of stored-grain insects. *J Stored*  
1197 *Prod Res*. 1998;34:27–32.
- 1198 75. Rajendran S. Phosphine resistance in stored grain insect pests in India. *Proc 7th Int*  
1199 *Work Conf Stored-Prod Prot*. 1998. p. 14–19.
- 1200 76. Zeng L. Development and countermeasures of phosphine resistance in stored grain  
1201 insects in Guangdong, China, 642–647. *Proc Seventh Int Work Conf Stored-Prod Prot Eds J*  
1202 *Zuxun Quan Yongsheng T Xianchang G Lianghua* 14–19 Oct 1998 Beijing China Sichuan  
1203 *Publ House Sci Technol Chengdu China*. 1999.
- 1204 77. Benhalima H, Chaudhry MQ, Mills KA, Price NR. Phosphine resistance in stored-product  
1205 insects collected from various grain storage facilities in Morocco. *J Stored Prod Res*.  
1206 2004;40:241–9.
- 1207 78. Pimentel MAG, Faroni LRD, Silva FH da, Batista MD, Guedes RNC. Spread of  
1208 phosphine resistance among brazilian populations of three species of stored product insects.  
1209 *Neotrop Entomol*. 2010;39:101–7.
- 1210 79. Nguyen TT, Collins PJ, Duong TM, Schlipalius DI, Ebert PR. Genetic conservation of  
1211 phosphine resistance in the rice weevil *Sitophilus oryzae* (L.). *J Hered*. 2016;107:228–37.
- 1212 80. Holloway JC, Falk MG, Emery RN, Collins PJ, Nayak MK. Resistance to phosphine in  
1213 *Sitophilus oryzae* in Australia: A national analysis of trends and frequencies over time and



- 1214 geographical spread. J Stored Prod Res. 2016;69:129–37.
- 1215 81. Agrafioti P, Athanassiou CG, Nayak MK. Detection of phosphine resistance in major  
1216 stored-product insects in Greece and evaluation of a field resistance test kit. J Stored Prod  
1217 Res. 2019;82:40–7.
- 1218 82. Carey AF, Carlson JR. Insect olfaction from model systems to disease control. Proc Natl  
1219 Acad Sci U S A. 2011;108:12987–95.
- 1220 83. Andersson MN, Newcomb RD. Pest control compounds targeting insect  
1221 chemoreceptors: Another silent spring? Front Ecol Evol. 2017;5:5.
- 1222 84. Leal WS. Odorant reception in insects: roles of receptors, binding proteins, and  
1223 degrading enzymes. Annu Rev Entomol. 2013;58:373–91.
- 1224 85. Hassanali A, Herren H, Khan Z, Pickett J, Woodcock C. Integrated pest management:  
1225 The push-pull approach for controlling insect pests and weeds of cereals, and its potential  
1226 for other agricultural systems including animal husbandry. Philos Trans R Soc Lond B Biol  
1227 Sci. 2008;363:611–21.
- 1228 86. Hatano E, Saveer AM, Borrero-Echeverry F, Strauch M, Zakir A, Bengtsson M, et al. A  
1229 herbivore-induced plant volatile interferes with host plant and mate location in moths through  
1230 suppression of olfactory signalling pathways. BMC Biol. 2015;13:75.
- 1231 87. Ukeh DA, Woodcock CM, Pickett JA, Birkett MA. Identification of host kairomones from  
1232 maize, *Zea mays*, for the maize weevil, *Sitophilus zeamais*. J Chem Ecol. 2012;38:1402–9.
- 1233 88. Germinara GS, De Cristofaro A, Rotundo G. Behavioral responses of adult *Sitophilus*  
1234 *granarius* to individual cereal volatiles. J Chem Ecol. 2008;34:523–9.
- 1235 89. Phillips JK, Walgenbach CA, Klein JA, Burkholder WE, Schmuff NR, Fales HM. (R (\*),S  
1236 (\*))-5-hydroxy-4-methyl-3-heptanone male-produced aggregation pheromone of *Sitophilus*  
1237 *oryzae* (L.) and *S. zeamais* motsch. J Chem Ecol. 1985;11:1263–74.
- 1238 90. Schmuff NR, Phillips JK, Burkholder WE, Fales HM, Chen C-W, Roller PP, et al. The  
1239 chemical identification of the rice weevil and maize weevil aggregation pheromone.  
