# Insights into the genomic evolution of insects from cricket genomes

Guillem Ylla[1]*, Taro Nakamura[1,2], Takehiko Itoh[3], Rei Kajitani[3], Atsushi Toyoda[4,5], Sayuri Tomonari[6], Tetsuya Bando[7], Yoshiyasu Ishimaru[6], Takahito Watanabe[6], Masao Fuketa[8], Yuji Matsuoka[6,9], Austen A. Barnett[1,10], Sumihare Noji[6], Taro Mito[6]*, Cassandra G. Extavour[1,11]*

1. Department of Organismic and Evolutionary Biology, Harvard University, Cambridge MA, USA
2. Current address: National Institute for Basic Biology, Okazaki, Japan
3. School of Life Science and Technology, Tokyo Institute of Technology, Tokyo, Japan
4. Comparative Genomics Laboratory, National Institute of Genetics, Shizuoka, Japan
5. Advanced Genomics Center, National Institute of Genetics, Shizuoka, Japan
6. Department of Bioscience and Bioindustry, Tokushima University, Tokushima, Japan
7. Graduate School of Medicine, Pharmacology and Dentistry, Okayama University, Okayama, Japan
8. Graduate School of Advanced Technology and Science, Tokushima University, Tokushima, Japan
9. Current address: Department of Biological Sciences, National University of Singapore, Singapore
10. Current address: Department of Natural Sciences, DeSales University, Center Valley PA, USA
11. Department of Molecular and Cellular Biology, Harvard University, Cambridge MA, USA

* correspondence to guillemyllabou@gmail.com, mito.taro@tokushima-u.ac.jp and extavour@oeb.harvard.edu

## Abstract

Most of our knowledge of insect genomes comes from Holometabolous species, which undergo the complete metamorphosis and have genomes under 2Gb with little signs of DNA methylation. In contrast, Hemiemetabolous insects undergo the ancestral incomplete metamorphosis and have larger genomes with high levels of DNA methylation. Hemimetabolous species from the Orthopteran order (grasshoppers and crickets) have some of the largest insect genomes. What drives the evolution of these unusual insect genome sizes, remains unknown. Here we report the sequencing, assembly and annotation of the 1.66-Gb genome of the Mediterranean field cricket *Gryllus bimaculatus*, and the annotation of the 1.60-Gb genome of the Hawaiian cricket *Laupala kohalensis*. We compare these two cricket genomes with those of 14 additional insects, and find evidence that hemimetabolous genomes expanded due to transposable element activity. Based on the ratio of observed to expected CpG sites, we find higher conservation and stronger purifying selection of methylated genes than non-methylated genes. Finally, our analysis suggests an expansion of the *pickpocket* class V gene family in crickets, which we speculate might play a role in the evolution of cricket courtship, including their characteristic chirping.

## Introduction

Much of what we know about insect genome structure and evolution comes from examination of the genomes of insects belonging to a single clade, the Holometabola. This group includes species such as flies and beetles, and is characterized by undergoing complete, or holometabolous, metamorphosis, in which the product of embryogenesis is a larva, which then undergoes an immobile stage called a pupa or chrysalis, during which the larval body plan is abandoned and the new, adult body plan is established. Following the pupal stage, the adult winged insect emerges[1]. This clade of insects includes nearly 90% of extant described insect species[2]. Members of this clade have become prominent model organisms for laboratory research, including the genetic model *Drosophila melanogaster*. Thus, a large proportion of our knowledge of insect biology, genetics, development, and evolution is based on studies of this clade.

Before the evolution of holometabolous metamorphosis, insects developed through incomplete or hemimetabolous metamorphosis. This mode of development is characterized by a generation of the final adult body plan during embryogenesis, followed by gradual physical growth of the hatchling through nymphal stages until the last transition to the sexually mature, winged adult, without major changes in body plan from hatchling to adult[1]. Many extant species maintained this presumed ancestral type of metamorphosis, including crickets, cockroaches, and aphids. Among hemimetabolous insects, most of our current genomic data is from the order Hemiptera (true bugs), which is the sister group to the Holometabola. For the remaining 15 hemimetabolous orders, genomic data remain scarce.

Based on data available to date, genome size and genome methylation show unexplained variation across insects. While most holometabolan species have relatively small genomes (0.2-1.5 pg), hemimetabolous species, and specifically polyneopterans (a taxon comprising 10 major hemimetabolous orders of winged insects with fan-like extensions of the hind wings), display a much larger range of genome sizes (up to 8 pg)[3]. This has led to the hypothesis that there is a genome size threshold at 2 pg (~2 Gb) for holometabolan insect genomes[3]. Studying genome size evolution the polyneopterans order Orthoptera (crickets, grasshoppers, locusts, and katydids) offers a valuable opportunity to investigate potential mechanisms of genome size evolution, as it includes species that have similar predicted gene counts, but have genomes ranging from 1.25 Gb to 16.56 Gb[4]. With respect to the level of CpG DNA methylation, only a few holometabolous species display evidence of genome wide DNA methylation at CpG sites, whereas 30 out of 34 studied polyneopteran species do[5,6]. However, the role of DNA methylation in polyneopteran species, and why it appears to have been lost in many holometabolans, is not clear.

Here, we present the 1.66-Gb genome assembly and annotation of *G. bimaculatus* (Orthoptera), commonly known as the two-spotted cricket, a name derived from the two yellow spots found on the base of the forewings of this species (**Figure 1A**). We also report the first genome annotation for a second cricket species, the Hawaiian cricket *Laupala*

85 *kohalensis*, whose genome assembly was recently made public[7]. *G. bimaculatus* has been
86 widely used as a laboratory research model for decades, in scientific fields including
87 neurobiology and neuroethology[8,9], evo-devo[10], developmental biology[11], and
88 regeneration[12]. Technical advantages of this cricket species as a research model include the
89 fact that *G. bimaculatus* does not require cold temperatures or diapause to complete its life
90 cycle, it is easy to rear in laboratories since it can be fed with generic insect or other pet
91 foods, it is amenable to RNA interference (RNAi) and targeted genome editing[13], stable
92 germline transgenic lines can be established[14], and it has an extensive list of available
93 experimental protocols ranging from behavioral to functional genetic analyses[15].

94 Comparing the two cricket genomes annotated here, with those of 14 other insect species,
95 allowed us to identify three interesting features of these cricket genomes, some of which may
96 relate to their unique biology. First, the differential transposable element (TE) composition
97 between the two cricket species suggests abundant TE activity since they diverged from a
98 last common ancestor, which our results suggest occurred circa 89.2 million years ago (Mya).
99 Second, based on gene CpG depletion, an indirect but robust method to identify typically
100 methylated genes[5,16], we find higher conservation of typically methylated genes than of non-
101 methylated genes. Finally, our gene family expansion analysis reveals an expansion of the
102 *pickpocket* class V gene family in the lineage leading to crickets, which we speculate might
103 play a relevant role in cricket courtship behavior, including their characteristic chirping.

## Results

### *Gryllus bimaculatus* genome assembly

106 We sequenced, assembled, and annotated the 1.66-Gb haploid genome of the white eyed
107 mutant strain[12] of the cricket *G. bimaculatus* **(Figure 1A)**. 50% of the genome is contained
108 within the 71 longest scaffolds (L50), the shortest of them having a length of 6.3 Mb (N50),
109 and 90% of the genome is contained within 307 scaffolds (L90). In comparison to other
110 polyneopteran genomes, our assembly displays high quality in terms of contiguity (N50 and
111 L50), and completeness (BUSCO scores) (**Supplementary Table 1**). Notably, the complete
112 BUSCO scores[17] of this genome assembly at the arthropod and insect levels are 98.50%
113 (C:98.5% [S:97.2%, D:1.3%], F:0.4%, M:1.1%, n:1066) and 97.00% (C:97.0% [S:95.2%,
114 D:1.8%], F:0.8%, M:2.2%, n:1658) respectively, indicating high completeness of this genome
115 assembly **(Table 1)**. The low percentage of duplicated BUSCO genes (1.31%-1.81%)
116 suggests that putative artifactual genomic duplication due to mis-assembly of heterozygotic
117 regions is unlikely.

118

119 ***Table 1:*** Gryllus bimaculatus *genome assembly statistics.*

| Number of Scaffolds | 47,877 |
|---|---|

| | |
|---|---|
| Genome Length (nt) | 1,658,007,496 |
| Genome Length (Gb) | 1.66 |
| Avg. scaffold size (Kb) | 34.63 |
| N50 (Mb) | 6.29 |
| N90 (Mb) | 1.04 |
| L50 | 71 |
| L90 | 307 |
| Complete BUSCO Score – Arthropoda | 98.50% |
| Complete BUSCO Score – Insecta | 97.00% |

120

## Annotation of two cricket genomes

122 The publicly available 1.6-Gb genome assembly of the Hawaiian cricket *L. kohalensis*[7],

123 although having lower assembly quality scores (N50=0.58 Mb, L90 = 3,483) than that of *G.*

124 *bimaculatus*, scores high in terms of completeness, with BUSCO scores of 99.3% at the

125 arthropod level and 97.80% at the insect level **(Supplementary Table 1)**.

126 Using three iterations of the MAKER2 pipeline[18], in which we combined *ab-initio* and

127 evidence-based gene models, we annotated the protein-coding genes in both cricket

128 genomes (**Supplementary Figures 1 & 2** ). We identified 17,871 coding genes and 28,529

129 predicted transcripts for *G. bimaculatus*, and 12,767 coding genes and 13,078 transcripts for

130 *L. kohalensis* **(Table 2)**.

131 To obtain functional insights into the annotated genes, we ran InterProScan[19] for all

132 predicted protein sequences and retrieved their InterPro ID, PFAM domains, and Gene

133 Ontology (GO) terms **(Table 2)**. In addition, we retrieved the best significant BLASTP hit

134 (E-value < 1e-6) for 70-90% of the proteins. Taken together, these methods predicted

135 functions for 75% and 94% of the proteins annotated for *G. bimaculatus* and *L. kohalensis*

136 respectively. We created a novel graphic interface through which interested readers can

137 access, search, BLAST and download the genome data and annotations

138 (http://gbimaculatusgenome.rc.fas.harvard.edu).

139

140 **Table 2:** *Genome annotation summary for the crickets* G. bimaculatus *and* L. kohalensis

| | *G. bimaculatus* | *L. kohalensis* |
|---|---|---|
| Annotated Protein-Coding Genes | 17,871 | 12,767 |

| | | |
|---|---|---|
| Annotated Transcripts | 28,529 | 13,078 |
| % With InterPro ID | 59.56% | 72.52% |
| % With GO-terms | 38.66% | 47.03% |
| % With PFAM motif | 62.44% | 76.59% |
| % With significant BLASTP hit | 73.64% | 93.23% |
| Complete BUSCO-proteome Score – Insecta | 90.50% | 87.20% |
| Repetitive content | 33.69% | 35.51% |
| TE content | 28.94% | 34.50% |
| GC level | 39.93% | 35.58% |

## Abundant Repetitive DNA

We used RepeatMasker[20] to determine the degree of repetitive content in the cricket genomes, using specific custom repeat libraries for each species. This approach identified 33.69% of the *G. bimaculatus* genome, and 35.51% of the *L. kohalensis* genome, as repetitive content **(Supplementary File 1)**. In *G. bimaculatus* the repetitive content density was similar throughout the genome, with the exception of scaffolds shorter than 1Mb (L90), which make up 10% of the genome and have a high density of repetitive content and low gene density **(Figure 1B)**. Because the repetitive content makes genome assemblies more challenging, as observed for the shortest scaffolds of *G. bimaculatus*, we cannot rule out the possibility that the lower contiguity of the *L. kohalensis* genome could lead us to underestimate its repetitive content. This caveat notwithstanding, we observed that transposable elements (TEs) accounted for 28.94% of the *G. bimaculatus* genome, and for 34.50% of the *L. kohalensis* genome. Although the overall proportion of genome made up of TEs was similar between the two cricket species, the proportion of each specific TE class varied greatly **(Figure 1C)**. In *L. kohalensis* the most abundant TE type was long interspersed elements (LINEs), accounting for 20.21% of the genome and 58.58% of the total TE content, while in *G. bimaculatus* LINEs made up only 8.88% of the genome and 30.68% of the total TE content. The specific LINE subtypes LINE1 and LINE3 appeared at a similar frequency in both cricket genomes (<0.5%), while the LINE2 subtype was over five times more represented in *L. kohalensis,* covering 10% of the genome (167 Mb). On the other hand, DNA transposons accounted for 8.61% of the *G. bimaculatus* genome, but only for 3.91% of the *L. kohalensis* genome.

