

Predicting Insect Invasiveness with Whole-Genome Sequencing Data

Cong Huang

Zhejiang University

Nianwan Yang

Chinese Academy of Agricultural Sciences <https://orcid.org/0000-0002-6014-3383>

Shuping Wang

Shanghai Customs College

Xiaodan Fan

Chinese University of Hong Kong <https://orcid.org/0000-0002-2744-9030>

Cong Pian

Zhejiang University <https://orcid.org/0000-0001-7401-2926>

Jiapeng Luo

Zhejiang University

Xi Li

Zhejiang University

Kun Lang

Zhejiang University

Longsheng Xing

Chinese Academy of Agricultural Sciences

Mingxing Jiang

Zhejiang University

Wanxue Liu

Chinese Academy of Agricultural Sciences

Wanqiang Qian

Chinese Academy of Agricultural Sciences

Daniel Simberloff

University of Tennessee

Fanghao Wan

Chinese Academy of Agricultural Sciences

Fei Li (✉ lifei18@zju.edu.cn)

Zhejiang University <https://orcid.org/0000-0002-8410-5250>

Keywords: Insect pest, Invasiveness, Genome features, Comparative genomics, Invasiveness index

Posted Date: December 10th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-39430/v2>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 Predicting Insect Invasiveness with Whole-Genome Sequencing

2 Data

3 Cong Huang^{1,3,8#}, Nianwan Yang^{2#}, Shuping Wang⁴, Xiaodan Fan⁵, Cong Pian¹, Jiapeng
4 Luo¹, Xi Li⁶, Kun Lang¹, Longsheng Xing³, Mingxing Jiang¹, Wanxue Liu², Wanqiang Qian^{3*},
5 Daniel Simberloff^{7*}, Fanghao Wan^{2,3,9*}, Fei Li^{1*}

6 ¹Ministry of Agriculture Key Laboratory of Molecular Biology of Crop Pathogens and Insect
7 Pests, Institute of Insect Sciences, College of Agriculture and Biotechnology, Zhejiang
8 University, Hangzhou, 310058, China

9 ²State Key Laboratory for Biology of Plant Diseases and Insect Pests, Institute of Plant
10 Protection, Chinese Academy of Agricultural Sciences, Beijing, 100193, China

11 ³Shenzhen Branch, Guangdong Laboratory for Lingnan Modern Agriculture, Genome Analysis
12 Laboratory of the Ministry of Agriculture, Agricultural Genomics Institute at Shenzhen, Chinese
13 Academy of Agricultural Sciences, Shenzhen, 518120, China.

14 ⁴Technical Centre for Animal Plant and Food Inspection and Quarantine, Shanghai Customs,
15 Shanghai, 200135, China

16 ⁵Department of Statistics, The Chinese University of Hong Kong, Hong Kong SAR, China

17 ⁶College of Computer Science, Zhejiang University, Hangzhou, 310027, China

18 ⁷Department of Ecology & Evolutionary Biology, The University of Tennessee, Knoxville, TN,
19 37996, USA.

20 ⁸Plant Protection College, Hunan Agricultural University, Changsha, 410128, China

21 ⁹College of Plant Health and Medicine, Qingdao Agricultural University, Qingdao, 266109,
22 China

23

24 # These authors contributed equally

25 *Corresponding authors, Dr. Fei Li, Dr. Fanghao Wan, Dr. Daniel Simberloff, Dr. Wanqiang
26 Qian, **Email:** lifei18@zju.edu.cn; [wanfanghao@caas.cn](mailto:wangfanghao@caas.cn); dsimberloff@utk.edu;
27 qianwanqiang@caas.cn

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

Abstract

Background: Invasive alien insects threaten agriculture, biodiversity, and human livelihoods globally. Unfortunately, insect invasiveness still cannot be reliably predicted. Empirical policies of insect pest quarantine and inspection are mainly designed against species that are already problematic.

Results: We conducted a comparative genomic analysis of 37 invasive insect species and six non-invasive insect species, showing that the gene families associated with defense, protein and nucleic acid metabolism, chemosensory function, and transcriptional regulation were significantly expanded in invasive insects, suggesting that enhanced abilities in self-protection, nutrition exploitation, and locating food or mates are intrinsic features conferring invasiveness in insects. By using these intrinsic genome features, we proposed an invasiveness index and estimated the invasiveness of 99 other insect species with genome data, classifying them as highly, moderately, or minimally invasive. Insects possessing all these aforementioned enhanced abilities are predicted to be highly invasive, and vice versa. Next, a logistic-regression classifier was trained to predict insect invasiveness, achieving 93.2% accuracy.

Conclusions: We present evidence that several traits may confer invasiveness in insects and these features can be used to predict insect invasiveness accurately, and we quantify insect invasiveness with an invasiveness index.

Keywords

Insect pest; Invasiveness; Genome features; Comparative genomics; Invasiveness index

51 **Background**

52 Invasive species threaten agriculture, biodiversity, and human livelihoods. The estimated
53 global economic loss to invasive insects is US \$70.0 billion annually [1]. Increased
54 globalization and connectedness via trade, as well as environmental changes owing to
55 climate change, will most likely significantly increase invasive species threats [2, 3]. Many
56 studies have shown some common properties in invasive species. Invasive plants have
57 syndromes including: a tendency to be annual or biennial; increased plant height and specific
58 leaf area; hermaphroditism with longer and earlier flowering; clonal growth and monoecy; and
59 higher fecundity [4, 5]. Invasive birds tend to: be large; prioritize future over current
60 reproduction; be less migratory; and be widespread in the source region [6]. In contrast,
61 invasive freshwater fishes have smaller body size, fast reproduction, high activity, and
62 boldness, and are omnivorous with high physiological tolerances [7]. Invasive insects are
63 reported to tend to have some intrinsic features such as parthenogenesis, high dispersal
64 ability, a dormant or resilient stage, and a longer adult stage [8].

65 These observations have yielded at least four hypotheses about invasive species: 1)
66 enemy release hypothesis – escape from natural enemies in the original habitat [9]; 2)
67 increased competitive ability hypothesis – efficient transfer of biological resources from
68 enemy defense to growth and reproduction [10]; 3) novel weapons hypothesis – carrying
69 parasites such as microsporidia that negatively affect or kill native species but not the
70 invading species [11]; and 4) inherent superiority hypothesis – invasive species have intrinsic
71 traits superior to those of non-invasive species, at least in new regions [12-15]. We only focus
72 on the last hypothesis.

73 Accurately identifying invasiveness-related traits and predicting invasiveness of a species
74 is important for pest risk assessments and developing national quarantine policies. However,
75 the traits identified as associated with invasiveness are quite controversial and do not
76 characterize all invasive species, especially in insects [16]. The controversy has hampered
77 development of a highly accurate method of predicting invasiveness, though much effort has
78 been expended on this project, such as Invasive Species Predictive Schemes (ISPS) [17] and
79 the SCOPE project [5]. A problem is that an invasion consists of several distinct stages [18],
80 and traits that would lead a species to pass successfully through one stage may not be the
81 same traits that would conduce to success at a different stage. With respect to risk, different
82 stages are at issue. The first stage, transport and initial introduction, consists of a propagule
83 arriving with human assistance in a new, distant site. Whatever traits a species has that
84 facilitate its association with a transport vector (such as ballast water, shipping containers, or
85 agricultural products) increase risk of transport [19]. In this paper, we focus on the next
86 stages, establishment and spread. Once a species has arrived in a new region, do particular
87 traits increase the risk that it will persist and spread? The propensity to establish and spread
88 once introduced is what we define as “invasiveness” in this study, although the term
89 “invasiveness” has also been used at times in the literature to refer to the first stage, simply
90 arriving in a new region.

91 As the cost of whole-genome sequencing has decreased dramatically, hundreds of insect
92 genomes have been sequenced [20], providing an opportunity to conduct comparative

93 genomic analysis in insects. Comparative genomics is a powerful tool and it was already
94 successfully applied in revealing the molecular mechanism of ancient maize adaptation to
95 temperate [21] and tracking the genomic changes associated with dietary specialization for
96 mammals with interesting results [22]. Recently, it has been reported that human physical
97 traits can be predicted by using whole-genome sequencing data with a high accuracy [23].
98 Since invasive insects distinguish with the others by invasiveness related traits [12], it is
99 worthwhile to explore whether these invasiveness related traits could be predicted by
100 genomic features. To this end, we focused on identifying invasive features at the gene family
101 level. By conducting comparative genomic analysis with 37 invasive insects and six non-
102 invasive insects, we proposed an invasiveness index to quantify insect invasiveness and
103 developed a machine-learning algorithm (Determining Invasiveness based on Genome
104 Sequences, DIGS) to predict insect invasiveness.

105

106 **Results**

107 **Selection of 37 invasive insects and six non-invasive insects**

108 From the 142 insect species for which have complete genome data and annotations are
109 available (Supplementary Table 1), we excluded ones lacking high-quality genome
110 assemblies as well as those that are not confirmed as having been introduced anywhere (89
111 species) and thus could not be classified as either invasive or non-invasive. We confirmed: 1)
112 37 insects known to be invasive by literature references, including nine Diptera, two
113 Coleoptera, fifteen Hymenoptera, five Hemiptera, and six Lepidoptera; 2) six non-invasive
114 insects (two Diptera, one Hemiptera, and three Lepidoptera) according to the criterion of
115 having been introduced to non-native regions but not spreading or exhibiting any signs of
116 invasion in the introduced regions (Supplementary Table 2). From these 43 insect species, we
117 identified 183 single-copy orthologous genes by all-vs-all BLASTP [24] against all proteins in
118 the OrthoMCL [25]. We constructed a phylogenetic tree using these single-copy orthologous
119 genes to infer the evolutionary relationship of these species.