1240 Tetrahedron Lett. 1984;25:1533–4.
- 1241 91. Mitchell RF, Schneider TM, Schwartz AM, Andersson MN, McKenna DD. The diversity

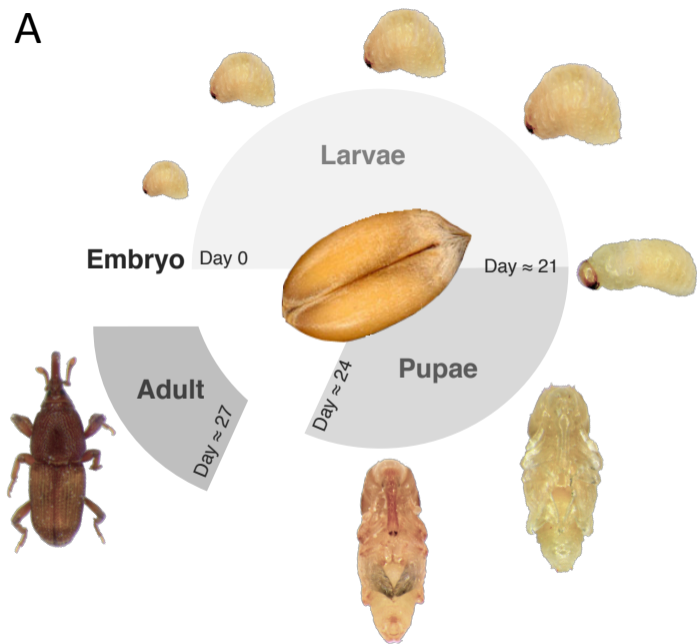
- 1242 and evolution of odorant receptors in beetles (Coleoptera). *Insect Mol Biol.* 2020;29:77–91.
- 1243 92. Makałowski W., Gotea V., Pande A., Makałowska I. Transposable elements:  
1244 Classification, identification, and their use as a tool for comparative genomics. In: Anisimova  
1245 M, editor. *Evolutionary Genomics. Methods Mol Biol*, 2019;1910.
- 1246 93. Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, et al. RepeatModeler2  
1247 for automated genomic discovery of transposable element families. *Proc Natl Acad Sci U S*  
1248 *A.* 2020;117:9451–7.
- 1249 94. Ou S, Su W, Liao Y, Chougule K, Agda JRA, Hellinga AJ, et al. Benchmarking  
1250 transposable element annotation methods for creation of a streamlined, comprehensive  
1251 pipeline. *Genome Biol.* 2019;20:275.
- 1252 95. Hernandez-Hernandez EM, Fernández-Medina RD, Navarro-Escalante L, Nuñez J,  
1253 Benavides-Machado P, Carareto CMA. Genome-wide analysis of transposable elements in  
1254 the coffee berry borer *Hypothenemus hampei* (Coleoptera: Curculionidae): description of  
1255 novel families. *Mol Genet Genomics.* 2017;292:565–83.
- 1256 96. Ic A, Es M, Rc M, Gl W. Diverse mobilome of *Dichotomius (Luederwaldtinia) schiffleri*  
1257 (Coleoptera: Scarabaeidae) reveals long-range horizontal transfer events of DNA  
1258 transposons. *Mol Genet Genomics.* 2020;295(6):1339-1353.
- 1259 97. Feschotte C, Zhang X, Wessler SR. Miniature inverted-repeat transposable elements  
1260 and their relationship to established DNA transposons. *Mob DNA II.* 2002;1147–58.
- 1261 98. Feschotte C, Mouchès C. Recent amplification of miniature inverted-repeat transposable  
1262 elements in the vector mosquito *Culex pipiens*: characterization of the Mimo family. *Gene.*  
1263 2000;250:109–16.
- 1264 99. Feschotte C, Swamy L, Wessler SR. Genome-wide analysis of mariner-like transposable  
1265 elements in rice reveals complex relationships with stowaway miniature inverted repeat  
1266 transposable elements (MITEs). *Genetics.* 2003;163:747–58.