## DNA methylation

CpG depletion, calculated as the ratio between observed versus the expected incidence of a cytosine followed by a guanine ($CpG_{o/e}$), is considered a reliable indicator of DNA methylation. This is because spontaneous C to T mutations occur more frequently on

167   methylated CpGs than unmethylated CpGs[16]. Thus, genomic regions that undergo
168   methylation are eventually CpG-depleted. We calculated the $CpG_{o/e}$ value for each predicted
169   protein-coding gene for the two cricket species. In both species, we observed a clear bimodal
170   distribution of $CpG_{o/e}$ values (**Figure 2A**). One interpretation of this distribution is that the
171   peak corresponding to lower $CpG_{o/e}$ values contains genes that are typically methylated, and
172   the peak of higher $CpG_{o/e}$ contains genes that do not undergo DNA methylation. Under this
173   interpretation, some genes have non-random differential DNA methylation in crickets. To
174   quantify the genes in the two putative methylation categories, we set a $CpG_{o/e}$ threshold as
175   the value of the point of intersection between the two normal distributions (**Figure 2A**).
176   After applying this cutoff, 44% of *G. bimaculatus* genes and 45% of *L. kohalensis* genes were
177   identified as CpG-depleted.

178   A GO enrichment analysis of the genes above and below the $CpG_{o/e}$ threshold defined above
179   revealed clear differences in the predicted functions of genes belonging to each of the two
180   categories. Strikingly, however, genes in each threshold category had functional similarities
181   across the two cricket species (**Figure 3**). Genes with low $CpG_{o/e}$ values, which are likely
182   those undergoing methylation, were enriched for functions related to DNA replication and
183   regulation of gene expression (including transcriptional, translational, and epigenetic
184   regulation), while genes with high $CpG_{o/e}$ values, suggesting little or no methylation, tended
185   to have functions related to metabolism, catabolism, and sensory systems.

186   To assess whether the predicted distinct functions of high- and low- $CpG_{o/e}$ value genes were
187   specific to crickets, or were a potentially more general trend of insects with DNA methylation
188   systems, we analyzed the predicted functions of genes with different $CpGo/e$ values in the
189   honeybee *Apis mellifera*, the first insect for which evidence for DNA methylation was
190   robustly described and studied[21,22], and the thrips *Frankliniella occidentalis*. We found that
191   in both *F. occidentalis* and *A. mellifera,* CpG-depleted genes were enriched for similar
192   functions as those observed in cricket CpG-depleted genes (**Figure 3 and Supplementary**
193   **Figure 3**). Specifically, 23GO terms were significantly enriched in all four studied insects,
194   and 15 additional GO terms were significantly enriched in the three hemimetabolous insects.
195   In the same way, high $CpG_{o/e}$ genes in all four insects were enriched for similar functions (8
196   GO-terms commonly enriched in all insects; **Supplementary Figure 3**).

197   Additionally, we observed that the proportion of species-specific genes was higher within
198   the high $CpG_{o/e}$ peak for all four insects (**Figure 2C**). In contrast, 86-96% of the genes
199   belonging to the low $CpG_{o/e}$ peak had an orthologous gene in at least one of the other studied
200   insect species. Furthermore, we observed 2,182 orthogroups whose members always
201   belonged to the low $CpG_{o/e}$ peak in all four species, and 728 orthogroups whose members
202   always belonged to the high $CpG_{o/e}$ peak, indicating that orthologous genes are likely to share
203   methylation state across these four insect species (**Figure 2B and Supplementary Figure**
204   **4**). Interestingly, 666 genes belonged to the low $CpG_{o/e}$ peak in the three hemimetabolous

205 species (*G. bimaculatus, L. kohalensis*, and *F. occidentallis)*, but to the high $CpG_{o/e}$ peak in the
206 holometabolous *A. mellifera*.

207 Taken together, these results suggest that genes that are typically methylated tend to be
208 more conserved across species, which could imply low evolutionary rates and strong
209 selective pressure. To test this hypothesized relationship between low $CpG_{o/e}$ and low
210 evolutionary rates, we compared the dN/dS values of 1-to-1 orthologous genes belonging to
211 the same $CpG_{o/e}$ peak between the two cricket species. We found that CpG-depleted genes in
212 both crickets had significantly lower dN/dS values than non-CpG-depleted genes (p-
213 value<0.05; **Figure 2D**), consistent with stronger purifying selection on CpG-depleted genes.

## Phylogenetics and gene family expansions

215 To study the genome evolution of these cricket lineages, we compared the two cricket
216 genomes with those of 14 additional insects, including members of all major insect lineages
217 with special emphasis on hemimetabolous species. For each of these 16 insect genomes, we
218 retrieved the longest protein per gene and grouped them into orthogroups (OGs), which we
219 called "gene families" for the purpose of this analysis. The 732 OGs containing a single
220 protein per insect, namely single copy orthologs, were used to infer a phylogenetic tree for
221 these 16 species. The obtained species tree topology was in accordance with the currently
222 understood insect phylogeny[23]. Then, we used the Misof et al. (2014) dated phylogeny to
223 calibrate our tree on four different nodes, which allowed us to estimate that the two cricket
224 species diverged circa 89.2 million years ago.

225 Our gene family expansion/contraction analysis using 59,516 OGs identified 18 gene families
226 that were significantly expanded (p-value<0.01) in the lineage leading to crickets. In
227 addition, we identified a further 34 and 33 gene family expansions specific to *G. bimaculatus*
228 and *L. kohalensis* respectively. Functional analysis of these expanded gene families
229 (**Supplementary File 2**) revealed that the cricket-specific gene family expansions included
230 *pickpocket* genes, which are involved in mechanosensation in *Drosophila melanogaster* as
231 described in the following section.

232

## Expansion of *pickpocket* genes

234 In *D. melanogaster*, the complete *pickpocket* gene repertoire is composed of 6 classes
235 containing 31 genes. We found cricket orthologs of all 31 *pickpocket* genes across seven of
236 our OGs, and each OG predominantly contained members of a single *pickpocket* class. We
237 used all the genes belonging to these 7 OGs to build a *pickpocket* gene tree, using the
238 predicted *pickpocket* orthologs from 16 insect species (**Figure 3**; **Supplementary Table 2**).
239 This gene tree allowed us to classify the different *pickpocket* genes in each of the 16 species.

240 One orthogroup, which contained eight members of the *pickpocket* gene family of *D.*
241 *melanogaster,* appeared to be significantly expanded to 14 or 15 members in crickets.

242 Following the classification of *pickpocket* genes used in *Drosophila spp.*[24] we determined that
243 the specific gene family expanded in crickets was *pickpocket* class V **(Figure 3).** In *D.*
244 *melanogaster* this class contains eight genes: *ppk* (*ppk1*), *rpk* (*ppk2*), *ppk5, ppk8, ppk12,*
245 *ppk17, ppk26,* and *ppk28* [24]. Our analysis suggests that the class V gene family contains 15
246 and 14 genes in *G bimaculatus* and *L. kohalensis* respectively. In contrast, their closest
247 analyzed relative, the locust *Locusta migratoria,* has only five such genes.

248

249 The *pickpocket* genes in crickets tended to be grouped in genomic clusters **(Figure 1B)**. For
250 instance, in *G. bimaculatus* nine of the 15 class V *pickpocket* genes were clustered within a
251 region of 900Kb, and four other genes appeared in two groups of two. In the *L. kohalensis*
252 genome, although this genome is more fragmented than that of *G. bimaculatus*
253 (**Supplementary Table 1**), we observed five clusters containing between two and five genes
254 each.

255 In *D. melanogaster,* the *pickpocket* gene *ppk1* belongs to class V and is involved in functions
256 related to stimulus perception and mechanotransduction[25]. For example, in larvae, this gene
257 is required for mechanical nociception[26], and for coordinating rhythmic locomotion[27]. *ppk* is
258 expressed in sensory neurons that also express the male sexual behavior determiner *fruitless*
259 *(fru)* [28-30].

260 To determine whether *pickpocket* genes in crickets are also expressed in the nervous system,
261 we checked for evidence of expression of *pickpocket* genes in the publicly available RNA-seq
262 libraries for the *G. bimaculatus* prothoracic ganglion[9]. This analysis detected expression (>20
263 transcripts per kilobase million, TPMs) of five *pickpocket* genes, four of them belonging to
264 class V, in the *G. bimaculatus* nervous system. In the same ganglionic RNA-seq libraries, we
265 also detected the expression of *fru* **(Supplementary Table 3)**. Out of the four *pickpocket*
266 genes, only one was detected in embryonic RNA-seq libraries. All four genes together with
267 *fru* were detected in wild type leg transcriptomes, and their expression was found to be
268 higher than wild type in a transcriptome from regenerating legs **(Supplementary Table 4)**.

269

## Discussion

### The importance of cricket genomes

272 Sequencing and analyzing genomes from underrepresented clades allow us to get a more
273 complete picture of genome diversity across the tree of life, and can provide insights
274 regarding their evolution. Since the first sequenced insect genome, that of *D. melanogaster*,
275 was made publicly available in 2000[31], the field of holometabolous genomics has flourished,
276 and this clade became the main source of subsequent genomic information for insects. The
277 first hemimetabolous genome was not available until ten years later, with the publication of

278    the genome sequence and annotation of the Pea aphid (*Acyrthosiphon pisum*)[32]. When even
279    more recently, polyneopteran genome sequences became available[33-36], some of their
280    distinct characteristics, such as their length and DNA methylation profiles, began to be
281    appreciated. Genome data are also very important as they can help establish an animal
282    species as tractable experimental models. *G. bimaculatus* is a common laboratory research
283    animal used in neuroethology, developmental and regeneration biology studies[12,15]. It is our
284    hope that the availability of the annotated genome presented here will encourage other
285    researchers to adopt this cricket as a model organism, and facilitate development of new
286    molecular genetic manipulation tools.

287    Moreover, we note that crickets are currently in focus as a source of animal protein for
288    human consumption and for vertebrate livestock. Crickets possess high nutritional value,
289    having a high proportion of protein for their body weight (>55%), and containing the
290    essential linoleic acid as their most predominant fatty acid[37-39]. Specifically, the cricket *G.*
291    *bimaculatus* has traditionally been consumed in different parts of the world including
292    northeast Thailand, which recorded 20,000 insect farmers in 2011[40]. Studies have reported
293    no evidence for toxicological effects related to oral consumption of *G. bimaculatus* by
294    humans[41,42], neither were genotoxic effects detected using three different mutagenicity
295    tests[43]. A rare but known health risk associated with cricket consumption, however, is
296    sensitivity and allergy to crickets[44,45]. Nevertheless, not only is the cricket *G. bimaculatus*
297    considered generally safe for human consumption, several studies also suggest that
298    introducing crickets into one's diet may confer multiple health benefits[46-48].Crickets might
299    therefore be part of the solution to the problem of feeding a worldwide growing population
300    in a sustainable way. However, most of the crops and livestock that humans eat have been
301    domesticated and subjected to strong artificial selection for hundreds or even thousands of
302    years to improve their characteristics most desirable for humans, including size, growth rate,
303    stress resistance, and organoleptic properties[49-52]. In contrast, to our knowledge, crickets
304    have never been selected based on any food-related characteristic. The advent of genetic
305    engineering techniques has accelerated domestication of some organisms[53]. These
306    techniques have been used, for instance, to improve the nutritional value of different crops,
307    or to make them tolerant to pests and climate stress[49,54]. Crickets are naturally nutritionally
308    rich[39], but in principle, their nutritional value could be further improved, for example by
309    increasing vitamin content or Omega-3 fatty acids proportion. In addition, other issues that
310    present challenges to cricket farming could potentially be addressed by targeted genome
311    modification, which can be achieved in *G. bimaculatus* using Zinc finger nucleases, TALENs,
312    or CRISPR/Cas9. These challenges include sensitivity to common insect viruses, aggressive
313    behavior resulting in cannibalism, complex mating rituals, and relatively slow growth rate.

314

## Comparing cricket genomes to other insect genomes

The annotation of these two cricket genomes was done by combining *de novo* gene models, homology-based methods, and the available RNA-seq and ESTs. This pipeline allowed us to predict 17,871 genes in the *G. bimaculatus* genome, similar to the number of genes reported for other hemimetabolous insect genomes including the locust *L. migratoria* (17,307)[33] and the termites *Cryptotermes secundus* (18,162)[55], *Macrotermes natalensis* (16,140)[36] and *Zootermopsis nevadensis*, (15,459)[35]. The slightly lower number of protein-coding genes annotated in *L. kohalensis* (12,767) may be due to the lesser amount of RNA-seq data available for this species, leading to higher assembly fragmentation, which challenges gene annotation. Nevertheless, the BUSCO scores are similar between the two crickets, and the proportion of annotated proteins with putative orthologous genes in other species (proteins with significant BLAST hits; see methods) for *L. kohalensis* is higher than for *G. bimaculatus*. This suggests the possibility that we may have successfully annotated most conserved genes, but that highly derived or species-specific genes might be missing from our annotations.

## TEs and genome size evolution

Approximately 35% of the genome of both crickets corresponds to repetitive content. This is substantially less than the 60% reported for the genome of *L. migratoria*[33]. This locust genome is one of the largest sequenced insect genomes to date (6.5 Gb) but has a very similar number of annotated genes (17,307) to those we report for crickets. We hypothesize that the large genome size difference between these orthopteran species is due to the TE content, which has also been correlated with genome size in multiple eukaryote species[56,57].