120

121 **The general genome features are not related to invasiveness**

122 We calculated general genome features of these 43 insects including genome size, GC
123 content, gene number, amount of repeat sequences, number of expanded gene families, and
124 number of contracted gene families (Fig. 1). None of these features differed significantly
125 between invasive and non-invasive insects, indicating that high invasiveness might be
126 ascribed only to several key gene families closely associated with invasive traits, rather than
127 to general genome features (Supplementary Fig. 1 and Supplementary Table 7).

128

129 **Identifying gene families associated with insect invasiveness**

130 It has been reported that invasive insects share some traits, such as nutrition acquisition
131 advantage, advanced defense systems, and high reproductive ability [8]. For the inherent
132 superiority hypothesis to be valid, we reasoned that gene families conferring functions related
133 to invasiveness should be positively selected and most likely expanded. To this end, we
134 analyzed the expansions and contractions of gene families in a phylogenetic context in the 43

135 insects using the program CAFÉ (v3.0) [26]. We found 36 gene families to have expanded in
 136 at least 13 of the 37 invasive species. This was determined with the additional criterion: the
 137 ratio of 'number of invasive species in which the gene family expanded' to 'number of non-
 138 invasive species in which the gene family expanded' exceeded 12. The criterion was
 139 determined by testing the accuracy of invasiveness classification for a range of ratios from 4
 140 to 13; among these, a ratio of 12 achieved the highest accuracy (Supplementary Table 8). The
 141 gene families expanded more frequently in invasive species by this criterion were treated as
 142 candidate gene families and were grouped into four categories based on their functional
 143 associations: 1) associated with defense; 2) associated with energy; 3) associated with
 144 chemosensory function; 4) associated with transcriptional regulation (Supplementary Table 8).

145 Next, to evaluate the contribution of these candidate gene families to invasiveness, we
 146 used a two-step logistic regression procedure (see Materials and Methods) to select the gene
 147 families whose expansion might contribute to invasiveness and to determine the relative
 148 weights of their contributions (based on the expanded gene numbers in Supplementary Table
 149 3). The results show that in total 14 gene families are associated with invasiveness
 150 (Supplementary Table 4). The expansion pattern of these gene families varied in different
 151 invasive insects, suggesting that a variety of traits have conferred invasiveness in insects
 152 (Fig. 2). For example, invasive hymenopterans have enhanced defense ability and advanced
 153 chemosensory function, while invasive lepidopterans have enhanced abilities of defense and
 154 energy metabolism, and particularly transcriptional regulation.

155

156 **Invasiveness index for insect invasiveness**

157 We next seek to estimate the invasiveness of insects by using the expansion indexes of four
 158 function groups. We calculated the expansion indexes for each function group in each insect
 159 species, which involves the weighting coefficients of the 14 gene families that resulted from
 160 the first-step logistic regression (see methods) and the corresponding expansion gene
 161 numbers (Supplementary Table 9). We built the invasiveness index of a species by using the
 162 second-step logistic regression model (see methods) that estimates the weighting coefficients
 163 of the four function groups with the following steps:

164 1) $m_j = 116.98y_{1j} + 12.98y_{2j} + 6.29y_{3j} + 6.12y_{4j} + 63.89$,

165 2) $n_j = \begin{cases} \log_{10}(|m_j|), & \text{if } m_j \geq 1 \\ 0, & \text{if } -1 < m_j < 1 \\ -\log_{10}(|m_j|), & \text{if } m_j \leq -1 \end{cases}$,

166 3) $z_j = 1 - \frac{1}{1 + e^{n_j}}$,

167 where z_j is the invasiveness index of the j th species and y_{1j} to y_{4j} are the expansion indexes of
 168 the four function groups of the j th species.

169 Then we calculated the invasiveness indexes of all 142 insects (Fig. 3, Fig. 4,
 170 Supplementary Table 9 and Supplementary Table 5). As expected, for the 37 invasive insects,
 171 all have high invasiveness indexes, by contrast, all six non-invasive insects have minimal
 172 invasiveness indexes (Fig. 3). For the 99 other insects, we classified their predicted
 173 invasiveness into three levels based on invasiveness indexes: high invasiveness (0.9 to 1),

174 moderate invasiveness (0.2 to 0.9), and minimal invasiveness (0 to 0.2) (Fig. 4). Among these
175 99 insects, ten species have been reported to be highly invasive; eight of these were
176 assigned a high invasiveness index and one was assigned a moderate invasiveness index,
177 the remaining one was assigned a minimal invasiveness index (Fig. 4).

178 The results showed that these four aspects of capacities were generally essential for
179 high invasiveness: defense, energy, chemosensory function, and transcriptional regulation.
180 Highly invasive insects tend to have high expansion indexes in all four function categories;
181 insects that are minimally invasive tend to have low expansion indexes in all four function
182 groups. Among the rest 89 insect species for which we lack adequate data on invasiveness,
183 we predicted species to be moderately invasive if they have high expansion indexes in some
184 function categories but low expansion indexes in others, while we predicted them to be highly
185 or minimally invasive according the rules just stated (Fig. 4).

186

187 **Classifying insect invasiveness by machine learning**

188 Having identified putative inherent genome features associated with insect invasiveness, we
189 adopted these features to develop a machine learning algorithm DIGS, in order to classify
190 insects in terms of invasiveness. DIGS used a random forest algorithm for feature selection
191 and then used a logistic regression model to construct a classifier; six-fold cross-validation
192 was used to train the DIGS classifier.

193 In each cross-validation, we used the R package “Boruta” [27] to evaluate the
194 contributions of the 36 candidate gene families to invasiveness. Boruta is designed as a
195 wrapper with a random forest classification. The key gene families selected by Boruta were
196 next used to construct a logistic regression model to estimate classification efficiency. As a
197 result, two gene families, pao retrotransposon peptidase, and putative nuclease HARBI, were
198 stably selected to predict insect invasiveness (Supplementary Table 6).

199 Given the high average accuracy and the balance between positive and negative
200 samples of the 43 species (Supplementary Table 6), these two gene families were used to
201 construct a logistic classifier to classify invasiveness. The results indicated that DIGS
202 performs well in classifying insect invasiveness, with an average accuracy of 93.2%.
203 Sensitivity, specificity, and precision were 88.1%, 100%, and 100%, respectively
204 (Supplementary Table 10). Based on the analysis of a ROC (Receiver Operating
205 Characteristic) curve, the AUC (Area under the Curve of ROC) [28] was 0.953, suggesting
206 good performance by the DIGS classifier (Fig. 5). Next, we used this classifier to predict
207 invasiveness of the other 99 insects, those not used for training. With a cutoff of 0.5, 56
208 species (56.6%) were classified as invasive.

209 Because we have developed two systems separately to evaluate insect invasiveness, we
210 compared the consistency of the invasiveness index and the DIGS. Fig. 6a showed that
211 94.1% of highly invasive insects as determined by the invasiveness index were predicted to
212 be invasive by DIGS, whereas 85.3% of insects predicted not to be invasive by the
213 invasiveness index were predicted not to be invasive by DIGS (Fig. 6a). On the other hand, of
214 those 56 insects predicted to be invasive by DIGS, 57.1% were classified as highly invasive
215 and 33.9% as moderately invasive by the invasiveness index. Of the 43 insects predicted not

216 to be invasive by DIGS, the invasiveness index analysis predicted 67.4% would not be
217 invasive (Fig. 6b). These results showed substantial consistency between the invasiveness
218 index and DIGS, we can use either to qualitative or quantitative estimate the insect
219 invasiveness.

220

221 **Discussion**

222 In the process of biological invasions, many invasive species experienced a "lag". It is a
223 period of time that after they arrived at a new region while the populations are small and
224 geographically restricted. This is so-called the "lag-time". The risks of such species at this
225 stage are often ignored or underestimated, however, no matter how long they will take before
226 causing serious trouble, they will successfully invaded at the end [29]. For example, the
227 Brazilian pepper (*Schinus terebinthifolius*), which was present as a restricted ornamental for
228 at least 50 years before its rapid invasion everywhere [30]. Therefore, one of the most
229 significant challenges regarding biological invasions is to predict the risk of successful
230 invasion. Meta-analysis is increasingly used to predict the invasion risk of non-native species
231 globally. For different invasive stages, the variables are under consideration differs: for the
232 introduction risk, variables often include the quantity of international trade [2, 3] and the
233 capacity of the transport vectors to assist arrival [2]; for the establishment risk, variables
234 include disturbance factors [2], biodiversity indices [31], and the similarity of biotic and abiotic
235 conditions between the native locations and the location of a newly arrived invasive alien
236 species (IAS) [3]. However, the species characteristics that increase the risk of population
237 establishment and spread once the species is introduced, which we define as "invasiveness"
238 in this study, are still not well determined. This lacuna contributes to ineffective management
239 and slow responses to newly arrived IAS [32]. Here, we have attempted to fill this gap by
240 presenting a new approach based on machine learning and genome data to predict high-risk
241 invasive insect species.

242 Based on the hypothesis that invasive insects tend to share particular invasiveness-
243 related traits, we conducted a comparative genomic analysis of invasive and non-invasive
244 insect species to identify gene families commonly expanded in invasive species. This strategy
245 yielded 14 gene families associated with insect invasiveness. These 14 gene families can be
246 grouped into four function categories: defense, energy, chemosensory function, and
247 transcriptional regulation. These results support the previous hypothesis that invasive species
248 should have certain intrinsic traits that are superior in new locations to those of non-invasive
249 species. Invasive insects tend to have high abilities to exploit nutrition and to defend
250 themselves, as well as advanced chemosensory abilities. We proposed an invasiveness
251 index to quantify invasiveness by weighting the abilities of the aforementioned four groups.
252 This invasiveness index is the first metric for evaluating insect invasiveness at the genome
253 level, and it should aid risk assessment and provide strong theoretical support for quarantine
254 policy decisions [33].