- 1267 100. Lu C, Chen J, Zhang Y, Hu Q, Su W, Kuang H. Miniature inverted-repeat transposable  
1268 elements (MITEs) have been accumulated through amplification bursts and play important  
1269 roles in gene expression and species diversity in *Oryza sativa*. *Mol Biol Evol.* 2012;29:1005–

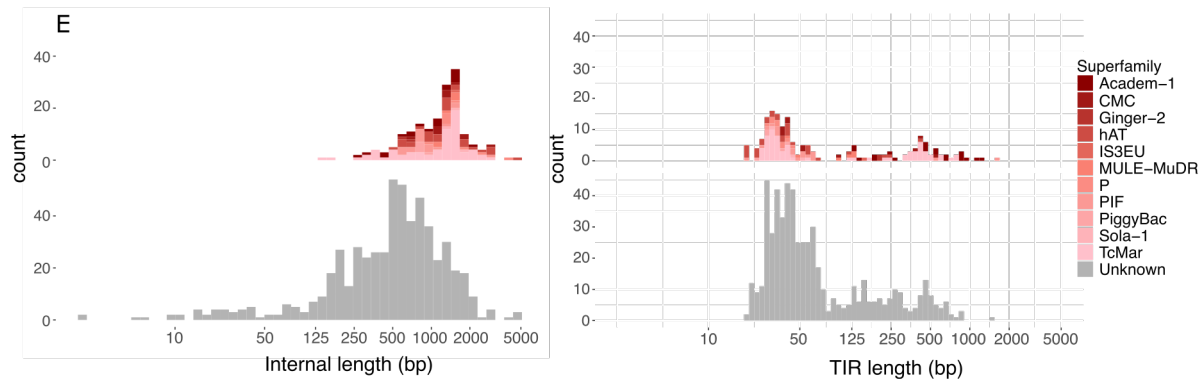
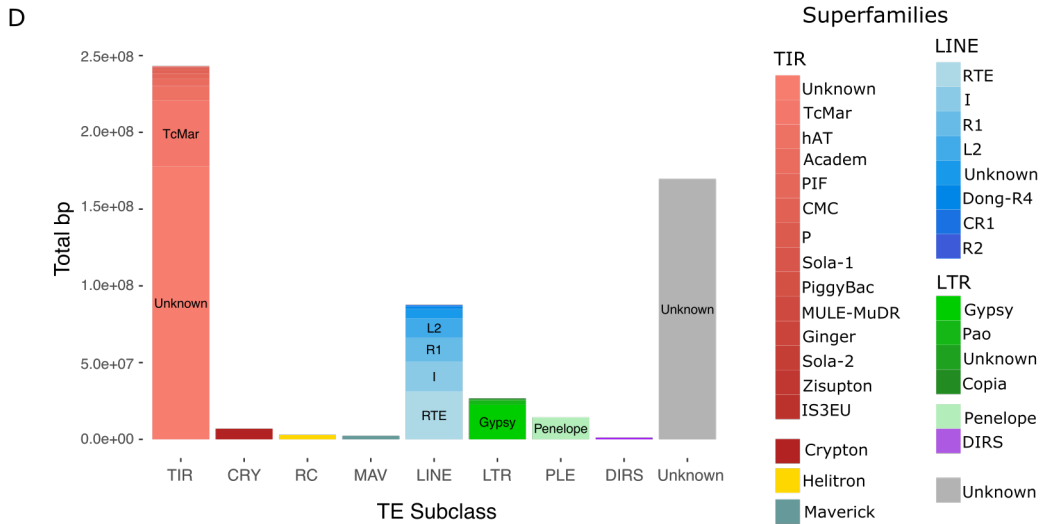
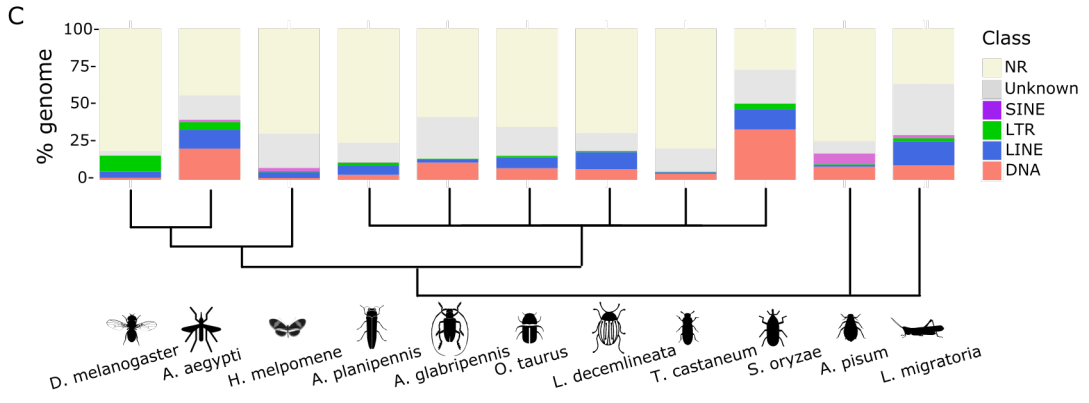
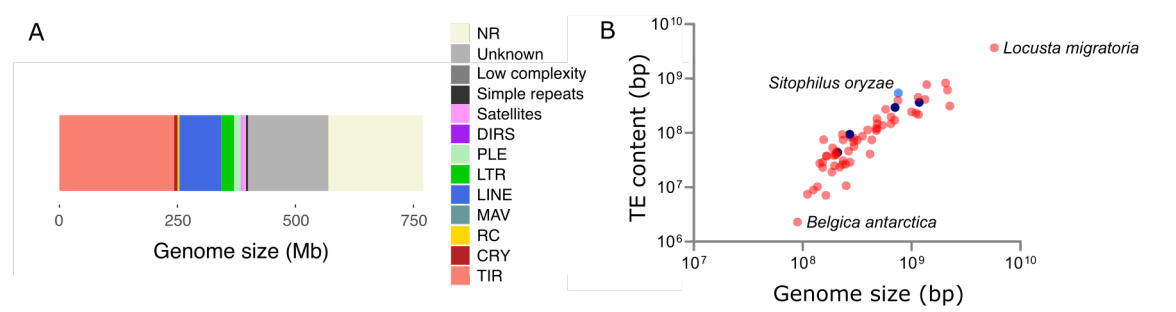
- 1270 17.