Furthermore, we hypothesize that the differences in the TE composition between the two crickets are the result of abundant and independent TE activity since their divergence around 89.2 Mya. This, together with the absence of evidence for large genome duplication events in this lineage, leads us to hypothesize that the ancestral orthopteran genome was shorter than those of the crickets studied here (1.6 Gb for *G. bimaculatus* and 1.59 Gb for *L. kohalensis*) which are in the lowest range of orthopteran genome sizes[58]. In summary, we propose that the wide range of genome sizes within Orthoptera, reaching as high as 8.55 Gb in the locust *Schistocerca gregaria,* and 16.56 Gb in the grasshopper *Podisma pedestris*[4,59], is likely due to TE activity since the time of the last orthopteran ancestor. These observations are consistent with the results reported by Palacios-Gimenez, et al.[60] of massive and independent recent TE accumulation in four chromosome races of the grasshopper *Vandiemenella viatica*.

There is a clear tendency of polyneopteran genomes to be much longer than those of the holometabolous genomes **(Figure 4)**. Two currently competing hypotheses are that (1) the ancestral insect genome was small, and was expanded outside of Holometabola, and (2) the

352  ancestral insect genome was large, and it was compressed in the Holometabola[3]. Our
353  observations are consistent with the first of these hypotheses.

354

## DNA Methylation

356  Most holometabolan species, including well-studied insects like *D. melanogaster* and
357  *Tribolium castaneum,* do not perform DNA methylation, or they do it at very low levels[6,61].
358  The honeybee *A. mellifera* was one of the first insects for which functional DNA methylation
359  was described[21]. Although this DNA modification was initially proposed to be associated
360  with the eusociality of these bees[22], subsequent studies showed that DNA methylation is
361  widespread and present in different insect lineages independently of social behavior [5]. DNA
362  methylation also occurs in other non-insect arthropods[62].

363  While the precise role of DNA methylation in gene expression regulation remains unclear,
364  our analysis suggests that cricket CpG-depleted genes (putatively hypermethylated genes)
365  show signs of purifying selection, tend to have orthologs in other insects, and are involved in
366  basic biological functions related to DNA replication and the regulation of gene expression.
367  These enriched functions are in agreement with previous observations that DNA methylated
368  genes in arthropods tend to perform housekeeping functions[6,63]. These predicted functions
369  differ from those of the non-CpG depleted genes (putatively hypomethylated genes), which
370  appear to be involved in signaling pathways, metabolism, and catabolism. These predicted
371  functional categories may be conserved from crickets over circa 345 million years of
372  evolution, as we also detect the same pattern in the honeybee and a thrips species.

373  Taken together, these observations suggest a potential relationship between DNA
374  methylation, sequence conservation, and function for many cricket genes. Nevertheless,
375  based on our data, we cannot determine whether the methylated genes are highly conserved
376  because they are methylated, or because they perform basic functions that may be regulated
377  by DNA methylation events. In the cockroach *Blattella germanica*, DNA methyltransferase
378  enzymes and genes with low $CpG_{o/e}$ values show an expression peak during the maternal to
379  zygotic transition[64], and functional analysis has shown that the DNA methyltransferase 1 is
380  essential for early embryo development in this cockroach[65]. These results in cockroaches,
381  together with our observations, leads us to speculate that at least in Polyneopteran species,
382  DNA methylation might contribute to the maternal zygotic transition by regulating essential
383  genes involved in DNA replication, transcription, and translation.

## *pickpocket* gene expansion

385  The *pickpocket* genes belong to the Degenerin/epithelial Na+ channel (DEG/ENaC) family,
386  which were first identified in *Caenorhabditis elegans* as involved in mechanotransduction[25].
387  The same family of ion channels was later found in many multicellular animals, with a
388  diverse range of functions related to mechanoreception and fluid–electrolyte homeostasis[66].

389   Most of the information on their roles in insects comes from studies in *D. melanogaster*. In
390   this fruit fly, *pickpocket* genes are involved in neural functions including NaCl taste[67],
391   pheromone detection [68], courtship behavior [69], and liquid clearance in the larval trachea[66].

392   In *D. melanogaster* adults, the abdominal ganglia mediate courtship and postmating
393   behaviors through neurons expressing *ppk* and *fru*[28-30]. In *D. melanogaster* larvae, *ppk*
394   expression in dendritic neurons is required to control the coordination of rhythmic
395   locomotion[27]. In crickets, the abdominal ganglia are responsible for determining song
396   rhythm[70]. Moreover, we find that in *G. bimaculatus*, both *ppk* and *fru* gene expression are
397   detectable in the adult prothoracic ganglion. These observations suggest the possibility that
398   class V *pickpocket* genes could be involved in song rhythm determination in crickets through
399   their expression in abdominal ganglia.

400   This possibility is consistent with the results of multiple quantitative trait locus (QTL)
401   studies done in cricket species from the genus *Laupala*, which identified genomic regions
402   associated with mating song rhythm variations and female acoustic preference[71]. The 179
403   scaffolds that the authors reported being within one logarithm of the odds (LOD) of the seven
404   QTL peaks, contained five *pickpocket* genes, three of them from class V and two from class
405   IV. One of the two class IV genes also appears within a QTL peak of a second experiment[7,72].
406   Xu and Shaw [73] found that a scaffold in a region of LOD score 1.5 of one of their minor linkage
407   groups (LG3) contains *slowpoke*, a gene that affects song interpulse interval in *D.*
408   *melanogaster,* and this scaffold also contains two class III pickpocket genes (**Supplementary**
409   **Table 5**).

410   In summary, the roles of *pickpocket* genes in controlling rhythmic locomotion, courtship
411   behavior, and pheromone detection in *D. melanogaster*, their appearance in genomic regions
412   associated with song rhythm variation in *Laupala*, and their expression in *G. bimaculatus*
413   abdominal ganglia, lead us to speculate that the expanded *pickpocket* gene family in cricket
414   genomes could be playing a role in regulating rhythmic wing movements and sound
415   perception, both of which are necessary for mating[15]. We note that Xu and Shaw [73]
416   hypothesized that song production in crickets is likely to be regulated by ion channels, and
417   that locomotion, neural modulation, and muscle development are all involved in singing[73].
418   However, further experiments, which could take advantage of the existing RNAi and genome
419   modification protocols for *G. bimaculatus*[13], will be required to test this hypothesis.

420

421   In conclusion, the *G. bimaculatus* genome assembly and annotation presented here is a
422   source of information and an essential tool that we anticipate will enhance the status of this
423   cricket as a modern functional genetics research model. This genome may also prove useful
424   to the agricultural sector, and could allow improvement of cricket nutritional value,
425   productivity, and reduction of allergen content. Annotating a second cricket genome, that of
426   *L. kohalensis*, and comparing the two genomes, allowed us to unveil possible

427 synapomorphies of cricket genomes, and to suggest potentially general evolutionary trends
428 of insect genomes.

429

## Materials and Methods

### DNA isolation

432 The *G. bimaculatus* white-eyed mutant strain was reared at Tokushima University, at 29±1
433 °C and 30-50% humidity under a 10-h light, 14-h dark photoperiod. Testes of a single male
434 adult of the *G. bimaculatus* white-eyed mutant strain were used for DNA isolation and short-
435 read sequencing. We used DNA from testes of an additional single individual to make a long
436 read PacBio sequencing library to close gaps in the genome assembly. Because sex
437 differentiation in the cricket *G. bimaculatus* is determined by the XX/XO system[74], genomic
438 DNA extracted from males contains the full set of chromosomes; males were therefore
439 chosen for genomic DNA isolation.

440

### Genome Assembly

442 Paired-end libraries were generated with insert sizes of 375 and 500 bp, and mate-pair
443 libraries were generated with insert sizes of 3, 5, 10, and 20kb. Libraries were sequenced
444 using the Illumina HiSeq 2000 and HiSeq 2500 sequencing platforms. This yielded a total of
445 127.4 Gb of short read paired-end data, that was subsequently assembled using the *de novo*
446 assembler Platanus (v. 1.2.1)[75]. Scaffolding and gap closing were performed using total 138.2
447 Gb of mate-pair data. A further gap closing step was performed using long reads generated
448 by the PacBio RS system. The 4.3 Gb of PacBio subread data were used to fill gaps in the
449 assembly using PBjelly (v. 15.8.24)[76].

450

### Repetitive Content Masking

452 We generated a custom repeat library for each of the two cricket genomes by combining the
453 outputs from homology-based and *de novo* repeat identifiers, including the LTRdigest
454 together with LTRharvest[77], RepeatModeler/RepeatClassifier
455 (www.repeatmasker.org/RepeatModeler), MITE tracker[78], TransposonPSI
456 (http://transposonpsi.sourceforge.net), and the databases SINEBase[79] and RepBase[80]. We
457 removed redundancies from the library by merging sequences that were greater than 80%
458 similar with usearch [81], and classified them with RepeatClassifier. Sequences classified as
459 "unknown" were searched with BLASTX against the 9,229 reviewed proteins of insects from
460 UniProtKB/Swiss-Prot. Those sequences with a BLAST hit (E-value < 1e-10) against a

461  protein not annotated as a transposase, transposable element, copia protein, or transposon
462  were removed from the custom repeat library. The custom repeat library was provided to
463  RepeatMasker version open-4.0.5 to generate the repetitive content reports, and to the
464  MAKER2 pipeline to mask the genome.

## Protein-Coding Genes Annotation

466  We performed genome annotations through three iterations of the MAKER2 (v2.31.8)
467  pipeline[18] combining *ab-initio* gene models and evidence-based models. For the *G.*
468  *bimaculatus* genome annotation, we provided the MAKER2 pipeline with the 43,595 *G.*
469  *bimaculatus* nucleotide sequences from NCBI, an assembled developmental transcriptome
470  [82], an assembled prothoracic ganglion transcriptome[9], and a genome-guided transcriptome
471  generated with StringTie[83] using 30 RNA-seq libraries (accession numbers: DRA011174 and
472  DDBJ DRA11117) mapped to the genome with HISAT2[84]. As alternative ESTs and protein
473  sequences, we provided MAKER2 with 14,391 nucleotide sequences from *L. kohalensis*
474  available at NCBI, and an insect protein database obtained from UniProtKB/Swiss-Prot[85].

475  For the annotation of the *L. kohalensis* genome, we ran the MAKER2 pipeline with the 14,391
476  *L. kohalensis* nucleotide sequences from NCBI, the assembled *G. bimaculatus* developmental
477  and prothoracic ganglion transcriptomes described above, and the 43,595 NCBI nucleotide
478  sequences. As protein databases, we provided the insect proteins from UniProtKB/Swiss-
479  Prot plus the proteins that we annotated in the *G. bimaculatus* genome.

480  For both crickets, we generated *ab-initio* gene models with GeneMark-ES[86] in self-training
481  mode, and with Augustus[87] trained with BUSCO v3[17]. After each of the first two MAKER2
482  iterations, additional gene models were obtained with SNAP[88] trained with the annotated
483  genes.

484  Functional annotations were obtained using InterProScan[19], which retrieved the
485  InterProDomains, PFAM domains, and GO-terms. Additionally, we ran a series of BLAST
486  rounds from more specific to more generic databases, to assign a descriptor to each
487  transcript based on the best BLAST hit. The first round of BLAST was against the reviewed
488  insect proteins from UniProtKB/Swiss-Prot. Proteins with no significant BLAST hits (E-value
489  < 1e-6) went to a second round against all proteins from UniProtKB/TrEMBL, and those
490  without a hit with E-value<1e-6 were used in the final round of BLAST against all proteins
491  from UniProtKB/Swiss-Prot.

492  A detailed pipeline scheme is available in **Supplementary Figures 1 & 2,** and the
493  annotation scripts are available on GitHub
494  (https://github.com/guillemylla/Crickets_Genome_Annotation).

495

## Quality Assessment

Genome assembly statistics were obtained with assembly-stats (https://github.com/sanger-pathogens/assembly-stats). BUSCO (v3.1.0)[17] was used to assess the level of completeness of the genome assemblies ('-m geno') as well as that of the gene annotations ('-m prot') at both arthropod ('arthropoda_odb9') and insect ('insecta_odb9') levels.

## $CpG_{o/e}$ Analysis

We used the genome assemblies and their gene annotations from this study for the two cricket species, and retrieved publicly available annotated genomes from the other 14 insect species (**Supplementary Table 1** ). The gene annotation files (in gff format) were used to obtain the amino-acid and CDS sequences for each annotated protein-coding gene per genome using gffread, with options "-y" and "-x" respectively. The $CpG_{o/e}$ value per gene was computed as the observed frequency of CpGs ($f_{CpG}$) divided by the product of C and G frequencies ($f_C$ and $f_G$) $f_{CpG}/f_C*f_G$ in the longest CDS per gene for each of the 16 studied insects. $CpG_{o/e}$ values larger than zero and smaller than two were retained and represented as density plots (**Figures 2 & 4**).