255 We found that insects with high expansion indexes in all or most of the four groups tend
256 to be highly invasive, whereas insects without high expansion indexes or with only one
257 category with a high expansion index tend to be minimally invasive. This result appears

258 consistent with commonsensical notions of invasive ability. However, this fact does not
259 support the trade-off hypothesis, which assumes that invasive species must allocate limited
260 energy and resources to either growth or defense [34]. Because of physical and chemical
261 constraints, resource allocation limitations, antagonistic pleiotropy, and linkage disequilibrium,
262 the trade-off hypothesis suggests that having a advantages trait for one function may
263 simultaneously reduce the strength of other functions [35]. However, our comparative
264 genomics analysis shows that highly invasive insects have most or even all four types of
265 enhanced abilities.

266 We applied an invasiveness index to evaluate 99 insect species that have annotated
267 genome data. About half of these species were given a high invasiveness index value, and
268 these species were also generally predicted to be highly invasive by DIGS, supporting the
269 reliability of this metric. For testing our model in the performance of lag-time IAS, we used the
270 DIGS and invasiveness index formula to qualitative and quantitative estimate the
271 invasiveness of a recent invasive insect, *Spodoptera frugiperda*, which has gone through
272 outbreaks with irregular intervals in its native regions for two centuries before its successful
273 invasion in Africa and Asia in 2016 [36]. The result showed that it is classified as invasive
274 species by DIGS and its invasiveness index is as high as 0.96.

275 We noticed that some congeners were classified very differently in invasiveness. It is true
276 that congeneric species even some cryptic species differ in some gene families, such as
277 cytochrome P450 genes and UDP glycosyltransferases in two cryptic species of invasive
278 whitefly, *Bemisia tabaci* Middle East-Asia Minor 1 (MEAM1, or 'B') and Mediterranean (MED,
279 or 'Q') [37]. We emphasize that the present classifier was trained with a very small set of non-
280 invasive species, these were all in just three orders, and that all the invasive species are in
281 only five orders. It remains possible that the small sample size induced bias in feature
282 selection and that limited phylogenetic diversity of species with adequate bases for
283 classification as to invasiveness limits the domain of application of the striking result depicted
284 in Fig. 3. Although the sample size for the negative training set was particularly small, the
285 DIGS classifier still performed well, suggesting that invasive insects do have some inherent
286 superiorities in new settings compared with non-invasive insects. These differences in
287 inherent superiorities are reflected in the genome data. As the cost of genome sequencing
288 decreases, more insect genomes will be available with high assembly quality. With additional
289 genomes, larger positive and negative training datasets can be constructed to achieve better
290 classification efficiency and to determine the extent to which our invasiveness index applies
291 beyond the five orders for which we currently have data (the 99 species that we did not use
292 for training are also all in the same five orders). A more robust classifier can then be trained
293 with better, more robust performance. In addition, our approach to profiling invasive species
294 used only the features of gene family expansions. As genomic information accumulates, other
295 genome features such as single nucleotide polymorphisms (SNPs), copy number variants
296 (CNV), and gene expression level at the subspecies level can be obtained and may be useful
297 in predicting invasiveness, which should improve prediction accuracy and understanding of
298 the molecular basis of invasiveness. Of course, the fact that a species truly has a tendency to
299 become invasive once introduced does not mean that every introduction of that species will

300 lead to invasion. Aside from the fact that some probability exists for stochastic reasons alone
301 that any population will be lost when very small (e.g., in the earliest stages of establishment),
302 physical and biotic environmental factors differ among introductions and play some role in
303 whether an invasion actually occurs. This fact is reflected in many examples of species that
304 usually become invasive when introduced but fail to do so occasionally [38].

305

306 **Methods**

307 **Genome resources and species selection.**

308 We downloaded 142 insect genome assemblies and the corresponding annotation data,
309 including Coleoptera, Diptera, Hemiptera, Hymenoptera, and Lepidoptera, from the National
310 Center for Biotechnology Information (NCBI) [39], InsectBase [40], VectorBase [41],
311 Fireflybase (<http://www.fireflybase.org/>), Ensembl Genomes [42], GigaDB [43], Fourmidable
312 [44], MonarchBase [45], and AphidBase [46]
313 (Supplemental Table 1).

314 We used the genome characteristic value of Scaffold N50 to filter the species with low
315 quality genomes, which is positively associated with genome quality that the higher the better
316 [47]. Species with a genome assembly with Scaffold N50 < 400 Kb were excluded. When a
317 protein-coding gene has different alternative splicing forms, the longest transcript was
318 chosen.

319 We analyzed 43 insect species, including 37 confirmed as invasive species by literature
320 references, which were used as positive samples of invasive insects, and six confirmed as
321 non-invasive by literature references, which were used as negative samples of non-invasive
322 insects (Supplemental Table 2). The invasiveness index and the classifier were applied to the
323 remaining 99 insects.

324 **Gene family analysis.**

325 We used TreeFam [48], which considers phylogenetic relationship, to define gene families
326 that descended from a single gene of the most recent common ancestor. The annotated
327 protein-coding genes of 37 invasive insect species, six non-invasive insect species, and the
328 additional 99 insect species were used for the application as noted above.

329 **Reconstruction of phylogenetic tree.**

330 We performed phylogenetic analyses using proteins from all 43 invasive and non-invasive
331 species, and *Tetranychus urticae* was used as an outgroup. OrthoMCL [25] was used with
332 default parameters to identify gene groups based on sequence similarities resulting from an
333 all-against-all BLASTP search [24]. We found 183 single-copy orthologous genes shared by
334 all species. Multiple sequence alignments of orthologous genes from all species were
335 produced by MAFFT v7 [49] with default parameters, and the aligned results were trimmed by
336 trimAl [50] to remove low-quality regions with the parameter “-automated1”. Finally, we
337 merged all 183 trimmed single-copy orthologous genes for each species to create a super
338 gene [51]. RAxML [52] was then used with the LG+I+F model, which is calculated with
339 ProtTest in IQ-TREE [53], to estimate a maximum likelihood tree starting with 1000 bootstraps
340 followed by likelihood optimization.

341 **Estimation of divergence time.**

342 A nonparametric rate-smoothing method [54] and a semiparametric penalized likelihood (PL)
343 method [55] were used to estimate the divergence time with the software r8s (V1.7.1) [56]. An
344 optimal tree obtained by RAxML [52] was used as an input tree for the divergence time
345 estimation. The cross-validation approach (with parameters “cvstart=0, cvinc=1, cvnum=18”)
346 was used to determine the optimal level of rate- smoothing of the PL analyses with smoothing
347 parameters varying from 1 to 1e17. The result showed that a smoothing parameter of 1 was
348 optimum for these data. To estimate divergence time, we calculated ages of nodes within the
349 phylogeny based on calibration points. Our calibration points were: 1) the most recent
350 common ancestor of the clade including *Papilio polytes* and *Plutella xylostella*, constrained to
351 be 140 Mya (million years ago); 2) the most recent common ancestor of the clade including
352 *Bombyx mori* and *Manduca sexta*, constrained to be 39.8 Mya; and 3) the most recent
353 common ancestor of the clade including *Aedes albopictus* and *Drosophila biarmipes*,
354 constrained to be 157.8 Mya [51, 57].

355 **Gene gain and loss.**

356 To identify gene family evolution as a stochastic birth and death process, we applied the
357 likelihood model originally implemented in the software CAFÉ (v3.0) [26]. Phylogenetic tree
358 topology and branch lengths were taken into account to infer the significance of change in
359 gene family size in each branch. The gene number of gene families defined by TreeFam [48]
360 in each insect and the phylogenetic tree corrected by r8s which is used to estimate the
361 absolute rates of molecular evolution and divergence times on a phylogenetic tree [56] were
362 used as input files for CAFÉ 3.0 [26].

363 **Comparative analysis of genomic features.**

364 We calculated the genome features of all insects, including genome size, GC%, number of
365 protein-coding genes, length of repetitive sequences, and gene number of expansion or
366 contraction. We used t-tests (we also did permutation tests with the same result) to compare
367 differences in genome features between all invasive and non-invasive insects, as well as fly
368 and lepidopteran species separately. We identified repetitive sequences using the
369 RepeatMasker [58] pipeline with “ncbi” set as the search engine and “insects” for the
370 parameter (-species). In addition, RepBase [59] was provided as a custom library to locate
371 associated repetitive elements in genomes of each species.

372 **Commonly expanded gene families of invasive insects.**

373 The expanded and contracted gene families in each species, as well as the expanded and
374 contracted gene numbers of each gene family, were extracted by using an in-house Perl
375 script (<https://github.com/hc18/extract-commonly-expanded-gene-families>) and calculated by
376 comparing a species with its parent node in the phylogenetic tree (Supplementary Table 3).
377 We first tried to define the significant gene family expansion as the ratio of the proportions of
378 species with expanded gene family in invasive vs. in non-invasive. However, some gene
379 families are not expanded in non-invasive species, the denominator will be zero. So, we
380 subsequently defined the commonly expanded gene families of invasive insects with two
381 criteria: 1) expanded in at least one-third of invasive species (≥ 13), or less than or equal to
382 half of the total number of non-invasive species (≤ 3). 2) an expanding ratio, defined as the
383 ratio of the number of invasive species to the number of non-invasive species in which the

384 gene family is expanded, greater than 12 (the criterion was determined by testing the
 385 accuracy of invasiveness classification for a range of ratios from 4 to 13, among them, ratio of
 386 12 achieved the highest accuracy). This protocol generated a number of candidate gene
 387 families that might be related to invasiveness. Next, we annotated these candidate gene
 388 families using the corresponding protein sequences as queries to perform a BLASTP [24]
 389 search (e-value cutoff of 1e-5) against the UniProt [60] database. We then grouped these
 390 gene families according to their annotated functions.