- 1271 101. Feng Y. Plant MITEs: Useful tools for plant genetics and genomics. *Genomics*
- 1272 *Proteomics Bioinformatics*. 2003;1:90–100.
- 1273 102. Sela N, Kim E, Ast G. The role of transposable elements in the evolution of non-
- 1274 mammalian vertebrates and invertebrates. *Genome Biol*. 2010;11:R59.
- 1275 103. Petrov DA. DNA loss and evolution of genome size in *Drosophila*. *Genetica*. 2002
- 1276 May;115(1):81-91.
- 1277 104. Petrov DA, Hartl DL. High rate of DNA loss in the *Drosophila melanogaster* and
- 1278 *Drosophila virilis* species groups. *Mol Biol Evol*. 1998;15:293–302.
- 1279 105. Pasyukova EG, Nuzhdin SV. Doc and copia instability in an isogenic *Drosophila*
- 1280 *melanogaster* stock. *Mol Gen Genet*. 1993;240:302–6.
- 1281 106. Ashburner M, Bergman CM. *Drosophila melanogaster*: a case study of a model
- 1282 genomic sequence and its consequences. *Genome Res*. 2005;15:1661–7.
- 1283 107. Czech B, Hannon GJ. One loop to rule them all: The Ping-Pong cycle and piRNA-
- 1284 guided silencing. *Trends Biochem Sci*. 2016;41:324–37.
- 1285 108. Sienski G, Dönertas D, Brennecke J. Transcriptional silencing of transposons by Piwi
- 1286 and Maelstrom and its impact on chromatin state and gene expression. *Cell*. 2012;151:964–
- 1287 80.
- 1288 109. Andersen PR, Tirian L, Vunjak M, Brennecke J. A heterochromatin-dependent
- 1289 transcription machinery drives piRNA expression. *Nature*. 2017;549:54–9.
- 1290 110. Lynch M, Conery JS. The Origins of Genome Complexity. *Science*. 2003;302:1401–4.
- 1291 111. Peona V, Blom MPK, Xu L, Burri R, Sullivan S, Bunikis I, et al. Identifying the causes
- 1292 and consequences of assembly gaps using a multiplatform genome assembly of a bird-of-
- 1293 paradise. *Mol Ecol Resour*. 2021;21(1):263-286.
- 1294 112. Di Genova A, Buena-Atienza E, Ossowski S, Sagot M-F. Efficient hybrid *de novo*
- 1295 assembly of human genomes with WENGAN. *Nat Biotechnol*. 2020;1–9.
- 1296 113. Platt RN II, Blanco-Berdugo L, Ray DA. Accurate transposable element annotation is
- 1297 vital when analyzing new genome assemblies. *Genome Biol Evol*. 2016;8:403–10.

- 1298 114. Maire J, Parisot N, Galvao Ferrarini M, Vallier A, Gillet B, Hughes S, et al. Spatial and  
1299 morphological reorganization of endosymbiosis during metamorphosis accommodates adult  
1300 metabolic requirements in a weevil. *Proc Natl Acad Sci U S A*. 2020;117:19347–58.
- 1301 115. Login FH, Balmand S, Vallier A, Vincent-Monégat C, Vigneron A, Weiss-Gayet M, et al.  
1302 Antimicrobial peptides keep insect endosymbionts under control. *Science*. 2011;334:362–5.
- 1303 116. Nardon P. Obtention d'une souche asymbiotique chez le charançon *Sitophilus sasakii*  
1304 Tak: différentes méthodes d'obtention et comparaison avec la souche symbiotique d'origine.  
1305 *CR Acad Sci Paris D*. 1973;277:981–4.
- 1306 117. Li H. BFC: correcting Illumina sequencing errors. *Bioinformatics*. 2015;31:2885–7.
- 1307 118. Salmela L, Rivals E. LoRDEC: accurate and efficient long read error correction.  
1308 *Bioinformatics*. 2014;30:3506–14.
- 1309 119. Magoc T, Salzberg SL. FLASH: fast length adjustment of short reads to improve  
1310 genome assemblies. *Bioinformatics*. 2011;27:2957–63.
- 1311 120. Chikhi R, Rizk G. Space-efficient and exact de Bruijn graph representation based on a  
1312 Bloom filter. *Algorithms Mol Biol*. 2013;8(1):22.
- 1313 121. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An  
1314 integrated tool for comprehensive microbial variant detection and genome assembly  
1315 improvement. *PLoS One*. 2014;9:e112963.
- 1316 122. Di Genova A, Ruz GA, Sagot M-F, Maass A. Fast-SG: an alignment-free algorithm for  
1317 hybrid assembly. *GigaScience*. 2018;7(5):giy048.