The distributions of gene $CpG_{o/e}$ values per gene of the two crickets, the honeybee *A. mellifera*, and the thrips *F. occidentalis*, were fitted with a mixture of normal distributions using the mixtools R package[89]. This allowed us to obtain the mean of each distribution, the standard errors, and the interception point between the two distributions, which was used to categorize the genes into low $CpG_{o/e}$ and high $CpG_{o/e}$ bins. For these two bins of genes, we performed a GO-enrichment analysis (based on GO-terms previously obtained using InterProScan) of Biological Process terms using the TopGO package[90] with all genes as universe, minimum node size of 10, the weight01 algorithm and the Fisher statistic. The GO terms with a p-value<0.05 were considered significantly enriched. Those GO terms significantly enriched in at least one gene set are shown in **Supplementary Figure 3**, and a subset of them with p-value<0.0001 are shown in **Figure 3**. In both figures, the size of the circle represents the percentage of enriched genes inside the set compared to all genes with the given GO term.

For each of the genes belonging to low and high $CpG_{o/e}$ categories in each of the four insect species, we retrieved their orthogroup identifier from our gene family analysis, allowing us to assign putative methylation status to orthogroups in each insect. Then we used the UpSet R package[91] to compute and display the number of orthogroups exclusive to each combination as an UpSet plot.

## dN/dS Analysis

We first aligned the longest predicted protein product of the single-copy-orthologs of all protein-coding genes between the two crickets (N=5,728) with MUSCLE (v3.8.31). Then, the amino-acid alignments were transformed into codon-based nucleotide alignments using the

533 Pal2Nal software[92]. The resulting codon-based nucleotide alignments were used to calculate
534 the pairwise dN/dS for each gene pair with the yn00 algorithm implemented in the PAML
535 package[93]. Genes with dN or dS >2 were discarded from further analysis. The Wilcoxon-
536 Mann-Whitney statistical test was used to compare the dN/dS values between genes with
537 high and low $CpG_{o/e}$ values in both insects.

## Gene Family Expansions and Contractions

539 Using custom Python scripts (see
540 https://github.com/guillemylla/Crickets_Genome_Annotation) we obtained the longest
541 predicted protein product per gene in each of the 16 studied insect species and grouped them
542 into orthogroups (which we also refer to herein as "gene families") using OrthoFinder
543 v2.3.3[94]. The orthogroups (OGs) determined by OrthoFinder that contained a single gene per
544 insect, namely putative one-to-one orthologs, were used for phylogenetic reconstruction.
545 The proteins within each orthogroup were aligned with MUSCLE[95] and the alignments
546 trimmed with GBlocks (-t=p -b4=5 -b5=a)[96]. The trimmed alignments were concatenated
547 into a single meta-alignment that was used to infer the species tree with FastTree2
548 (FastTreeMP –gamma)[97].

549 To calibrate the species tree, we used the "chronos" function from the R package ape v5.3[98],
550 setting the common node between Blattodea and Orthoptera at 248 million years (my), the
551 origin of Holometabola at 345 my, the common node between Hemiptera and Thysanoptera
552 at 339 my, and the ancestor of hemimetabolous and holometabolous insects (root of the tree)
553 at between 385 and 395 my. These time points were obtained from a phylogeny published
554 that was calibrated with several fossils[23].

555 The gene family expansion/contraction analysis was done with the CAFE software[99]. We ran
556 CAFE using the calibrated species tree and the table generated by OrthoFinder with the
557 number of genes belonging to each orthogroup in each insect. Following the CAFE manual,
558 we first calculated the birth-death parameters with the orthogroups having less than 100
559 genes. We then corrected them by assembly quality and calculated the gene expansions and
560 contractions for both large (>100 genes) and small (≤100) gene families. This allowed us to
561 identify gene families that underwent a significant (p-value<0.01) gene family expansion or
562 contraction on each branch of the tree. We proceeded to obtain functional information from
563 those families expanded on our branches of interest (i.e. the origin of Orthoptera, the branch
564 leading to crickets, and the branches specific to each cricket species.). To functionally
565 annotate the orthogroups of interest, we first obtained the *D. melanogaster* identifiers of the
566 proteins within each orthogroup, and retrieved the FlyBase Symbol and the FlyBase gene
567 summary per gene using the FlyBase API[100]. Additionally, we ran InterProScan on all the
568 proteins of each orthogroup and retrieved all PFAM motifs and the GO terms together with
569 their descriptors. All of this information was summarized in tabulated files (**Supplementary
570 File 2**), which we used to identify gene expansions with potentially relevant functions for
571 insect evolution.

### *pickpocket* gene family expansion

Among the expanded gene families in crickets, we identified an orthogroup containing seven out of the eight *D. melanogaster pickpocket* class V genes, leading us to interpret that the *pickpocket* class V was significantly expanded in crickets. Subsequently, we retrieved the 6 additional orthogroups containing the completeset of *pickpocket* genes in *D. melanogaster*, and we assigned to each orthogroup the *pickpocket* class to which most of its *D. melanogaster* genes belonged according to Zelle and colleagues [24] (**Supplementary Table 2**). The protein sequences of all the members of the seven Pickpocket orthogroups were aligned with MUSCLE, and the *pickpocket* gene tree obtained with FastTree2 (FastTreeMP --gamma). The tips of the tree were colored based on the orthogroup to which they belong. A subset of the tree containing all the orthogroups that compose the entire *pickpocket* class V family was displayed as a circular cladogram (**Figure 3**), revealing an independent expansion of this family in *T. castaneum*.

To check for evidence of expression *pickpocket* genes in the cricket nervous system, we used the 21 RNA-seq libraries from prothoracic ganglion[9] of *G. bimaculatus* available at NCBI GEO (PRJNA376023). Reads were mapped against the *G. bimaculatus* genome with RSEM[101] using STAR[102] as the mapping algorithm, and the number of expected counts and TPMs were retrieved for each gene in each library. The TPMs of the *pickpocket* genes and *fruitless* are shown in **Supplementary Table 3**. Genes with a sum of more than 20 TPMs across all samples were considered to be expressed in *G. bimaculatus* prothoracic ganglion. We further analyzed the *pickpocket* expression in the aggregated embryo RNA-seq dataset (DRA011174) and normal and regenerating legs RNA-seq dataset[103] (DRR001985, DRR001986), using the same methodology.

## Acknowledgments

## Author contributions statement

GY, SN, TM and CE designed experiments; TI and AT conducted sequencing by HiSeq and assembling short reads using the Platanus assembler; ST, YI, TW, MF and YM performed DNA isolation, gap closing of contigs and manual annotation; GY, TN, ST, TB and AAB conducted all other experiments and analyses; TM and CE funded the project; GY and CE wrote the paper with input from all authors.

## Data availability

606

607 The genome sequencing reads, RNA-seq reads, and the genome assembly for *Gryllus*
608 *bimaculatus* were submitted to DDBJ and to NCBI under the accession number
609 (PRJDB10609). The genome assembly and annotations can also be accessed and browsed at
610 http://gbimaculatusgenome.rc.fas.harvard.edu.

## Code availability

611

612 The scripts used for genome annotation and analysis are available at GitHub
613 (https://github.com/guillemylla/Crickets_Genome_Annotation).

## Competing interests

614

615 The authors declare no competing interests.

## References

616

617

618 1    Belles, X. in *Encyclopedia of Life Sciences*   1-11 (John Wiley & Sons, Ltd, Chichester,
619      2011).
620 2    Engel, M. S. & Grimaldi, D. A. in *Nature* Vol. 427   627-630 (2004).
621 3    Gregory, T. R. Genome size and developmental complexity. *Genetica* **115**, 131-146,
622      doi:10.1023/a:1016032400147 (2002).
623 4    Camacho, J. P. *et al.* A step to the gigantic genome of the desert locust: chromosome
624      sizes and repeated DNAs. *Chromosoma* **124**, 263-275, doi:10.1007/s00412-014-
625      0499-0 (2015).
626 5    Bewick, A. J., Vogel, K. J., Moore, A. J. & Schmitz, R. J. in *Molecular Biology and*
627      *Evolution* Vol. 34   654-665 (2016).
628 6    Provataris, P., Meusemann, K., Niehuis, O., Grath, S. & Misof, B. Signatures of DNA
629      Methylation across Insects Suggest Reduced DNA Methylation Levels in
630      Holometabola. *Genome Biol Evol* **10**, 1185-1197, doi:10.1093/gbe/evy066 (2018).
631 7    Blankers, T., Oh, K. P., Bombarely, A. & Shaw, K. L. in *Genetics* Vol. 209   1329-1344
632      (Genetics, 2018).
633 8    Huber, F., Moore, T. E. T. E. & Loher, W. Cricket behavior and neurobiology. 565
634      (1989).
635 9    Fisher, H. P. *et al.* in *PLoS ONE* Vol. 13  (ed Silvana Allodi) e0199070 (Public Library
636      of Science, 2018).
637 10   Kainz, F., Ewen-Campen, B., Akam, M. & Extavour, C. G. Notch/Delta signalling is not
638      required for segment generation in the basally branching insect *Gryllus bimaculatus*.
639      *Development* **138**, 5015-5026, doi:10.1242/dev.073395 (2011).
640 11   Donoughe, S. & Extavour, C. G. in *Developmental Biology* Vol. 411   140-156
641      (Academic Press, 2015).
642 12   Mito, T. & Noji, S. in *CSH protocols* Vol. 2008   pdb.emo110 (Cold Spring Harbor
643      Laboratory Press, 2008).
644 13   Kulkarni, A. & Extavour, C. G. in *Evo-Devo: Non-model Species in Cell and*
645      *Developmental Biology* Vol. 68   183-216 (Springer, 2019).
646 14   Shinmyo, Y. *et al.* in *Development, Growth and Differentiation* Vol. 46   343-349 (John
647      Wiley & Sons, 2004).
648 15   Wilson Horch, H., Mito, T., Popadić, A., Ohuchi, H. & Noji, S. in *The Cricket as a Model*
649      *Organism: Development, Regeneration, and Behavior*   (eds Hadley Wilson Horch *et*
650      *al.*) 1-376 (Springer Japan, Tokyo, 2017).
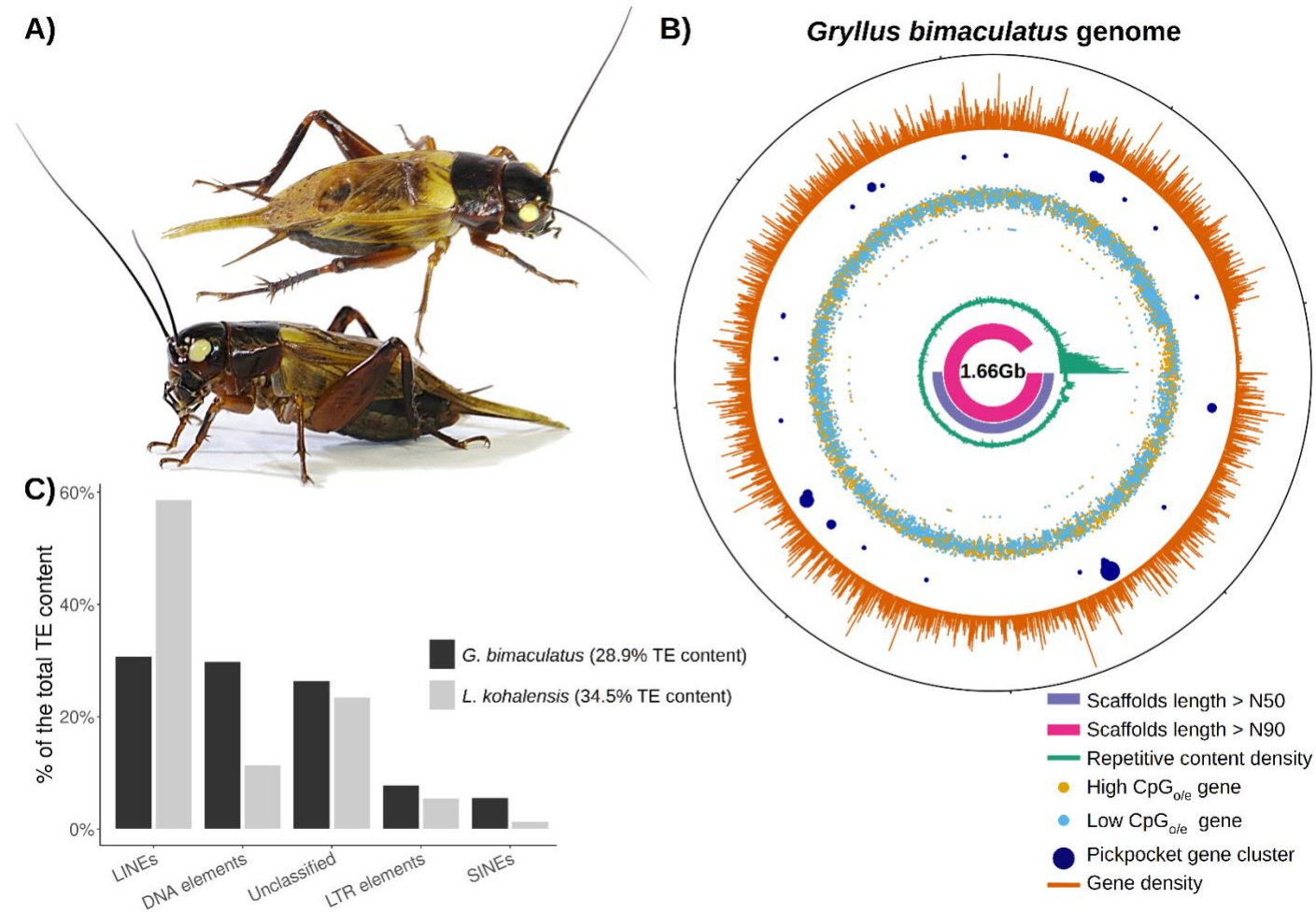651 16   Bird, A. P. in *Nucleic Acids Research* Vol. 8   1499-1504 (1980).