391 **Gene families associated with invasiveness and expansion index.**

392 To estimate the contribution of these gene families to invasiveness, we used the candidate
 393 gene families from each function group to train the first-step logistic regression model with
 394 70% of the 43 species as the training set and the other 30% as the testing set. The partition of
 395 training set and testing set was randomly done 30 times. The coefficient of each gene family
 396 in the logistic regression was regarded as its weight coefficient for invasiveness within the
 397 function group. If a gene family has a non-positive weight coefficient for more than one third of
 398 the partitions, we removed it and performed the same pipeline again until all remaining gene
 399 families had positive weight coefficients for more than half of the partitions. We next selected
 400 the model with the highest AUC [28] among all partitions for each function group. Finally, we
 401 used 14 gene families from four function groups and their weight coefficients to construct the
 402 expansion indexes (Supplementary table 4). The expansion index formula for a specific
 403 function group containing n expanded gene families associated with invasiveness is defined
 404 as follows:

405
$$y_j = \frac{\sum_{i=1}^n k_i x_{ij}}{\sum_{i=1}^n k_i}$$

406 where y_j is the expansion index of the function group for the j th species, i is the number of
 407 the i th invasiveness-related gene family in the function group, n is the total number of
 408 invasiveness-related gene families in the function group, x_{ij} is the expanded gene number of
 409 the i th invasiveness-related gene family in the function group of the j th species, and k_i is the
 410 weight coefficient, i.e., the logistic regression coefficient, of the i th invasiveness-related gene
 411 family to invasiveness.

412 **Invasiveness index formula.**

413 To estimate the weight coefficients of these expansion indexes of the five function groups to
 414 invasiveness, we used them to train the second-step logistic regression model with 70% of
 415 the 43 species as the training set and the remaining 30% as the testing set. A model with
 416 highest AUC [28] was fitted. Subsequently, the expansion indexes and their corresponding
 417 weight coefficients as well as the intercept in the logistic regression model were used to
 418 construct the invasiveness index formula in three steps:

419

420 1) $m_j = \sum_{i=1}^g k_i y_{ij} + b$,

421 2)
$$n_j = \begin{cases} \log_{10}(|m_j|), & \text{if } m_j \geq 1 \\ 0, & \text{if } -1 < m_j < 1 \\ -\log_{10}(|m_j|), & \text{if } m_j \leq -1 \end{cases},$$

422 3)
$$z_j = 1 - \frac{1}{1 + e^{n_j}},$$

423

424 where z_j is the invasiveness index of the j th species, y_{ij} is the expansion index of the i th
425 function group of the j th species, g is the total number of function groups (in this study, four),
426 k_i is the weight coefficient of the i th function group, and b is the constant in the logistic
427 regression. We used the first step to calculate the total weight m_j of four function groups that
428 contribute to the invasiveness classification. The second step was used to normalize the m_j ;
429 and we calculated the invasiveness index by a logistic formula in third step.

430 **Applying the invasiveness index formula to estimate invasiveness.**

431 To calculate the invasiveness indexes of the other 99 insects, we added one of the 99
432 species to the data set of the 44 species analyzed (37 invasive insects, six non-invasive
433 insects, and the outgroup *Tetranychus urticae*) at a time, then the same methods and
434 parameters were applied to find their single-copy orthologous genes, construct the
435 phylogenetic tree, correct divergence time, and calculate gene gain and loss for each species.
436 Finally, we calculated invasiveness indexes of all of the 99 species using the invasiveness
437 index formula from the above section (Supplementary Table 5).

438 **Invasiveness classification by the machine-learning algorithm.**

439 A machine-learning algorithm named DIGS was built for invasiveness classification.

440 First, we used a random forest algorithm for feature selection and used a logistic
441 regression to estimate the classification performance of features selected by this algorithm.
442 Six-fold cross-validation was used to estimate the accuracy and stability of feature selection.
443 To guarantee that each of the six non-invasive species would be allocated to the testing set
444 once, only one non-invasive species was allocated to the testing set in each iteration of cross-
445 validation. The remaining five non-invasive species were allocated to the training set.
446 According to the ratio of 1:5 to allocate species into testing and training sets of non-invasive
447 species, in each iteration of cross-validation, the 37 invasive species were randomly
448 distributed into six groups (each group has six or seven species), one group (six or seven
449 species) was allocated to the testing set, while the remaining five groups (a total of 31 or 30
450 species) were allocated to the training set. The R package “Boruta” [27] was used to perform
451 feature selection with the parameters of “ntree=1000” and “maxRuns=1000”, using all 36
452 candidate gene families of the training set. This algorithm selects features with a random
453 forest classification algorithm and a statistical test. Features that do not contribute more to the
454 classification information than random features were removed. We used the expanded or
455 contracted gene numbers of species as the input data. Two features were stably confirmed to
456 be important for invasiveness by Boruta with six-fold cross-validation (Supplementary Table
457 6).

458 We then used the two gene families confirmed as important to invasiveness by previous
459 steps to construct the logistic regression model as the DIGS classifier. To treat sample size
460 imbalance, the entire dataset of negative samples was duplicated six times before being used
461 for training or testing in the logistic regression model [61] as follows:

462
$$y_j = 1 - \frac{1}{1 + e^{\sum_{i=1}^n k_i x_{ij} + b}},$$

463 where y_j is the probability that the j th species belongs to the invasive set, i is the number of i th

464 invasiveness-related gene family confirmed by Boruta, n is the total number of invasiveness-
465 related gene families confirmed by Boruta (here $n=2$), x_{ij} is the expanded gene number of the
466 i th invasiveness-related gene family of the j th species, k_i is the weight coefficient of the i th
467 invasiveness-related gene family in the classification model. By the six-fold cross-validation in
468 DIGS, two features were stably estimated to associate with insect invasiveness, x_{1j} (pao
469 retrotransposon peptidase) and x_{2j} (putative nuclease HARB11) with the corresponding
470 coefficients of 0.31 and 1.86, respectively; b was equal to 1.20.

471 We then used the DIGS classifier to calculate the probabilities that the species in the
472 testing set are invasive. A species was predicted to be invasive if the probability exceeded 0.5
473 by DIGS; otherwise it was predicted not to be invasive.

474

475 **Abbreviations**

476 DIGS: Determining Invasiveness based on Genome Sequences; SNPs: single nucleotide
477 polymorphisms; RAxML: Random accelerated maximum likelihood; CNV: copy number
478 variants;

479 ISPS: Invasive Species Predictive Schemes; SCOPE: Scientific Committee on Problems of
480 the Environment; BLAST: Basic Local Alignment Search Tool; CAFÉ: Computational Analysis
481 of gene Family Evolution; ROC: Receiver Operating Characteristic; AUC: Area under the
482 Curve of ROC; GC: Guanine and Cytosine nucleotides; IAS: Invasive Alien Species; P450:
483 Cytocrome P450; UDP: Uridine Diphosphate; MEAM1: Middle East-Asia Minor 1; MED:
484 Mediterranean

485

486 **Declarations**

487 **Acknowledgements**

488 Not applicable.

489

490 **Authors' contributions**

491 F.L. conceived the work and designed the experiment plan; F.W., D.S., and W.Q. designed
492 and improved the experiment plans. D.S., N.Y., and C.H. determined the invasive and non-
493 invasive insects by reference mining; C.H., N.Y., and L.X. collected the genome data; C.H.
494 carried out machine learning classification of invasiveness; X.F., C.P., J.L., K.L., and X.L.
495 participated in the discussion of machine learning work. S.W., W.Q., L.X., M.J., and W.L.
496 participated in the discussion of insect invasiveness. F.L., N.Y., C.H., D.S., and F.W. wrote the
497 manuscript.

498

499 **Funding**

500 This work was supported by the National Key Research and Development Project of China
501 [2016YFC1200600, 2016YFC1201200, 2017YFC1200600]. The funders had no role in study
502 design, data collection, and analysis, or in the decision to publish or in preparing the
503 manuscript.

504

505 **Availability of data and materials**

506 All genomic data used in this study could be downloaded from the databases which have
507 been listed in the supplementary table 1.