- 1318 123. Mandric I, Zelikovsky A. ScaffMatch: scaffolding algorithm based on maximum weight  
1319 matching. *Bioinformatics*. 2015;31:2632–8.
- 1320 124. Xu G-C, Xu T-J, Zhu R, Zhang Y, Li S-Q, Wang H-W, et al. LR\_Gapcloser: a tiling path-  
1321 based gap closer that uses long reads to complete genome assembly. *GigaScience*.  
1322 2019;8(1):giy157.
- 1323 125. Song L, Shankar DS, Florea L. Rascaf: Improving genome assembly with RNA  
1324 sequencing data. *Plant Genome*. 2016;9:1–12.
- 1325 126. Paulino D, Warren RL, Vandervalk BP, Raymond A, Jackman SD, Birol I. Sealer: a

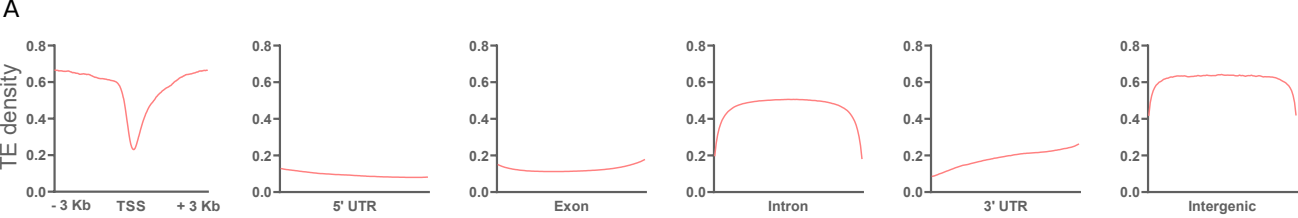
- 1326 scalable gap-closing application for finishing draft genomes. *BMC Bioinformatics*.  
1327 2015;16:230.
- 1328 127. Mikheenko A, Prijbelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome  
1329 assembly evaluation with QUASt-LG. *Bioinformatics*. 2018;34:i142–50.
- 1330 128. Kokot M, Długosz M, Deorowicz S. KMC 3: counting and manipulating k-mer statistics.  
1331 *Bioinformatics*. 2017;33:2759–61.
- 1332 129. Smit AF, Hubley R, Green P. RepeatMasker Open-4.0. <http://www.repeatmasker.org>.  
1333 2013;
- 1334 130. Crescente JM, Zavallo D, Helguera M, Vanzetti LS. MITE Tracker: an accurate  
1335 approach to identify miniature inverted-repeat transposable elements in large genomes.  
1336 *BMC Bioinformatics*. 2018;19:348.
- 1337 131. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:  
1338 improvements in performance and usability. *Mol Biol Evol*. 2013;30:772–80.
- 1339 132. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al.  
1340 Introducing mothur: Open-source, platform-independent, community-supported software for  
1341 describing and comparing microbial communities. *Appl Environ Microbiol*. 2009;75:7537–41.
- 1342 133. Novák P, Ávila Robledillo L, Koblížková A, Vrbová I, Neumann P, Macas J. TAREAN: a  
1343 computational tool for identification and characterization of satellite DNA from unassembled  
1344 short reads. *Nucleic Acids Res*. 2017;45:e111.
- 1345 134. Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. The Dfam community resource of  
1346 transposable element families, sequence models, and genome annotations. *Mob DNA*.  
1347 2021;12:2.
- 1348 135. Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in  
1349 eukaryotic genomes. *Mob DNA*. 2015;6:11.
- 1350 136. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:  
1351 architecture and applications. *BMC Bioinformatics*. 2009;10:421.
- 1352 137. Kapusta A, Suh A. Evolution of bird genomes—a transposon’s-eye view. *Ann N Y Acad*  
1353 *Sci*. 2017;1389:164–85.

- 1354 138. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic  
1355 features. *Bioinformatics*. 2010;26:841–2.
- 1356 139. Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, et al. deepTools2:  
1357 a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.*  
1358 2016;44:W160–5.
- 1359 140. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence  
1360 data. *Bioinformatics*. 2014;30:2114–20.
- 1361 141. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast  
1362 universal RNA-seq aligner. *Bioinformatics*. 2013;29:15–21.
- 1363 142. Liao Y, Smyth GK, Shi W. featureCounts: an efficient general purpose program for  
1364 assigning sequence reads to genomic features. *Bioinformatics*. 2014;30:923–30.
- 1365 143. Lerat E, Fablet M, Modolo L, Lopez-Maestre H, Vieira C. TEtools facilitates big data  
1366 expression analysis of transposable elements and reveals an antagonism between their  
1367 activity and that of piRNA genes. *Nucleic Acids Res.* 2017;45:e17.
- 1368 144. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
1369 RNA-seq data with DESeq2. *Genome Biol.* 2014;15:550.
- 1370 145. Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, et al. Phylogenomics  
1371 resolves the timing and pattern of insect evolution. *Science*. 2014;346:763–7.