652  17    Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. in
653        *Bioinformatics* Vol. 31  3210-3212 (Oxford University Press, 2015).
654  18    Holt, C. & Yandell, M. in *BMC Bioinformatics* Vol. 12  491 (BioMed Central, 2011).
655  19    Jones, P. *et al.* in *Bioinformatics*  1-5 (2014).
656  20    Smit, A., Hubley, R. & Grenn, P. RepeatMasker Open-4.0. *RepeatMasker Open-4.0.7.*
657        (2015).
658  21    Wang, Y. *et al.* in *Science* Vol. 314  645-647 (2006).
659  22    Elango, N., Hunt, B. G., Goodisman, M. A. D. & Yi, S. V. in *Proceedings of the National*
660        *Academy of Sciences of the United States of America* Vol. 106  11206-11211 (National
661        Academy of Sciences, 2009).
662  23    Misof, B. *et al.* in *Science* Vol. 346  763-767 (2014).
663  24    Zelle, K. M., Lu, B., Pyfrom, S. C. & Ben-Shahar, Y. in *G3: Genes, Genomes, Genetics* Vol.
664        3  441-450 (Genetics Society of America, 2013).
665  25    Adams, C. M. *et al.* in *The Journal of Cell Biology* Vol. 140  143-152 (Rockefeller
666        University Press, 1998).
667  26    Zhong, L., Hwang, R. Y. & Tracey, W. D. in *Current Biology* Vol. 20  429-434 (2010).
668  27    Ainsley, J. A. *et al.* in *Current Biology* Vol. 13  1557-1563 (Cell Press, 2003).
669  28    Häsemeyer, M., Yapici, N., Heberlein, U. & Dickson, B. J. in *Neuron* Vol. 61  511-518
670        (Cell Press, 2009).
671  29    Rezával, C. *et al.* in *Current Biology* Vol. 22  1155-1165 (Elsevier, 2012).
672  30    Pavlou, H. J. & Goodwin, S. F. in *Current opinion in neurobiology* Vol. 23  76-83
673        (Elsevier, 2013).
674  31    Adams, M. D. *et al.* in *Science* Vol. 287  2185-2195 (2000).
675  32    Elsik, C. G. The pea aphid genome sequence brings theories of insect defense into
676        question. *Genome Biol* **11**, 106, doi:10.1186/gb-2010-11-2-106 (2010).
677  33    Wang, X. *et al.* in *Nature communications* Vol. 5  2957 (Nature Publishing Group,
678        2014).
679  34    Harrison, M. C. *et al.* Hemimetabolous genomes reveal molecular basis of termite
680        eusociality. *Nat Ecol Evol* **2**, 557-566, doi:10.1038/s41559-017-0459-1 (2018).
681  35    Terrapon, N. *et al.* in *Nat Commun* Vol. 5  3636 (Nature Publishing Group, 2014).
682  36    Poulsen, M. *et al.* in *Proceedings of the National Academy of Sciences of the United*
683        *States of America* Vol. 111  14500-14505 (National Academy of Sciences, 2014).
684  37    Kouřimská, L. & Adámková, A. in *NFS Journal* Vol. 4  22-26 (Elsevier, 2016).
685  38    Van Huis, A. *et al. Edible insects: future prospects for food and feed security*.  (Food
686        and agriculture organization of the United Nations (FAO), 2013).
687  39    Ghosh, S., Lee, S.-M., Jung, C. & Meyer-Rochow, V. Nutritional composition of five
688        commercial edible insects in South Korea. *Journal of Asia-Pacific Entomology* **20**,
689        686-694 (2017).
690  40    Hanboonsong, Y., Jamjanya, T. & Durst, P. B. in *Office*  69 (Food and Agriculture
691        Organization of the United Nations, Regional Office for Asia and the Pacific, 2013).
692  41    Ryu, H. Y. *et al.* in *Toxicological Research* Vol. 32  159-173 (Korean Society of
693        Toxicology, 2016).
694  42    Ahn, M. Y., Han, J. W., Kim, S. J., Hwang, J. S. & Yun, E. Y. in *Toxicological Research* Vol.
695        27  231-240 (2011).
696  43    Mi, Y. A. *et al.* in *Journal of Toxicology and Environmental Health - Part A* Vol. 68
697        2111-2118 (2005).
698  44    Pener, M. P. in *Journal of Orthoptera Research* Vol. 25  91-95 (2016).
699  45    Ribeiro, J. C., Cunha, L. M., Sousa-Pinto, B. & Fonseca, J. in *Molecular Nutrition and*
700        *Food Research* Vol. 62  1700030 (2018).
701  46    Ahn, M. Y., Hwang, J. S., Yun, E. Y., Kim, M. J. & Park, K. K. in *Toxicological Research*
702        Vol. 31  173-180 (Korean Society of Toxicology, 2015).
703  47    Park, S. A., Lee, G. H., Lee, H. Y., Hoang, T. H. & Chae, H. J. in *Food Science and*
704        *Nutrition* Vol. 8  402-409 (2019).
705  48    Hwang, B. B. *et al.* in *Nutrients* Vol. 11  857 (2019).
706  49    Thrall, P. H., Bever, J. D. & Burdon, J. J. Evolutionary change in agriculture: the past,
707        present and future. *Evol Appl* **3**, 405-408, doi:10.1111/j.1752-4571.2010.00155.x
708        (2010).
709  50    Yamasaki, M. *et al.* A large-scale screen for artificial selection in maize identifies
710        candidate agronomic loci for domestication and crop improvement. *Plant Cell* **17**,
711        2859-2872, doi:10.1105/tpc.105.037242 (2005).

712   51   Chen, Y. H., Gols, R. & Benrey, B. Crop domestication and its impact on naturally
713        selected trophic interactions. *Annu Rev Entomol* **60**, 35-58, doi:10.1146/annurev-
714        ento-010814-020601 (2015).
715   52   Gepts, P. Crop domestication as a long-term selection experiment. *Plant breeding*
716        *reviews* **24**, 1-44 (2004).
717   53   Chen, K. & Gao, C. Targeted genome modification technologies and their applications
718        in crop improvements. *Plant Cell Reports* **33**, 575-583, doi:10.1007/s00299-013-
719        1539-6 (2014).
720   54   Qaim, M. The Economics of Genetically Modified Crops. *Annual Review of Resource*
721        *Economics* **1**, 665-694, doi:10.1146/annurev.resource.050708.144203 (2009).
722   55   Harrison, M. C. *et al.* Hemimetabolous genomes reveal molecular basis of termite
723        eusociality. *Nature ecology & evolution*, doi:10.1038/s41559-017-0459-1 (2018).
724   56   Kidwell, M. G. in *Genetica* Vol. 115   49-63 (Kluwer Academic Publishers, 2002).
725   57   Chénais, B., Caruso, A., Hiard, S. & Casse, N. in *Gene* Vol. 509   7-15 (Elsevier, 2012).
726   58   Hanrahan, S. J. & Johnston, J. S. in *Chromosome research : an international journal on*
727        *the molecular, supramolecular and evolutionary aspects of chromosome biology* Vol.
728        19   809-823 (2011).
729   59   Westerman, M., Barton, N. & Hewitt, G. M. Differences in DNA content between two
730        chromosomal races of the grasshopper Podisma pedestris. *Heredity* **58**, 221-228
731        (1987).
732   60   Palacios-Gimenez, O. M. *et al.* Comparative analysis of morabine grasshopper
733        genomes reveals highly abundant transposable elements and rapidly proliferating
734        satellite DNA repeats. *BMC Biology* **18**, doi:10.1186/s12915-020-00925-x (2020).
735   61   Lyko, F., Ramsahoye, B. H. & Jaenisch, R. DNA methylation in Drosophila
736        melanogaster. *Nature* **408**, 538-540, doi:10.1038/35046205 (2000).
737   62   Thomas, G. W. C. *et al.* Gene content evolution in the arthropods. *Genome Biol* **21**, 15,
738        doi:10.1186/s13059-019-1925-7 (2020).
739   63   Lewis, S. H. *et al.* Widespread conservation and lineage-specific diversification of
740        genome-wide DNA methylation patterns across arthropods. *PLoS Genet* **16**,
741        e1008864, doi:10.1371/journal.pgen.1008864 (2020).
742   64   Ylla, G., Piulachs, M.-D. & Belles, X. in *iScience* Vol. 4   164-179 (Elsevier, 2018).
743   65   Ventós-Alfonso, A., Ylla, G., Montañes, J.-C. & Belles, X. DNMT1 promotes genome
744        methylation and early embryo development in cockroaches. *iScience*, 101778,
745        doi:10.1016/j.isci.2020.101778 (2020).
746   66   Liu, L., Johnson, W. A. & Welsh, M. J. in *Proceedings of the National Academy of*
747        *Sciences of the United States of America* Vol. 100   2128-2133 (National Academy of
748        Sciences, 2003).
749   67   Lee, M. J. *et al.* in *Molecules and cells* Vol. 40   787-795 (Korean Society for Molecular
750        and Cellular Biology, 2017).
751   68   Averhoff, W. W., Richardson, R. H., Starostina, E., Kinser, R. D. & Pikielny, C. W. in
752        *Proceedings of the National Academy of Sciences of the United States of America* Vol.
753        73   591-593 (National Academy of Sciences, 1976).
754   69   Lu, B., LaMora, A., Sun, Y., Welsh, M. J. & Ben-Shahar, Y. in *PLoS Genetics* Vol. 8   (ed
755        Miriam B. Goodman) e1002587 (Public Library of Science, 2012).
756   70   Jacob, P. F. & Hedwig, B. in *Behavioural Brain Research* Vol. 309   51-66 (Elsevier,
757        2016).
758   71   Blankers, T., Oh, K. P. & Shaw, K. L. in *Genes* Vol. 9   346 (Multidisciplinary Digital
759        Publishing Institute, 2018).
760   72   Shaw, K. L. & Lesnick, S. C. in *Proceedings of the National Academy of Sciences* Vol.
761        106   9737-9742 (National Academy of Sciences, 2009).
762   73   Xu, M. & Shaw, K. L. in *Genetics* Vol. 211   1089-1104 (2019).
763   74   Yoshimura, A., Nakata, A., Mito, T. & Noji, S. in *Cytogenetic and Genome Research* Vol.
764        112   329-336 (2006).
765   75   Kajitani, R. *et al.* Efficient de novo assembly of highly heterozygous genomes from
766        whole-genome shotgun short reads. *Genome Res* **24**, 1384-1395,
767        doi:10.1101/gr.170720.113 (2014).
768   76   English, A. C. *et al.* Mind the gap: upgrading genomes with Pacific Biosciences RS
769        long-read sequencing technology. *PLoS One* **7**, e47768,
770        doi:10.1371/journal.pone.0047768 (2012).
771   77   Ellinghaus, D., Kurtz, S. & Willhoeft, U. in *BMC Bioinformatics* Vol. 9   18 (BioMed
772        Central, 2008).