508

509 **Ethics approval and consent to participate**

510 Not applicable.

511

512 **Consent for publication**

513 Not applicable.

514

515 **Competing interests**

516 The authors declare no conflict of interest related to the results reported in this study.

517

518 **References**

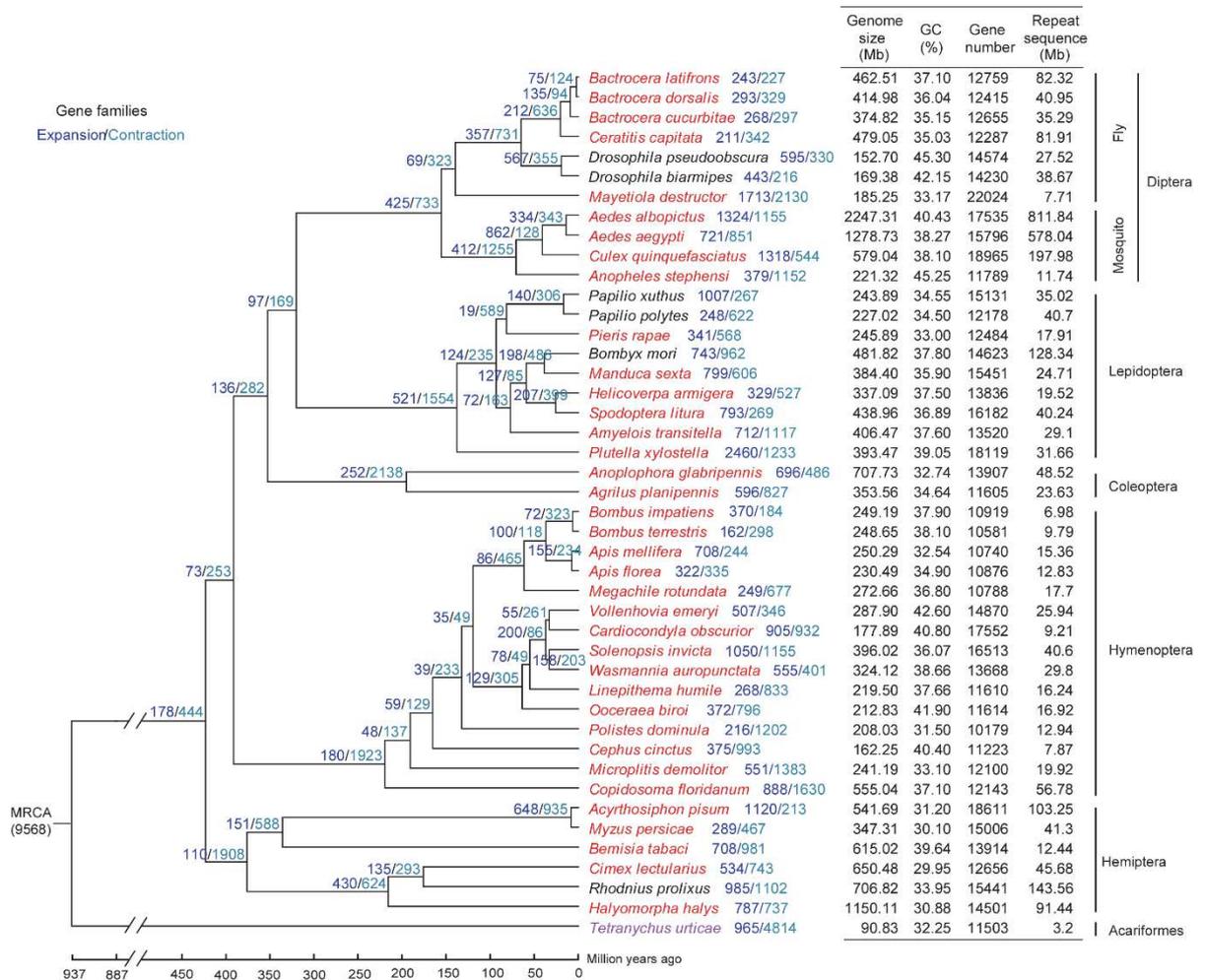
- 519 1. Bradshaw CJ, Leroy B, Bellard C, Roiz D, Albert C, Fournier A, Barbet-Massin M, Salles
520 JM, Simard F, Courchamp F: **Massive yet grossly underestimated global costs of**
521 **invasive insects**. *Nat Commun* 2016, **7**:12986.
- 522 2. Early R, Bradley BA, Dukes JS, Lawler JJ, Olden JD, Blumenthal DM, Gonzalez P,
523 Grosholz ED, Ibanez I, Miller LP *et al*: **Global threats from invasive alien species in**
524 **the twenty-first century and national response capacities**. *Nature Communications*
525 2016, **7**.
- 526 3. Paini DR, Sheppard AW, Cook DC, De Barro PJ, Worner SP, Thomas MB: **Global threat**
527 **to agriculture from invasive species**. *P Natl Acad Sci USA* 2016, **113**(27):7575-7579.
- 528 4. Moravcova L, Pysek P, Jarosik V, Havlickova V, Zakravsky P: **Reproductive**
529 **characteristics of neophytes in the Czech Republic: traits of invasive and non-**
530 **invasive species**. *Preslia* 2010, **82**(4):365-390.
- 531 5. Pyšek P, Richardson DM: **Traits associated with invasiveness in alien plants: where**
532 **do we stand?** In: *Biological invasions*. Edited by Caldwell MM, Heldmaier G, Jackson
533 RB, Lange OL, Mooney HA, Schulze ED, Sommer U. Berlin Heidelberg: Springer; 2008:
534 97-125.
- 535 6. Sol D, Maspons J, Vall-Llosera M, Bartomeus I, Garcia-Pena GE, Pinol J, Freckleton RP:
536 **Unraveling the life history of successful invaders**. *Science* 2012, **337**(6094):580-583.
- 537 7. Cote J, Fogarty S, Weinersmith K, Brodin T, Sih A: **Personality traits and dispersal**
538 **tendency in the invasive mosquitofish (*Gambusia affinis*)**. *P Roy Soc B-Biol Sci* 2010,
539 **277**(1687):1571-1579.
- 540 8. Ochocki BM, Miller TEX: **Rapid evolution of dispersal ability makes biological**
541 **invasions faster and more variable**. *Nature Communications* 2017, **8**.
- 542 9. Colautti RI, Ricciardi A, Grigorovich IA, MacIsaac HJ: **Is invasion success explained**
543 **by the enemy release hypothesis?** *Ecol Lett* 2004, **7**(8):721-733.
- 544 10. Callaway RM, Ridenour WM: **Novel weapons: invasive success and the evolution of**
545 **increased competitive ability**. *Front Ecol Environ* 2004, **2**(8):436-443.
- 546 11. Vilcinskas A, Stoecker K, Schmidtberg H, Rohrich CR, Vogel H: **Invasive harlequin**
547 **ladybird carries biological weapons against native competitors**. *Science* 2013,

- 548 **340**(6134):862-863.
- 549 12. Sax DF, Brown JH: **The paradox of invasion**. *Global Ecol Biogeogr* 2000, **9**(5):363-371.
- 550 13. Liu SS, De Barro PJ, Xu J, Luan JB, Zang LS, Ruan YM, Wan FH: **Asymmetric mating**
- 551 **interactions drive widespread invasion and displacement in a whitefly**. *Science*
- 552 2007, **318**(5857): 1769-1772.
- 553 14. Shavit O, Dafni A, Ne'eman G: **Competition between honeybees (*Apis mellifera*) and**
- 554 **native solitary bees in the Mediterranean region of Israel-Implications for**
- 555 **conservation**. *Israel Journal of Plant Sciences* 2009, **57**(3): 171-183.
- 556 15. Tillberg CV, Holway DA, LeBrun EG, Suarez AV: **Trophic ecology of invasive Argentine**
- 557 **ants in their native and introduced ranges**. *Proceeding of the National Academy of*
- 558 *Sciences of the USA* 2007, **104**(52): 20856-20861.
- 559 16. Alpert P, Bone E, Holzapfel C: **Invasiveness, invasibility and the role of**
- 560 **environmental stress in the spread of non-native plants**. *Perspectives in plant*
- 561 *ecology, evolution and systematics* 2000, **3**(1):52-66.
- 562 17. Whitney KD, Gabler CA: **Rapid evolution in introduced species, 'invasive traits' and**
- 563 **recipient communities: challenges for predicting invasive potential**. *Divers Distrib*
- 564 2008, **14**(4):569-580.
- 565 18. Blackburn TM, Pysek P, Bacher S, Carlton JT, Duncan RP, Jarosik V, Wilson JRU,
- 566 Richardson DM: **A proposed unified framework for biological invasions**. *Trends Ecol*
- 567 *Evol* 2011, **26**(7):333-339.
- 568 19. Capellini I, Baker J, Allen WL, Street SE, Venditti C: **The role of life history traits in**
- 569 **mammalian invasion success**. *Ecol Lett* 2015, **18**(10):1099-1107.
- 570 20. Rius M, Bourne S, Hornsby HG, Chapman MA: **Applications of next-generation**
- 571 **sequencing to the study of biological invasions**. *Curr Zool* 2015, **61**(3):488-504.
- 572 21. Swarts K, Gutaker RM, Benz B, Blake M, Bukowski R, Holland J, Kruse-Peebles M,
- 573 Lepak N, Prim L, Romay MC *et al*: **Genomic estimation of complex traits reveals**
- 574 **ancient maize adaptation to temperate North America**. *Science* 2017, **357**(6350):512-
- 575 515.
- 576 22. Kim S, Cho YS, Kim HM, Chung O, Kim H, Jho S, Seomun H, Kim J, Bang WY, Kim C *et*
- 577 *al*: **Comparison of carnivore, omnivore, and herbivore mammalian genomes with a**
- 578 **new leopard assembly**. *Genome Biol* 2016, **17**.
- 579 23. Lippert C, Sabatini R, Maher MC, Kang EY, Lee S, Arikan O, Harley A, Bernal A, Garst P,
- 580 Lavrenko V *et al*: **Identification of individuals by trait prediction using whole-**
- 581 **genome sequencing data**. *P Natl Acad Sci USA* 2017, **114**(38):10166-10171.
- 582 24. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL:
- 583 **BLAST+: architecture and applications**. *BMC Bioinformatics* 2009, **10**:421.
- 584 25. Li L, Stoeckert CJ, Roos DS: **OrthoMCL: Identification of ortholog groups for**
- 585 **eukaryotic genomes**. *Genome Res* 2003, **13**(9):2178-2189.
- 586 26. Han MV, Thomas GWC, Lugo-Martinez J, Hahn MW: **Estimating gene gain and loss**
- 587 **rates in the presence of error in genome assembly and annotation using CAFE 3**.
- 588 *Mol Biol Evol* 2013, **30**(8):1987-1997.
- 589 27. Kursa MB, Rudnicki WR: **Feature selection with the Boruta package**. *J Stat Softw*