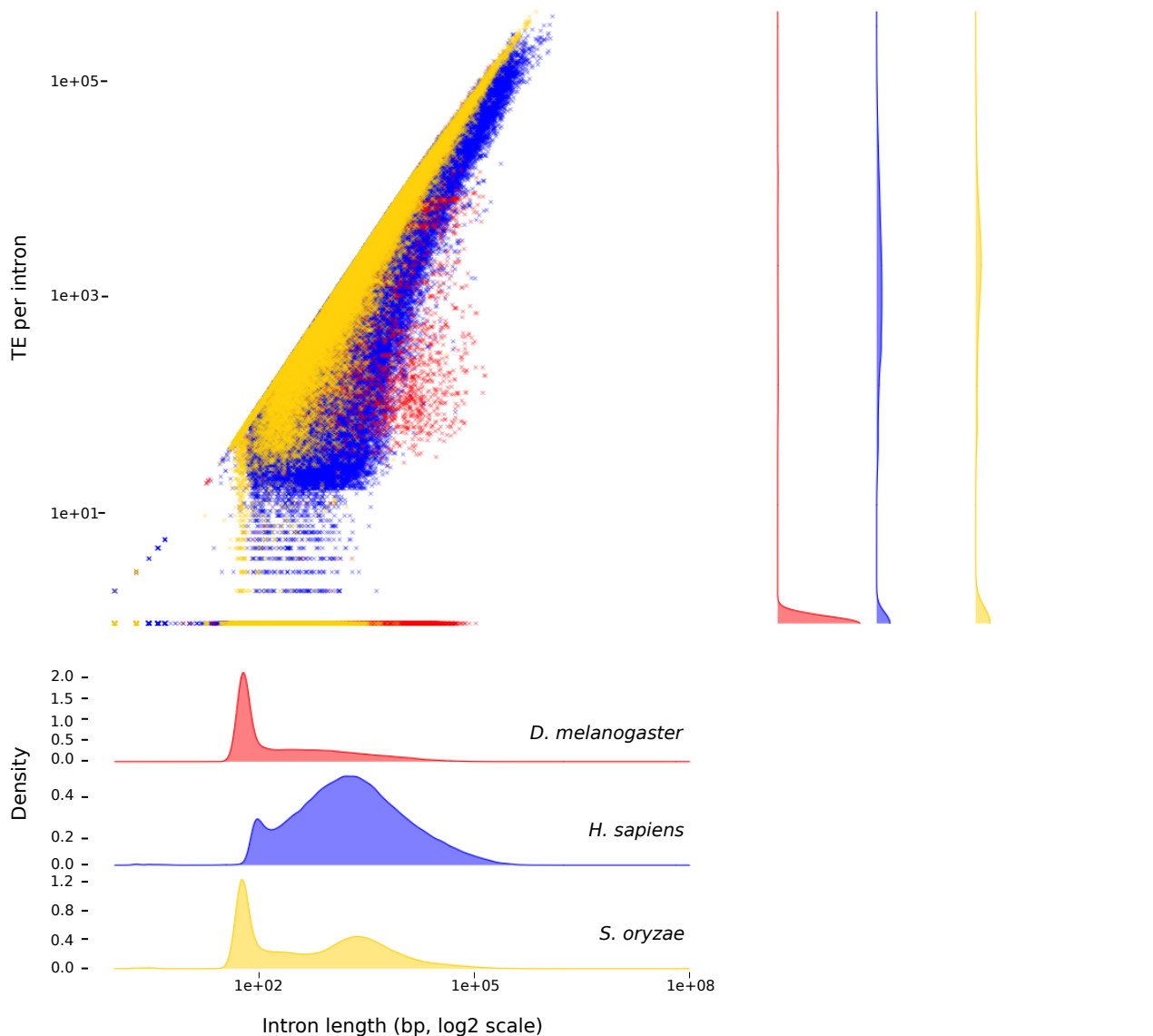
**A****B**

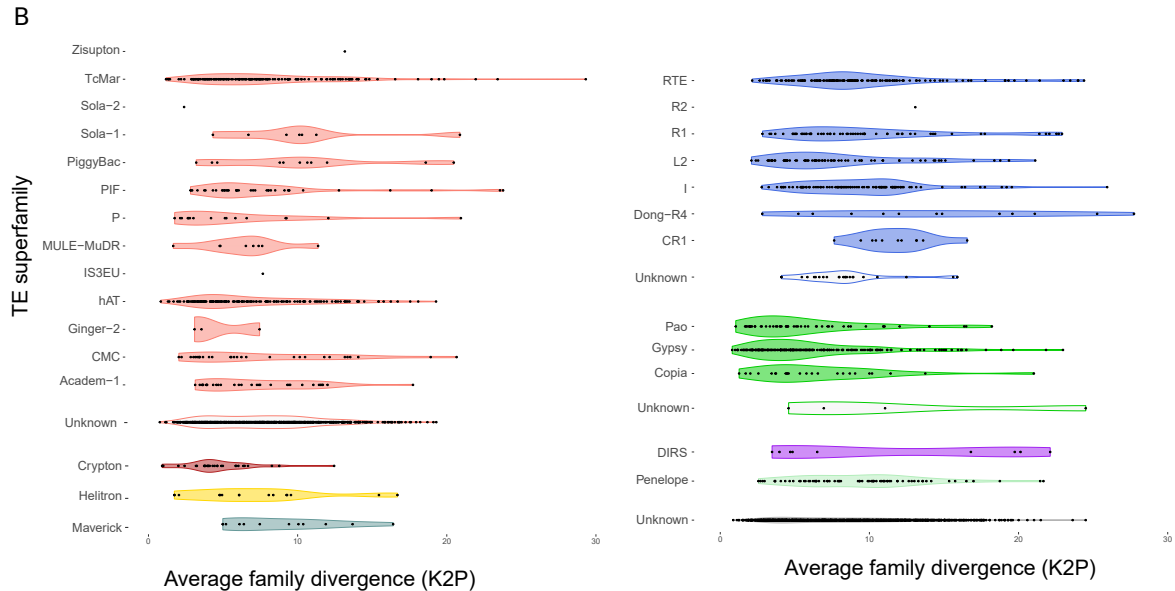