773  78  Crescente, J. M., Zavallo, D., Helguera, M. & Vanzetti, L. S. in *BMC Bioinformatics* Vol.
774      19   348 (BioMed Central, 2018).
775  79  Vassetzky, N. S. & Kramerov, D. A. in *Nucleic acids research* Vol. 41   D83-89 (Oxford
776      University Press, 2013).
777  80  Bao, W., Kojima, K. K. & Kohany, O. in *Mobile DNA* Vol. 6   11 (BioMed Central, 2015).
778  81  Edgar, R. C. in *Bioinformatics* Vol. 26   2460-2461 (2010).
779  82  Zeng, V. *et al.* in *PLoS ONE* Vol. 8   (ed Peter K. Dearden) e61479 (Public Library of
780      Science, 2013).
781  83  Pertea, M. *et al.* in *Nature Biotechnology* Vol. 33   290-295 (Nature Publishing Group,
782      2015).
783  84  Kim, D., Langmead, B. & Salzberg, S. L. HISAT: A fast spliced aligner with low
784      memory requirements. *Nature Methods*, doi:10.1038/nmeth.3317 (2015).
785  85  UniProt, C. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* **47**,
786      D506-D515, doi:10.1093/nar/gky1049 (2019).
787  86  Ter-Hovhannisyan, V., Lomsadze, A., Chernoff, Y. O. & Borodovsky, M. in *Genome
788      Research* Vol. 18   1979-1990 (Cold Spring Harbor Laboratory Press, 2008).
789  87  Stanke, M. & Waack, S. in *Bioinformatics* Vol. 19   ii215-ii225 (2003).
790  88  Korf, I. in *BMC Bioinformatics* Vol. 5   59 (BioMed Central, 2004).
791  89  Benaglia, T., Chauveau, D., Hunter, D. R. & Young, D. S. in *Journal of Statistical
792      Software* Vol. 32   1-29 (2009).
793  90  Alexa, A. & Rahnenfuhrer, J. topGO: Enrichment Analysis for Gene Ontology. *R
794      package version 2.36.0* (2019).
795  91  Lex, A., Gehlenborg, N., Strobelt, H., Vuillemot, R. & Pfister, H. in *IEEE Transactions on
796      Visualization and Computer Graphics* Vol. 20   1983-1992 (Europe PMC Funders,
797      2014).
798  92  Suyama, M., Torrents, D. & Bork, P. in *Nucleic Acids Research* Vol. 34   W609-W612
799      (2006).
800  93  Yang, Z. in *Molecular Biology and Evolution* Vol. 24   1586-1591 (2007).
801  94  Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for
802      comparative genomics. *Genome Biol* **20**, 238, doi:10.1186/s13059-019-1832-y
803      (2019).
804  95  Edgar, R. C. in *Nucleic Acids Research* Vol. 32   1792-1797 (Oxford University Press,
805      2004).
806  96  Castresana, J. in *Molecular Biology and Evolution* Vol. 17   540-552 (2000).
807  97  Price, M. N., Dehal, P. S. & Arkin, A. P. in *PLoS ONE* Vol. 5   (ed Art F. Y. Poon) e9490
808      (Public Library of Science, 2010).
809  98  Paradis, E. & Schliep, K. in *Bioinformatics* Vol. 35   (ed Russell Schwartz) 526-528
810      (2019).
811  99  De Bie, T., Cristianini, N., Demuth, J. P. & Hahn, M. W. in *Bioinformatics* Vol. 22   1269-
812      1271 (2006).
813  100 Thurmond, J. *et al.* FlyBase 2.0: the next generation. *Nucleic Acids Res* **47**, D759-
814      D765, doi:10.1093/nar/gky1003 (2019).
815  101 Li, B. & Dewey, C. N. in *BMC Bioinformatics* Vol. 12   323 (2011).
816  102 Dobin, A. *et al.* in *Bioinformatics* Vol. 29   15-21 (2013).
817  103 Bando, T. *et al.* Analysis of RNA-Seq data reveals involvement of JAK/STAT signalling
818      during leg regeneration in the cricket Gryllus bimaculatus. *Development* **140**, 959-
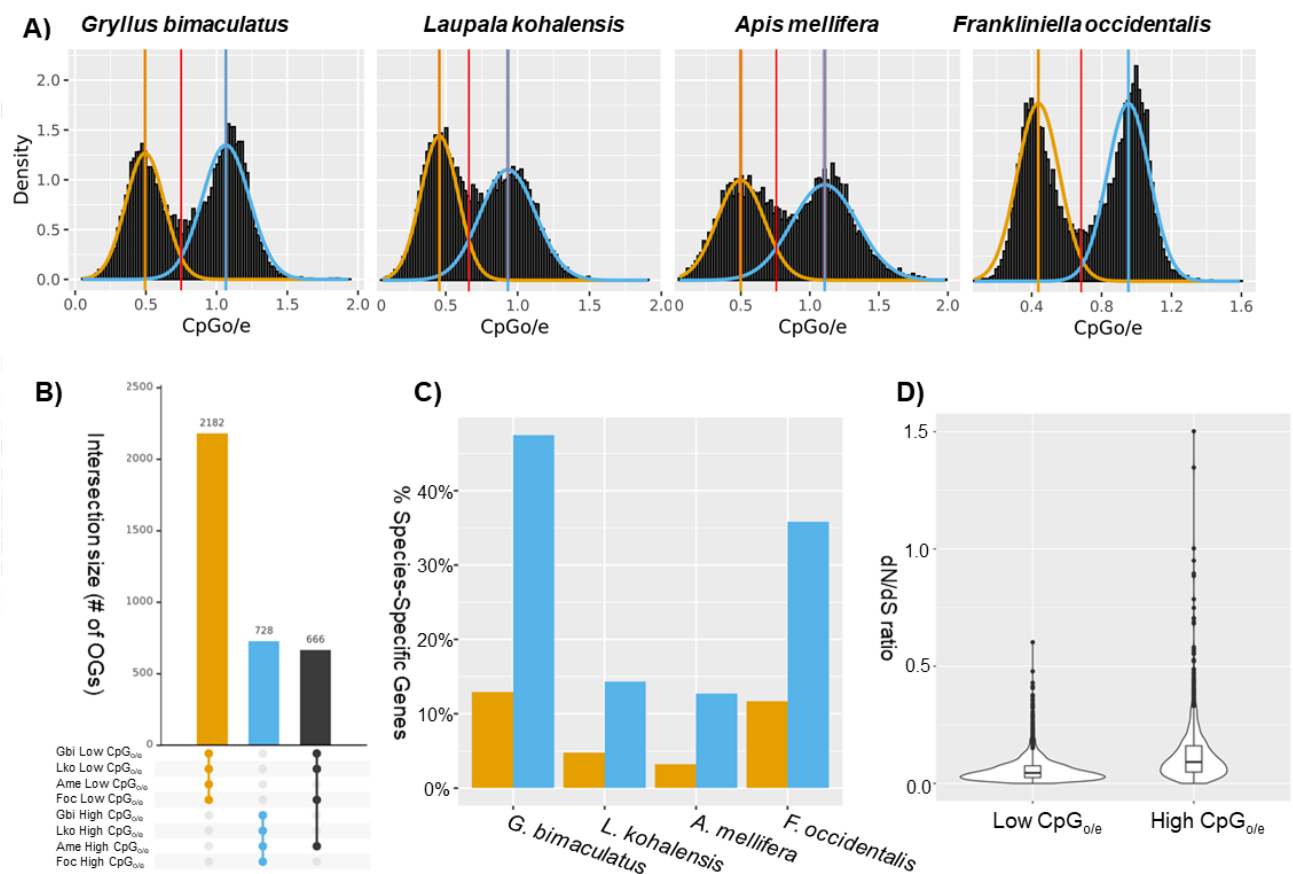819      964, doi:10.1242/dev.084590 (2013).
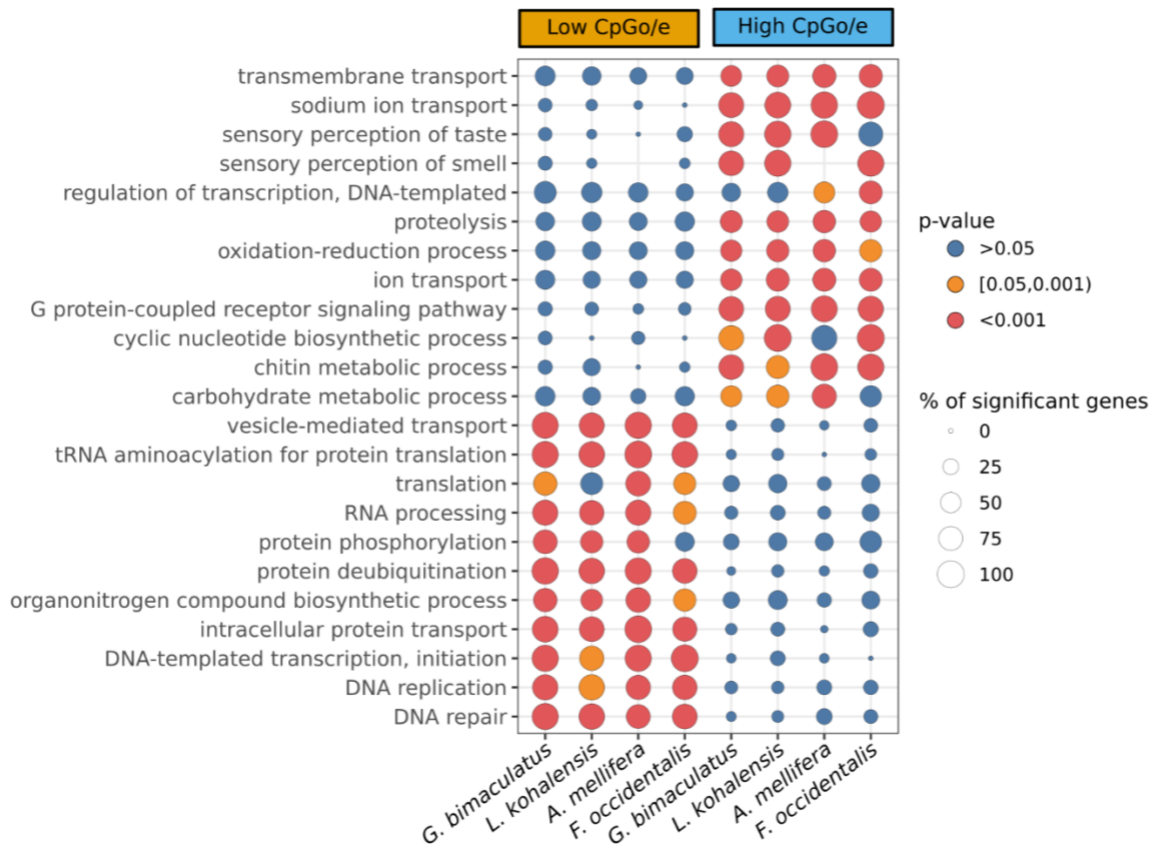
820

**Figure 1: The *G. bimaculatus* genome. A)** The cricket *G. bimaculatus* (top and side views of an adult male), commonly called the two-spotted cricket, owes its name to the two yellow spots on the base of the forewings. **B)** Circular representation of the *G. bimaculatus* genome, displaying the N50 (pink) and N90 (purple) scaffolds, repetitive content density (green), the high- (yellow) and low- (light blue) CpG$_{o/e}$ value genes, *pickpocket* gene clusters (dark blue), and gene density (orange). **C)** The proportion of the genome made up of transposable elements (TEs) is similar between *G. bimaculatus* and *L. kohalensis* (28.9% and 34.5% respectively), but the specific TE family composition varies widely.
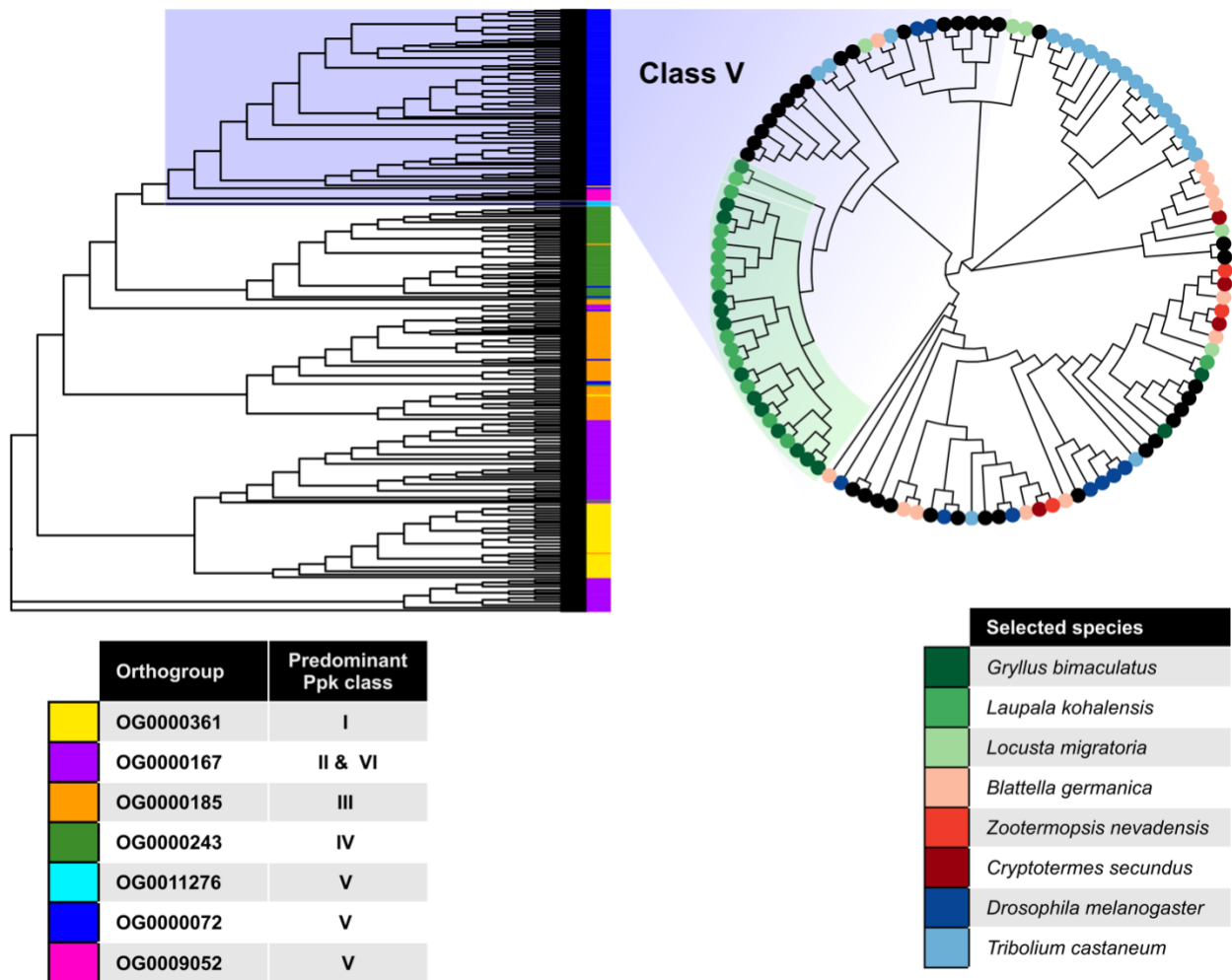
**Figure 2: CpG$_{o/e}$ distribution across insects**. **A)** The distribution of CpG$_{o/e}$ values within the CDS regions displays a bimodal distribution in the two crickets, as well as in the honeybee *A. mellifera* and the thrips *F. occidentalis*. We modeled each peak with a normal distribution and defined their intersection (red line) as a threshold to separate genes into low- and high- CpG$_{o/e}$ value categories represented in yellow and blue color respectively. **B)** UpSet plot showing the top three intersections (linked dots) in terms of the number of orthogroups (OGs) commonly present in the same category (low- and high- CpG$_{o/e}$) across the four insect species. The largest intersection corresponds to 2,182 OGs whose genes have low CpG$_{o/e}$ in the four insect species, followed by the 728 OGs whose genes have high CpG$_{o/e}$ levels in all four species, and 666 OGs whose genes have low CpG$_{o/e}$ in the three hemimetabolous species and high CpG$_{o/e}$ in the holometabolous honeybee. Extended plot with 50 intersections is shown in **Supplementary Figure 4**. **C)** Percentage of species-specific genes within low CpG$_{o/e}$ (yellow) and high CpG$_{o/e}$.(blue) in each insect, indicating that more such genes tend to have high CpG$_{o/e}$ values. **D)** One-to-one orthologous genes with low CpG$_{o/e}$ values in both crickets have significantly lower dN/dS values than genes with high CpG$_{o/e}$ values.