- 2010, **36**(11):1-13.
- 591 28. Ling CX, Huang J, Zhang H: **AUC: a statistically consistent and more discriminating**
592 **measure than accuracy**. In: *IJCAI: 2003*; 2003: 519-524.
- 593 29. Simberloff D, Gibbons L: **Now you see them, now you don't!-population crashes of**
594 **established introduced species**. *Biological Invasions* 2004, **6**(2): 161-172.
- 595 30. Ewel JJ: **Invasibility: lessons from South Florida**. In: *Ecology of Biological Invasions*
596 *of North America and Hawaii*. Edited by Mooney HA and Drake JA. New York: Springer-
597 Verlag; 1986: 214-230.
- 598 31. De Roy K, Marzorati M, Negroni A, Thas O, Balloi A, Fava F, Verstraete W, Daffonchio D,
599 Boon N: **Environmental conditions and community evenness determine the**
600 **outcome of biological invasion**. *Nature Communications* 2013, **4**.
- 601 32. Chown SL, Hodgins KA, Griffin PC, Oakeshott JG, Byrne M, Hoffmann AA: **Biological**
602 **invasions, climate change and genomics**. *Evolutionary Applications* 2015, **8**(1):23-46.
- 603 33. Worner SP, Gevrey M: **Modelling global insect pest species assemblages to**
604 **determine risk of invasion**. *J Appl Ecol* 2006, **43**(5):858-867.
- 605 34. Tilman D: **Causes, consequences and ethics of biodiversity**. *Nature* 2000,
606 **405**(6783):208-211.
- 607 35. Jessup CM, Bohannan BJM: **The shape of an ecological trade-off varies with**
608 **environment**. *Ecol Lett* 2008, **11**(9):947-959.
- 609 36. Jing DP, Guo JF, Jiang YY, Zhao JZ, Sethi A, He KL, Wang ZY: **Initial detections and**
610 **spread of invasive *Spodoptera frugiperda* in China and comparisons with other**
611 **noctuid larvae in cornfield using molecular techniques**. *Insect Science* 2019, **27**(4):
612 780-790.
- 613 37. Xie W, Yang X, Chen CH, Yang ZZ, Guo LT, Wang D, Huang JQ, Zhang HL, Wen YN,
614 Zhao JY *et al*: **The invasive MED/Q *Bemisia tabaci* genome: a tale of gene loss and**
615 **gene gain**. *Bmc Genomics* 2018, **19**.
- 616 38. Zenni RD, Nunez MA: **The elephant in the room: the role of failed invasions in**
617 **understanding invasion biology**. *Oikos* 2013, **122**(6):801-815.
- 618 39. Agarwala R, Barrett T, Beck J, Benson DA, Bollin C, Bolton E, Bourexis D, Brister JR,
619 Bryant SH, Canese K *et al*: **Database resources of the National Center for**
620 **Biotechnology Information**. *Nucleic Acids Res* 2018, **46**(D1):D8-D13.
- 621 40. Yin C, Shen G, Guo D, Wang S, Ma X, Xiao H, Liu J, Zhang Z, Liu Y, Zhang Y *et al*:
622 **InsectBase: a resource for insect genomes and transcriptomes**. *Nucleic Acids Res*
623 2016, **44**(D1):D801-807.
- 624 41. Giraldo-Calderon GI, Emrich SJ, MacCallum RM, Maslen G, Dialynas E, Topalis P, Ho N,
625 Gesing S, VectorBase C, Madey G *et al*: **VectorBase: an updated bioinformatics**
626 **resource for invertebrate vectors and other organisms related with human**
627 **diseases**. *Nucleic Acids Res* 2015, **43**(Database issue):D707-713.
- 628 42. Kersey PJ, Allen JE, Allot A, Barba M, Boddu S, Bolt BJ, Carvalho-Silva D, Christensen
629 M, Davis P, Grabmueller C *et al*: **Ensembl Genomes 2018: an integrated omics**
630 **infrastructure for non-vertebrate species**. *Nucleic Acids Res* 2018, **46**(D1):D802-
631 D808.

- 632 43. Sneddon TP, Li P, Edmunds SC: **GigaDB: announcing the GigaScience database.**
633 *Gigascience* 2012, **1**:11.
- 634 44. Wurm Y, Uva P, Ricci F, Wang J, Jemielity S, Iseli C, Falquet L, Keller L: **Fourmidable:**
635 **a database for ant genomics.** *Bmc Genomics* 2009, **10**:5.
- 636 45. Zhan S, Reppert SM: **MonarchBase: the monarch butterfly genome database.**
637 *Nucleic Acids Res* 2013, **41**(D1):D758-D763.
- 638 46. Legeai F, Shigenobu S, Gauthier JP, Colbourne J, Rispe C, Collin O, Richards S, Wilson
639 ACC, Murphy T, Tagu D: **AphidBase: a centralized bioinformatic resource for**
640 **annotation of the pea aphid genome.** *Insect Mol Biol* 2010, **19**:5-12.
- 641 47. Parks DH, Rinke C, Chuvochina M, Chaumeil PA, Woodcroft BJ, Evans PN, Hugenholtz
642 P, Tyson GW: **Recovery of nearly 8,000 metagenome-assembled genomes**
643 **substantially expands the tree of life.** *Nat Microbiol* 2017, **2**(11):1533-1542.
- 644 48. Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, Li RQ, Liu T, Zhang Z,
645 Bolund L *et al*: **TreeFam: a curated database of phylogenetic trees of animal gene**
646 **families.** *Nucleic Acids Res* 2006, **34**:D572-D580.
- 647 49. Katoh K, Standley DM: **MAFFT Multiple sequence alignment software version 7:**
648 **improvements in performance and usability.** *Mol Biol Evol* 2013, **30**(4):772-780.
- 649 50. Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T: **trimAl: a tool for automated**
650 **alignment trimming in large-scale phylogenetic analyses.** *Bioinformatics* 2009,
651 **25**(15):1972-1973.
- 652 51. Misof B, Liu SL, Meusemann K, Peters RS, Donath A, Mayer C, Frandsen PB, Ware J,
653 Flouri T, Beutel RG *et al*: **Phylogenomics resolves the timing and pattern of insect**
654 **evolution.** *Science* 2014, **346**(6210):763-767.
- 655 52. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis**
656 **of large phylogenies.** *Bioinformatics* 2014, **30**(9):1312-1313.
- 657 53. Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ: **IQ-TREE: a fast and effective**
658 **stochastic algorithm for estimating maximum-likelihood phylogenies.** *Mol Biol Evol*
659 2015, **32**(1):268-274.
- 660 54. Sanderson MJ: **A nonparametric approach to estimating divergence times in the**
661 **absence of rate constancy.** *Mol Biol Evol* 1997, **14**(12):1218-1231.
- 662 55. Sanderson MJ: **Estimating absolute rates of molecular evolution and divergence**
663 **times: a penalized likelihood approach.** *Mol Biol Evol* 2002, **19**(1):101-109.
- 664 56. Sanderson MJ: **r8s: inferring absolute rates of molecular evolution and divergence**
665 **times in the absence of a molecular clock.** *Bioinformatics* 2003, **19**(2):301-302.
- 666 57. McKenna DD, Scully ED, Pauchet Y, Hoover K, Kirsch R, Geib SM, Mitchell RF,
667 Waterhouse RM, Ahn SJ, Arsala D *et al*: **Genome of the Asian longhorned beetle**
668 **(*Anoplophora glabripennis*), a globally significant invasive species, reveals key**
669 **functional and evolutionary innovations at the beetle-plant interface.** *Genome Biol*
670 2016, **17**(1):227.
- 671 58. Smit AF, Hubble R, Green P: **2010 RepeatMasker Open-3.0.** URL: [http://www](http://www.repeatmasker.org)
672 [repeatmasker.org](http://www.repeatmasker.org) 1996.
- 673 59. Bao W, Kojima KK, Kohany O: **Repbase Update, a database of repetitive elements in**

674 **eukaryotic genomes.** *Mobile DNA* 2015, **6**:11.
675 60. UniProt Consortium T: **UniProt: the universal protein knowledgebase.** *Nucleic Acids*
676 *Res* 2018, **46**(5):2699.
677 61. Maloof MA: **Learning when data sets are imbalanced and when costs are unequal**
678 **and unknown.** In: *ICML-2003 workshop on learning from imbalanced data sets II: 2003*;
679 2003: 2-1.
680