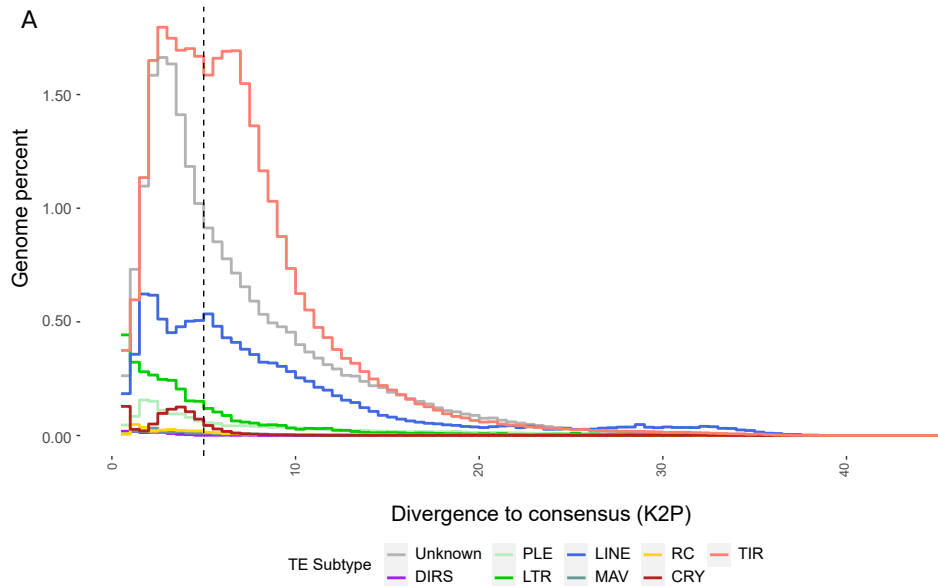


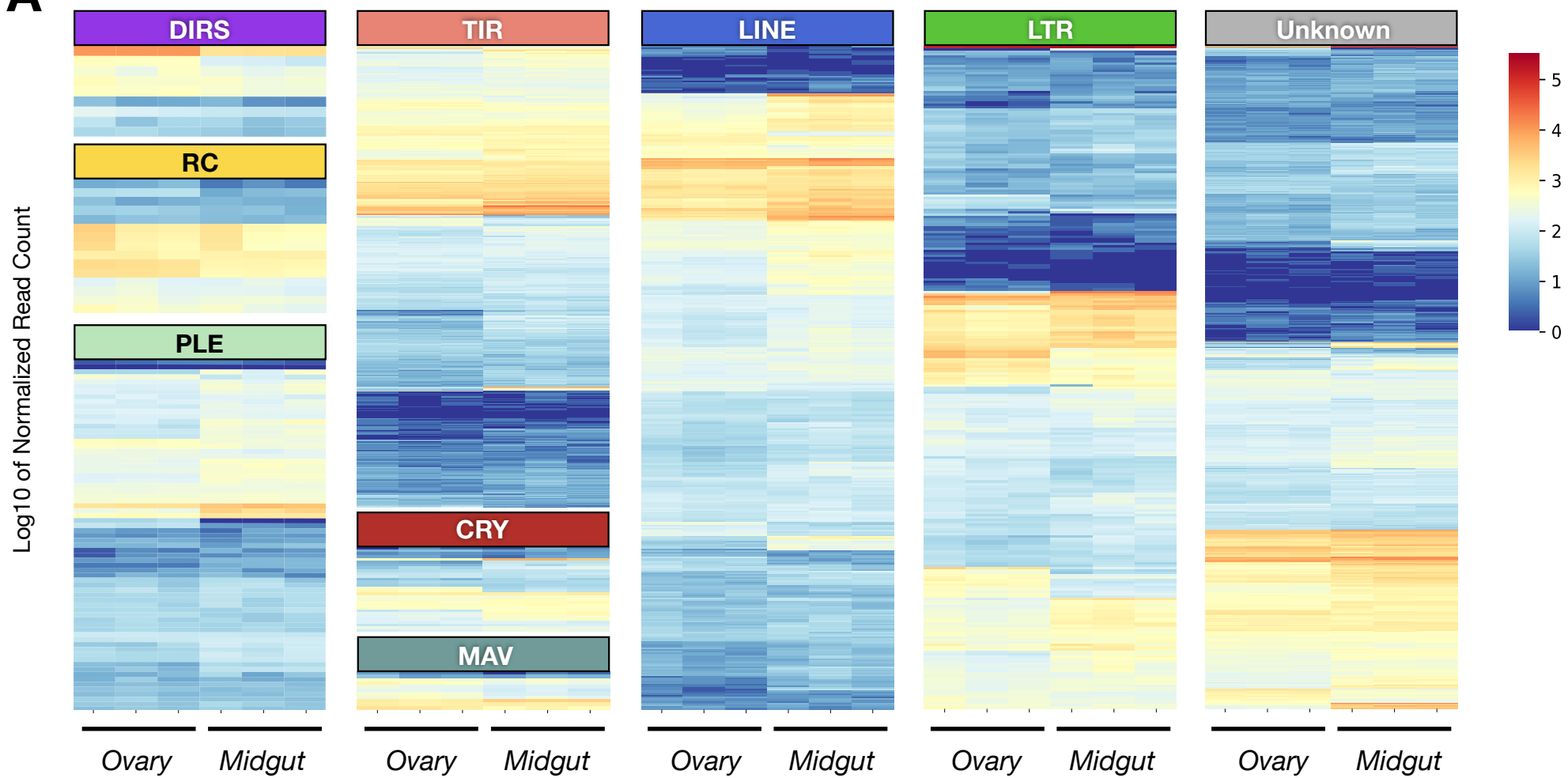