**Figure 3: Functional analysis of high- and low- CpG$_{o/e}$ genes:** Enriched GO terms with a p-value<0.00001 in at least one of the eight categories, which are high CpG$_{o/e}$ and low CpG$_{o/e}$ genes of *G. bimaculatus*, *L. kohalensis*, *F. occidentalis*, and *A. mellifera*. The dot diameter is proportional to the percentage of significant genes with the GO term within the gene set. The dot color represents the p-value level: blue >0.05, orange [0.05, 0.001), red <0.001. Extended figure with all significant GO terms (p-value<0.05) available as **Supplementary Figure 3**.

| Orthogroup | Predominant Ppk class |
|---|---|
| OG0000361 | I |
| OG0000167 | II & VI |
| OG0000185 | III |
| OG0000243 | IV |
| OG0011276 | V |
| OG0000072 | V |
| OG0009052 | V |

Selected species

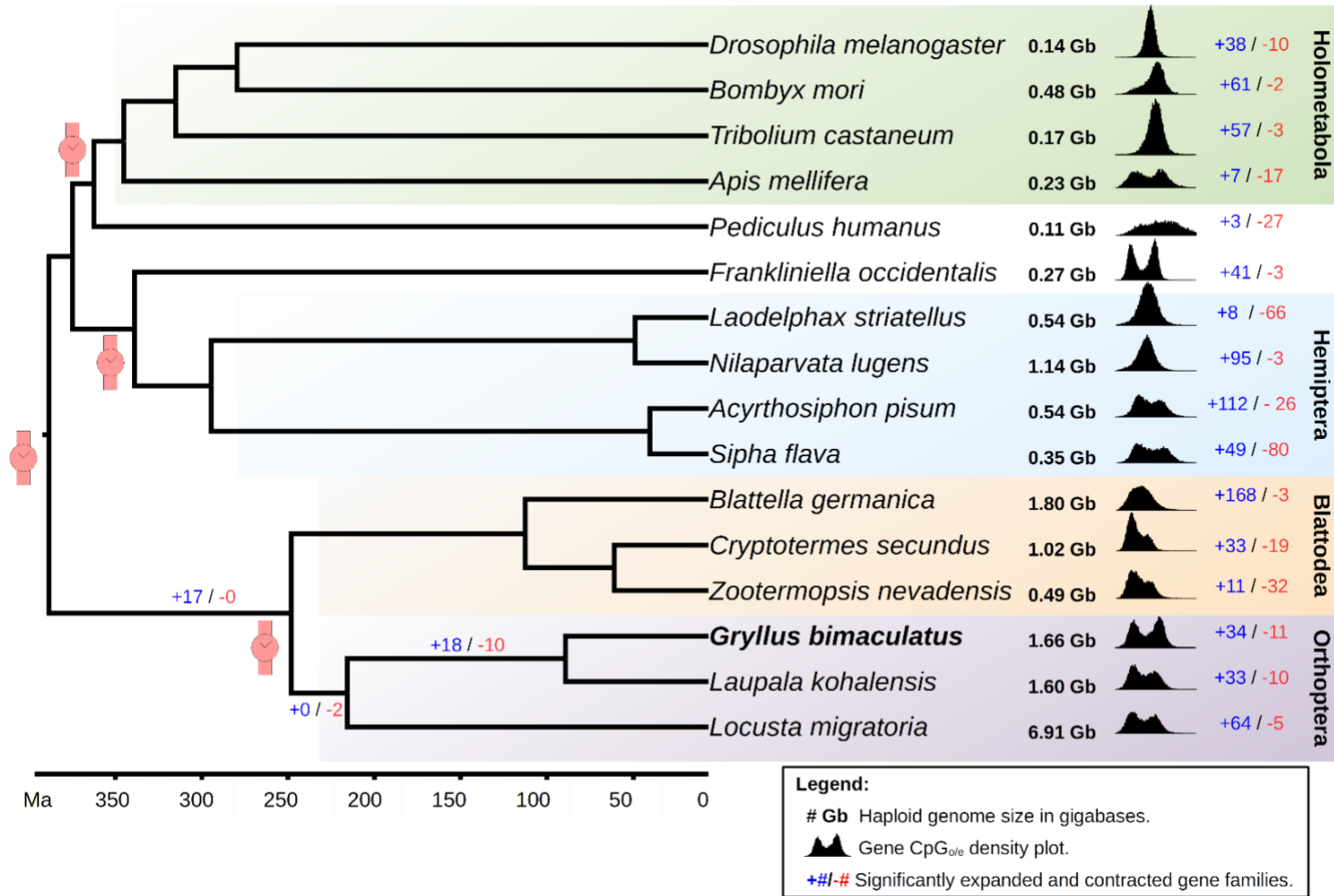| | |
|---|---|
| | *Gryllus bimaculatus* |
| | *Laupala kohalensis* |
| | *Locusta migratoria* |
| | *Blattella germanica* |
| | *Zootermopsis nevadensis* |
| | *Cryptotermes secundus* |
| | *Drosophila melanogaster* |
| | *Tribolium castaneum* |

854

855 **Figure 4: The *pickpocket* gene family class V is expanded in crickets**. *pickpocket* gene
856 tree with all the genes belonging to the seven OGs that contain the *D. melanogaster*
857 *pickpocket* genes. All OGs predominantly contain members of a single *ppk* family. The
858 OG0000167contains members of two *pickpocket* classes, II and VI. The orthogroup
859 OG0000072 containing most *pickpocket* class V genes (circular cladogram) was
860 significantly expanded in crickets relative to other insects.

**Figure 5: Cricket genomes in the context of insect evolution.** A phylogenetic tree including 16 insect species calibrated at four different time points (red watch symbols) based on Misof, et al. [23], suggests that *G. bimaculatus* and *L. kohalensis* diverged ca. 89.2 Mya. The number of expanded (blue text) and contracted (red text) gene families is shown for each insect, and for the branches leading to crickets. The density plots show the $CpG_{o/e}$ distribution for all genes for each species. The genome size in Gb was obtained from the genome fasta files (**Supplementary Table 1**).

## Supplementary Materials

**Supplementary Materials for**

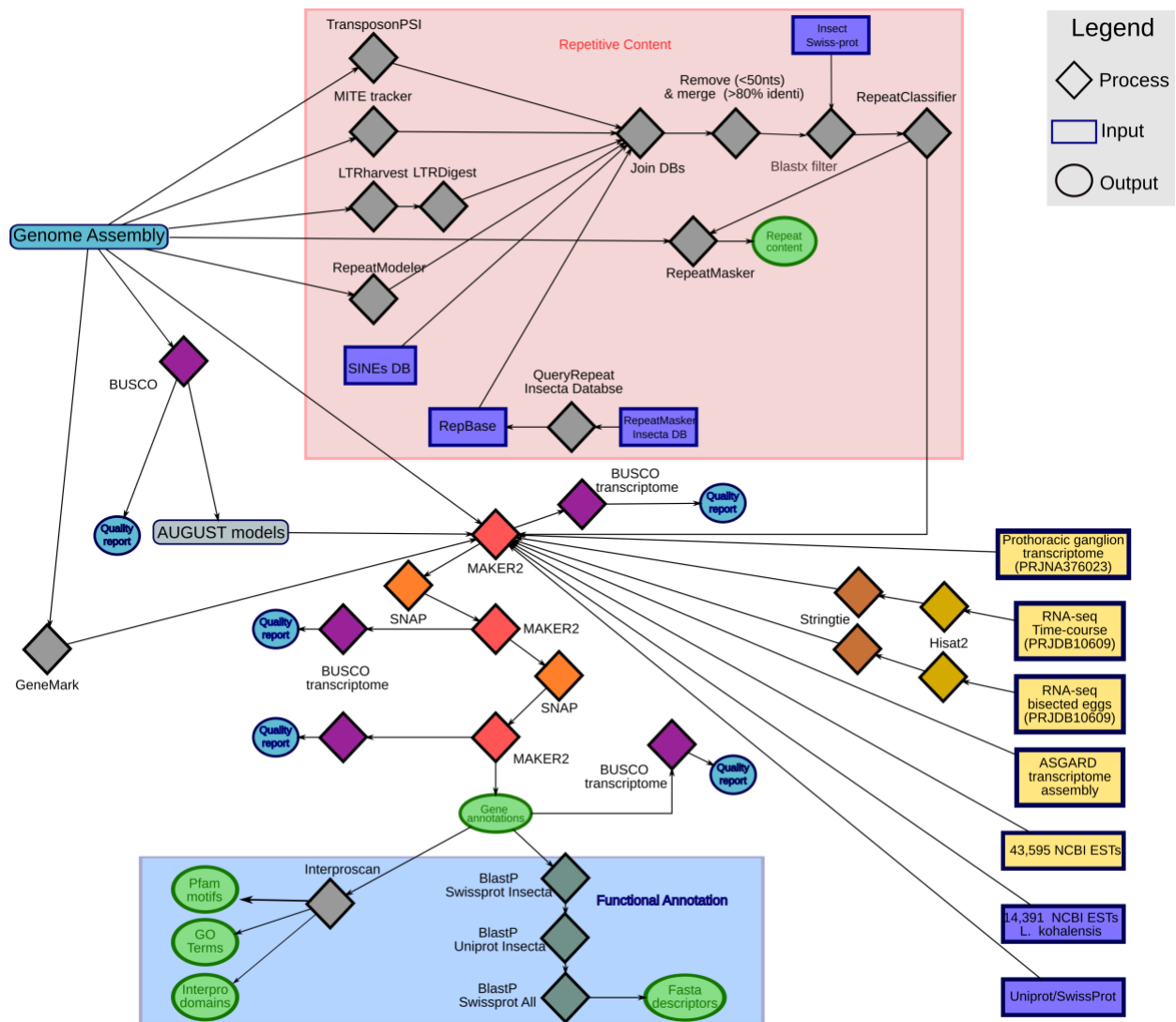# Insights into the genomic evolution of insect from Cricket genomes

Guillem Ylla, Taro Nakamura, Takehiko Itoh, Rei Kajitani, Atsushi Toyoda, Sayuri Tomonari, Tetsuya Bando, Yoshiyasu Ishimaru, Takahito Watanabe, Masao Fuketa, Yuji Matsuoka, Austen A. Barnett, Sumihare Noji, Taro Mito, Cassandra G. Extavour

These Supplementary Materials consist of the following:

- Supplementary Figures 1 – 4 (this document)
- Supplementary File 1 (this document)
- Supplementary File 2 ("Supplementary_File_2.xls")
- Supplementary Table 1 ("Supplementary_Table_1.xls")
- Supplementary Table 2 (this document)
- Supplementary Table 3 and 4 ("Supplementary_Table_3-4.xls")
- Supplementary Table 5 (this document)
- Supplementary References (this document)

885

886

887 **Supplementary Figure 1: Schematic of *G. bimaculatus* genome annotation pipeline**.
888 Rectangles represent data inputs: yellow rectangles represent *G. bimaculatus* data; purple
889 rectangles represent data from other species or databases. Diamonds represent computational
890 processes: gray diamonds indicate processes executed a single time; non-gray diamonds of the
891 same color indicate the same process. Circles indicate outputs: blue circles indicate quality
892 controls; green circles indicate annotations. Scripts available at GitHub
893 https://github.com/guillemylla/Crickets_Genome_Annotation.