682

683

684

685

686

687

688

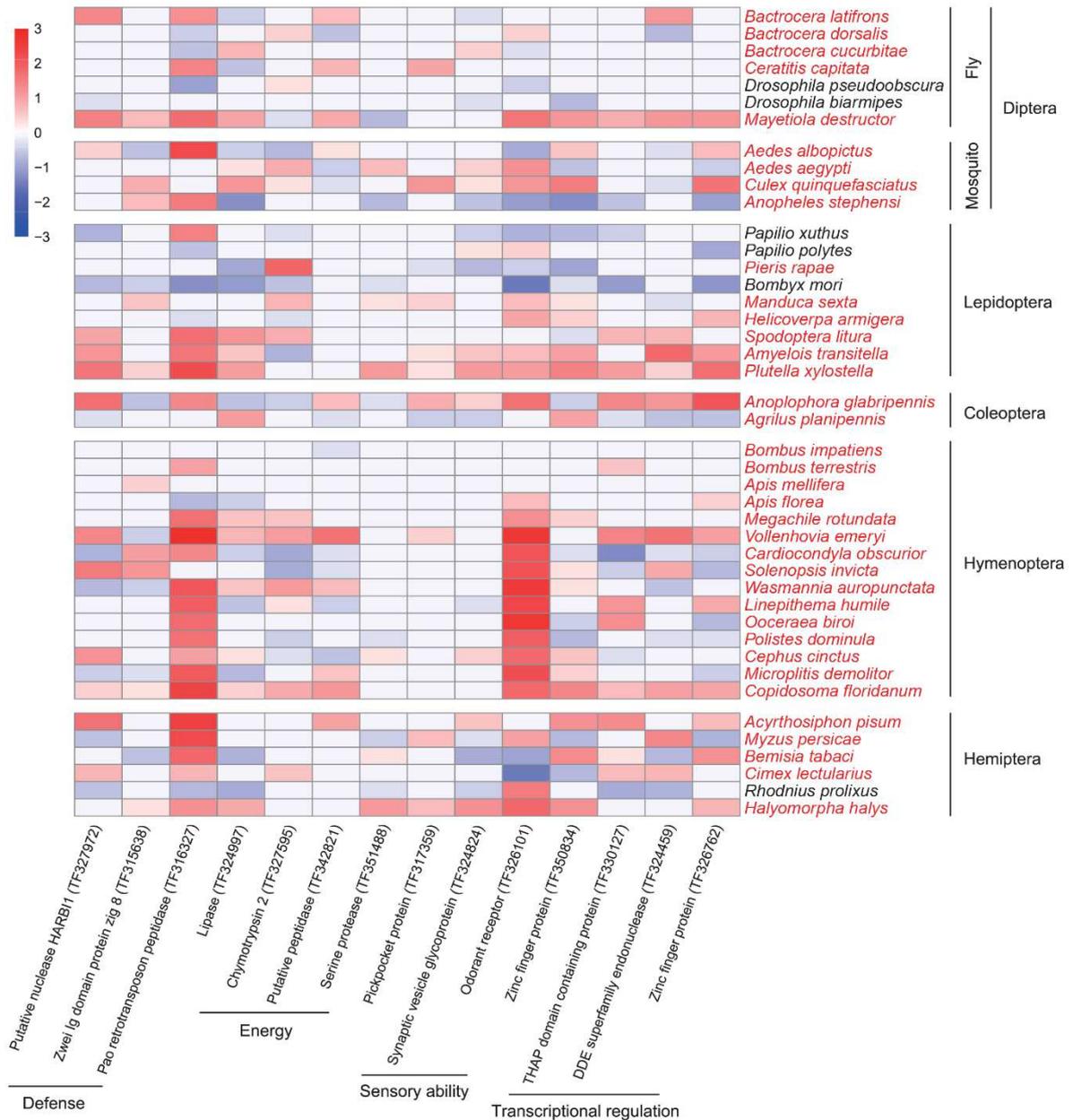
689

690

691

692

Figure 1. Phylogenetic tree and comparison of general genome features. The phylogenetic tree shows the topology and divergence time for 44 arthropods. The mite *Tetranychus urticae* was used as an outgroup. Numbers at branches and tips indicate the number of gene families that are expanded (blue color) or contracted (green color) as compared to the closest tip. MRCA = most recent common ancestor. The number in parentheses is the number of gene families in the MRCA as estimated with TreeFam software. Differences in genome size, GC content, gene number, and amount of repeat sequences between all invasive (red lettering) and non-invasive species (black lettering), as well as between the corresponding fly species and Lepidoptera each analyzed separately, are not significant by t-test.



693

694

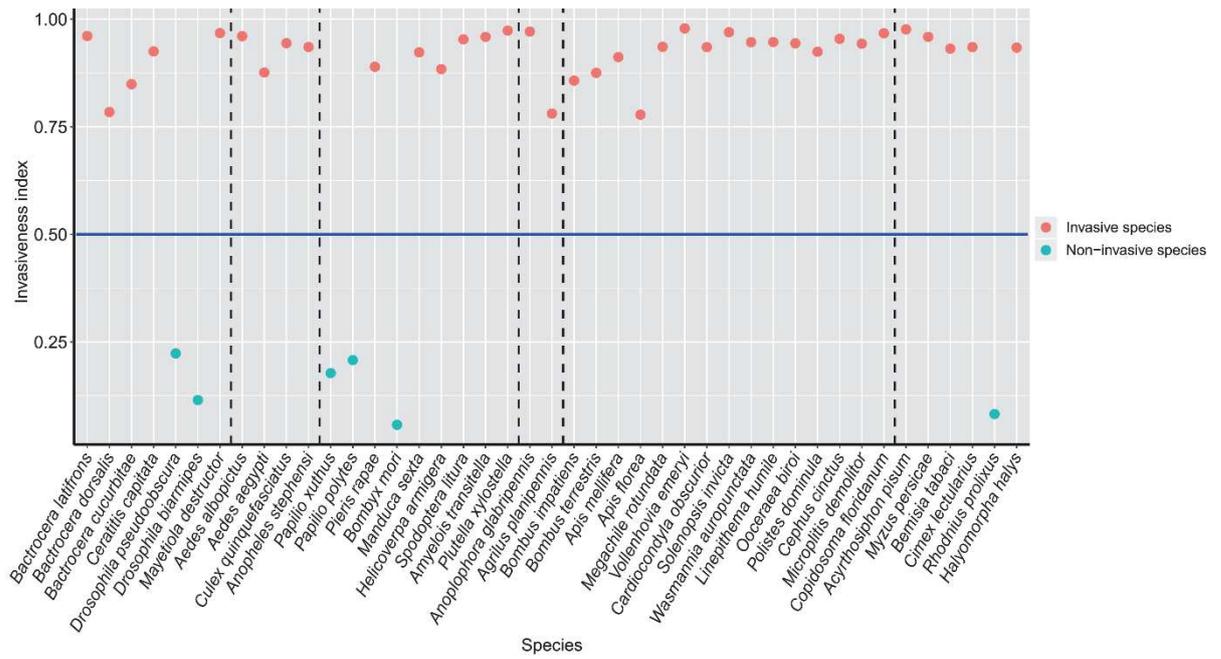
695

696

697

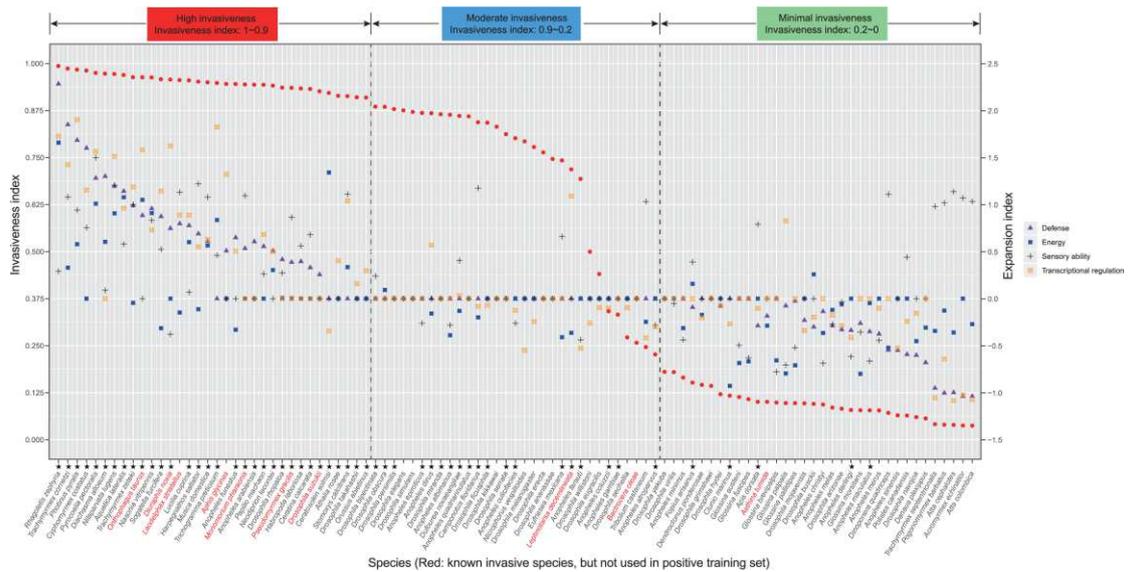
698

Figure 2. The comparison of expanded and contracted gene number in invasiveness-related gene families between invasive species (red lettering) and non-invasive species (black lettering). The expansion and contraction gene numbers were converted by \log_{10} . $y = \log_{10}(|x|)$ ($x \geq 1$) or $y = -\log_{10}(|x|)$ ($x \leq -1$), where x represents the expanded gene number ($x \geq 1$) or the contracted gene number ($x \leq -1$) and y was used in the heatmap.



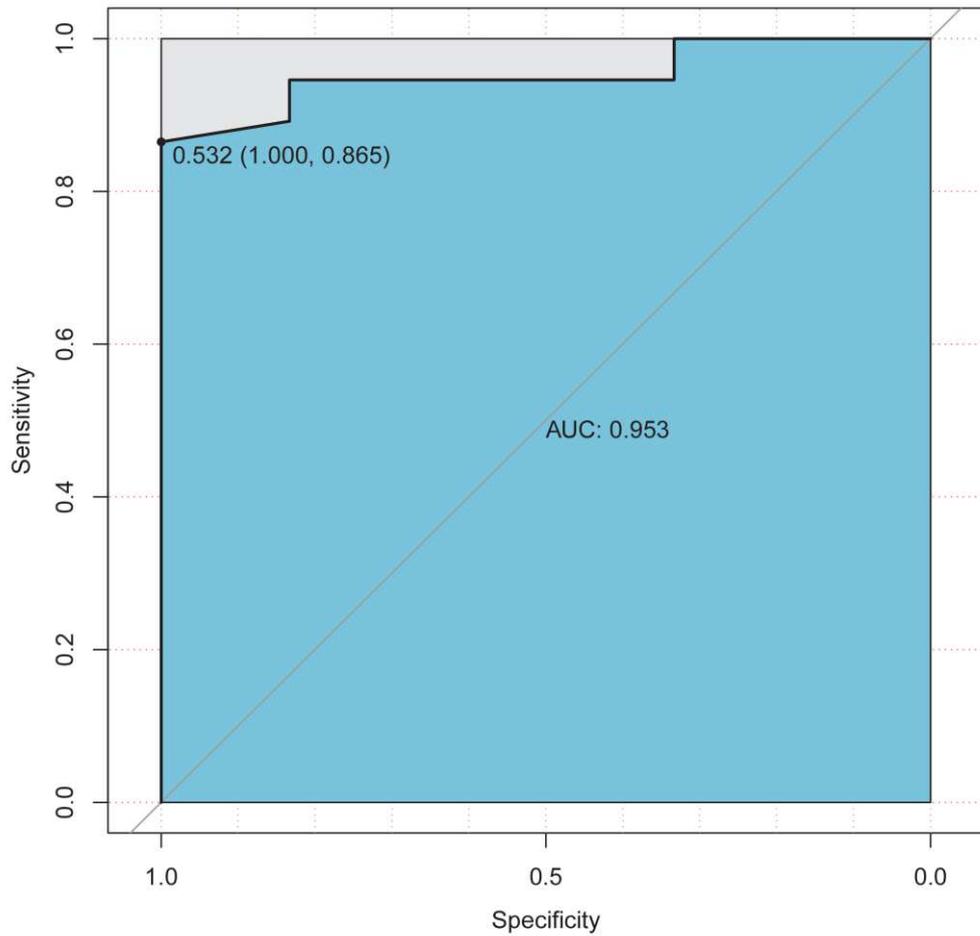
699
700
701
702
703

Figure 3. Invasiveness indexes of all forty-three species. Red dots represent invasive insects and green dots represent non-invasive insects. The dashed lines separate the species into different taxa: fly, mosquito, Lepidoptera, Coleoptera, Hymenoptera, and Hemiptera from left to right. The blue solid line represents the cutoff value of 0.5.