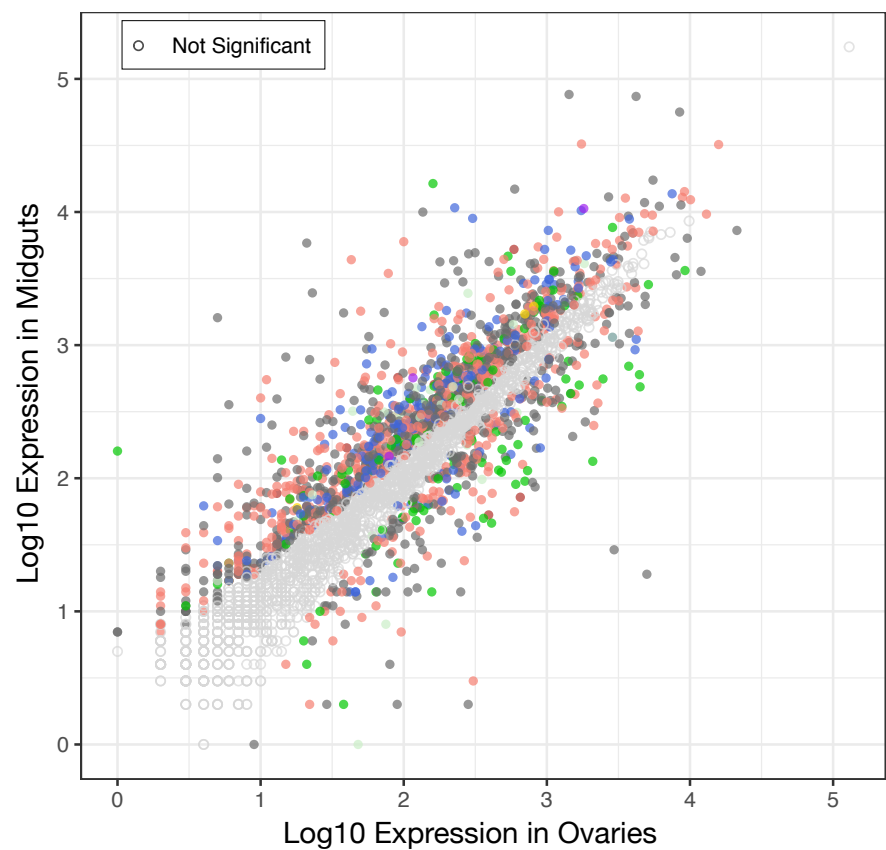




**B**





**A****B****C**