**Supplementary Figure 2: Scheme of *L. kohalensis* genome annotation pipeline**. All symbols as per **Supplementary Figure 1**.

**Supplementary Figure 3: Enriched GO-terms among genes with high or low CpGo/e levels**. This plot shows the enriched GO terms with p-value<0.05 in at least one of the eight categories which are the high CpGo/e and low CpGo/e genes of *G. bimaculatus* (Gbi), *L. kohalensis* (Lko), *F. occidentalis* (Foc), and *A. mellifera* (Ame). The dot diameter is proportional the percentage of significant genes with the GO term within the gene set. The dot color represents the p-value level: blue >0.05, orange [0.05, 0.001), red <0.001.

928



929

**Supplementary Figure 4: UpSet plot of orthologous genes within the high and low CpG$_{o/e}$ value categories.** Top 50 intersections of orthogroups (OGs) that are common across the eight different categories, which are the high CpG$_{o/e}$ and low CpG$_{o/e}$ genes for *G. bimaculatus* (Gbi), *L. kohalensis* (Lko), *F. Occidentalis* (Foc), and *A. mellifera* (Ame). Blue color indicates OGs that contain genes that only belong to the high CpG$_{o/e}$ peak, and yellow OGs contain genes that only belong to the low CpG$_{o/e}$ peak.

934 **Supplementary File 1**: **RepeatMasker summaries.** Report of the repeat content in the
935 genomes of *G. bimaculatus* and *L. kohalensis* generated by RepeatMasker using custom
936 libraries.

## *Gryllus bimaculatus*

```
==================================================
file name: Gbimaculatus_Gap_filled.fasta
sequences:          47877
total length: 1658007496 bp  (1601517380 bp excl N/X-runs)
GC level:          39.93 %
bases masked:   558652201 bp ( 33.69 %)
==================================================
              number of      length    percentage
              elements*    occupied   of sequence
--------------------------------------------------
SINEs:           138895    26406967 bp    1.59 %
      ALUs              6        9564 bp    0.00 %
      MIRs              0           0 bp    0.00 %

LINEs:           454301   147302087 bp    8.88 %
      LINE1          1803      826764 bp    0.05 %
      LINE2        115576    32029561 bp    1.93 %
      L3/CR1        18286     6358119 bp    0.38 %

LTR elements:    131656    36970251 bp    2.23 %
      ERVL             92       44183 bp    0.00 %
      ERVL-MaLRs        0           0 bp    0.00 %
      ERV_classI    11451     2441461 bp    0.15 %
      ERV_classII     980      401749 bp    0.02 %

DNA elements:    500741   142828465 bp    8.61 %
      hAT-Charlie   11512     4094376 bp    0.25 %
      TcMar-Tigger   2039      537995 bp    0.03 %

Unclassified:    367653   126552078 bp    7.63 %

Total interspersed repeats:480059848 bp    28.95 %


Small RNA:         2562     1002728 bp    0.06 %


Satellites:       31087     7528498 bp    0.45 %
Simple repeats:  769175    77632578 bp    4.68 %
Low complexity:   85129     6215377 bp    0.37 %
==================================================
```

## *Laupala kohalensis*

```
==================================================
file name: GCA_002313205.1_ASM231320v1_genomic.fna
sequences:          148784
total length: 1595214429 bp  (1563778341 bp excl N/X-runs)
GC level:          35.58 %
```

937

```
bases masked:  566518287 bp ( 35.51 %)
==================================================
                number of      length   percentage
                elements*     occupied  of sequence
--------------------------------------------------
SINEs:             29510     7083717 bp    0.44 %
        ALUs         304      101257 bp    0.01 %
        MIRs        1248      430584 bp    0.03 %

LINEs:           1035151   322470849 bp   20.21 %
        LINE1        941      367057 bp    0.02 %
        LINE2     584526   167380843 bp   10.49 %
        L3/CR1     10257     4624100 bp    0.29 %

LTR elements:      57347    29690552 bp    1.86 %
        ERVL         231       43500 bp    0.00 %
        ERVL-MaLRs     0           0 bp    0.00 %
        ERV_classI  1821      585650 bp    0.04 %
        ERV_classII  389      125302 bp    0.01 %

DNA elements:     189815    62384975 bp    3.91 %
      hAT-Charlie  15008     5154516 bp    0.32 %
      TcMar-Tigger  8896     2459752 bp    0.15 %

Unclassified:     409303   128822550 bp    8.08 %

Total interspersed repeats:550452643 bp   34.51 %


Small RNA:         13816     3005585 bp    0.19 %

Satellites:         2088      882748 bp    0.06 %
Simple repeats:   307925    19782955 bp    1.24 %
Low complexity:    48386     2381730 bp    0.15 %
==================================================
```

938

939

940

941 **Supplementary File 2: Gene family expansions in crickets.** Gene families (Orthogroups)
942 significantly expanded in the lineage leading to crickets (tab 1), expanded in *G. bimaculatus*
943 (tab 2), and expanded in *L. kohalensis* (tab 3). For each expanded orthogroup (OG), we
944 report the expansion size as the number of genes gained, and the functional information
945 about the OG. The functional information consists of the list of PFAMs and GO terms
946 associated with the genes within the OG, and the list of *D. melanogaster* genes within the OG
947 with their FlyBase summaries.

948

949 *See file "Supplementary_File_2_GeneExpansions.xls"*

950

951 **Supplementary Table 1: Genome assembly information for the 16 insect genomes**
952 **analyzed.** For each genome, we show the database that the assembly was retrieved from,
953 the assembly file name, the accession code, the assembly statistics obtained with assembly-
954 stats software (https://github.com/sanger-pathogens/assembly-stats) and the BUSCO
955 v3.1.0 reports at Arthropoda and Insecta levels.

956

957 *See file "Supplementary_Table_1_GenomeStats.xls"*

958

959

960 **Supplementary Table 2**: The orthogroups (OG) containing the 31 *D. melanogaster*
961 *pickpocket* genes, with their FlyBase ID, symbol, and class according to Zelle et al. (2013).

962

| OG | FlyBase ID | gene symbol | Zelle et al. (2013) class |
|---|---|---|---|
| OG0000361.fa | FBgn0034965 | *ppk29* | I |
| OG0000361.fa | FBgn0039424 | *ppk15* | I |
| OG0000361.fa | FBgn0051065 | *ppk31* | I |
| OG0000361.fa | FBgn0053508 | *ppk13* | I |
| OG0009052.fa | FBgn0032602 | *ppk17* | V |
| OG0000185.fa | FBgn0039675 | *ppk21* | III |
| OG0000185.fa | FBgn0039677 | *ppk30* | III |
| OG0000185.fa | FBgn0039679 | *ppk19* | III |
| OG0000185.fa | FBgn0065109 | *ppk11* | IV |
| OG0000185.fa | FBgn0039676 | *ppk20* | III |
| OG0000185.fa | FBgn0031802 | *ppk7* | III |
| OG0000185.fa | FBgn0031803 | *ppk14* | III |
| OG0000072.fa | FBgn0022981 | *rpk / ppk2* | V |
| OG0000072.fa | FBgn0034730 | *ppk12* | V |
| OG0000072.fa | FBgn0052792 | *ppk8* | V |
| OG0000072.fa | FBgn0053289 | *ppk5* | V |
| OG0000072.fa | FBgn0020258 | *ppk / ppk1* | V |
| OG0000072.fa | FBgn0265001 | *ppk18* | IV |
| OG0000072.fa | FBgn0030795 | *ppk28* | V |
| OG0000072.fa | FBgn0035785 | *ppk26* | V |
| OG0011276.fa | FBgn0035458 | *ppk27* | IV |
| OG0000243.fa | FBgn0034489 | *ppk6* | IV |
| OG0000243.fa | FBgn0039839 | *ppk24* | IV |
| OG0000243.fa | FBgn0051105 | *ppk22* | IV |
| OG0000243.fa | FBgn0065108 | *ppk16* | IV |
| OG0000243.fa | FBgn0024319 | *Nach / ppk4* | IV |
| OG0000167.fa | FBgn0050181 | *ppk3* | II |
| OG0000167.fa | FBgn0053349 | *ppk25* | II |
| OG0000167.fa | FBgn0065110 | *ppk10* | II |
| OG0000167.fa | FBgn0085398 | *ppk9* | II |
| OG0000167.fa | FBgn0030844 | *ppk23* | VI |

963

964

965 **Supplementary Table 3: *pickpocket* gene expression levels in the *G. bimaculatus***
966 **prothoracic ganglion.** Expression in TPMs of *fruitless* and *pickpocket* genes in each RNA-
967 seq library generated from adult male prothoracic ganglia previously generated by Fisher
968 and colleagues (2018). Genes with read sum across samples > 20 TPMs across samples are
969 highlighted.

970

971 **Supplementary Table 4: *pickpocket* gene expression levels in the *G. bimaculatus***
972 **embryo and regenerating legs.** Expression in TPMs of *fruitless* and *pickpocket* genes in
973 the aggregated embryo RNA-seq dataset, control legs and regenerating legs. Genes with
974 read sum across samples > 20 TPMs across samples in the prothoracic ganglion
975 (Supplementary Table 3) are highlighted.

976

977 See file "Supplementary_Table_3-4.xls"

978 **Supplementary Table 5: *pickpocket* genes present in previous QTL analyses examining the genetic basis for sound-based cricket courtship behavior**
979 **variation.** Genomic position information for the *L. kohalensis pickpocket* genes found in linkage groups (LG) in previously published QTL analyses (Blankers, Oh & Shaw
980 2018; Shaw & Lesnick, 2009) examining mating song rhythm variations and female acoustic preference in the genus *Laupala*.

981

| Scaff names Shaw | Scaff Names NCBI | start | end | width | strand | Name | Ppk class | Table S3 and S6 (Blankers, Oh, & Shaw, 2018) | | Table S4 (Blankers, Oh, Bombarely, & Shaw, 2018) | Table 2 (Xu and Shaw, 2019) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | LG | proximity | LG | LG |
| Lko057S000409 | NNCF01126148.1 | 1083057 | 1116038 | 32982 | + | Lko_01144 | Class IV | 1 | LOD1 | 1 | |
| Lko057S000550 | NNCF01126289.1 | 666338 | 667949 | 1612 | - | Lko_06470 | Class IV | 3 | LOD2 | | |
| Lko057S005538 | NNCF01131273.1 | 20948 | 31450 | 10503 | - | Lko_31867 | Class V | 4 | LOD1 | | |
| Lko057S005538 | NNCF01131273.1 | 6676 | 8154 | 1479 | - | Lko_31866 | Class V | 4 | LOD1 | | |
| Lko057S005538 | NNCF01131273.1 | 43198 | 60736 | 17539 | - | Lko_31869 | Class V | 4 | LOD1 | | |
| Lko057S000206 | NNCF01125945.1 | 353321 | 357106 | 3786 | - | Lko_06341 | Class III | | | | 3 |
| Lko057S000206 | NNCF01125945.1 | 404113 | 432386 | 28274 | - | Lko_06342 | Class III | | | | 3 |

982

## Supplementary References

Blankers, T., Oh, K. P., Bombarely, A., & Shaw, K. L. (2018). The genomic architecture of a rapid Island radiation: Recombination rate variation, chromosome structure, and genome assembly of the Hawaiian cricket *Laupala*. In *Genetics* (Vol. 209, pp. 1329-1344).

Blankers, T., Oh, K. P., & Shaw, K. L. (2018). The genetics of a behavioral speciation phenotype in an Island system. In *Genes* (Vol. 9, pp. 346)

Fisher, H. P., Pascual, M. G., Jimenez, S. I., Michaelson, D. A., Joncas, C. T., Quenzer, E. D., . . . Horch, H. W. (2018). De novo assembly of a transcriptome for the cricket *Gryllus bimaculatus* prothoracic ganglion: An invertebrate model for investigating adult central nervous system compensatory plasticity. In *PLoS One* (Vol. 13, pp. e0199070).

Shaw, K. L., & Lesnick, S. C. (2009). Genomic linkage of male song and female acoustic preference QTL underlying a rapid species radiation. In *Proceedings of the National Academy of Sciences* (Vol. 106, pp. 9737-9742).

Xu, M., & Shaw, K. L. (2019). The genetics of mating song evolution underlying rapid speciation: Linking quantitative variation to candidate genes for behavioral isolation. In *Genetics* (Vol. 211, pp. 1089-1104).

Zelle, K. M., Lu, B., Pyfrom, S. C., & Ben-Shahar, Y. (2013). The genetic architecture of degenerin/epithelial sodium channels in *Drosophila*. In *G3: Genes, Genomes, Genetics* (Vol. 3, pp. 441-450).