704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720

Figure 4. Invasiveness indexes and gene family expansion indexes of the other 99 insect species. The symbol “●” represents the invasiveness index. Three levels of invasiveness (high, moderate, and low) were classified by the invasiveness index cutoffs at 0.9 and 0.2. Fourteen identified invasiveness-related gene families are categorized into four function groups as defense, energy, chemosensory function, and transcriptional regulation. The symbol “▲” represents the expansion index of gene families in defense function group, the symbol “■” represents the expansion index of gene families in the energy function group, the symbol “+” represents the expansion index of gene families in the chemosensory function group, and the symbol “⊠” represents the expansion index of gene families in the transcriptional regulation function group. The species in red lettering were confirmed to be invasive but excluded in the 43-species sample set because of their relatively low-quality genome assemblies (a scaffold N50 < 400 Kb), while the ones in black were species with no evidence to confirm them as either invasive nor non-invasive (generally because they have not been confirmed to have been introduced anywhere). The symbol “★” above a species’ scientific name means this species was predicted to be invasive by DIGS (Determine Invasiveness based on Genome Sequences).



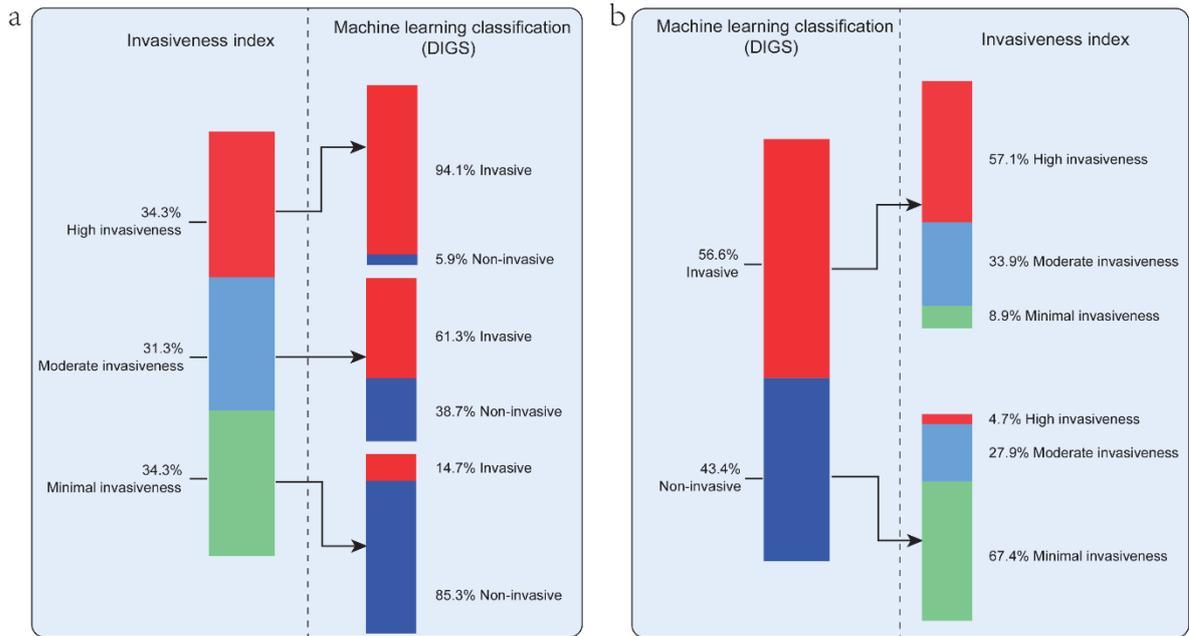
721

722

723

724

Figure 5. The ROC curve of the DIGS. ROC curve: the receiver operating characteristic curve, AUC: the area under the curve of ROC, sensitivity: true positive rate, specificity: true negative rate, DIGS: Determine Invasiveness based on Genome Sequences.



725
 726
 727
 728
 729
 730
 731
 732
 733

Figure 6. Invasiveness classification by DIGS and invasiveness level assessment by invasiveness index. a) The percentage of invasive and non-invasive species classified by DIGS in each invasiveness level assessed by the invasiveness index was calculated. b) The proportions of species with different levels of invasiveness in invasive and non-invasive categories classified by DIGS is shown.

Figures

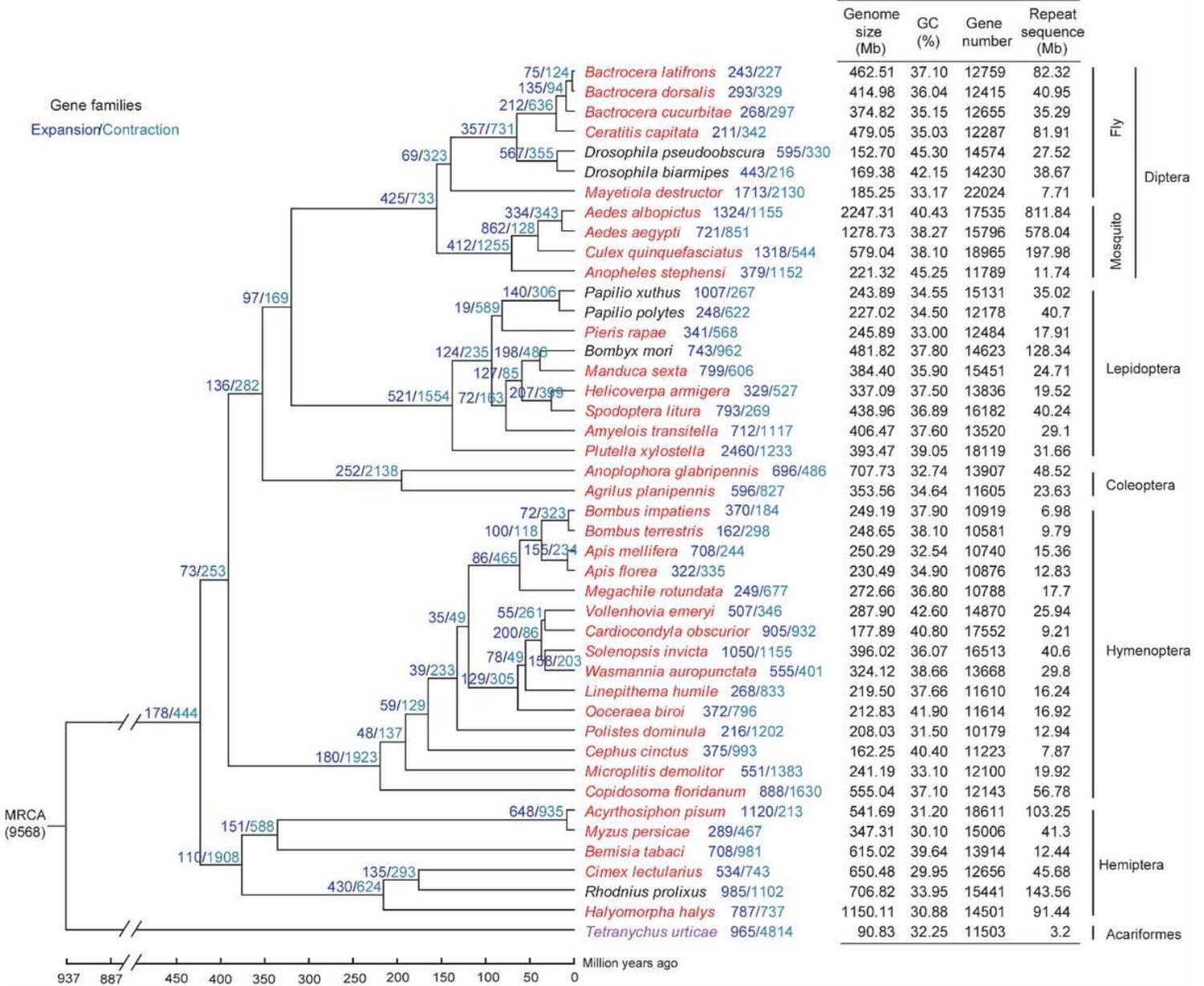


Figure 1

Phylogenetic tree and comparison of general genome features. The phylogenetic tree shows the topology and divergence time for 44 arthropods. The mite *Tetranychus urticae* was used as an outgroup. Numbers at branches and tips indicate the number of gene families that are expanded (blue color) or contracted (green color) as compared to the closest tip. MRCA = most recent common ancestor. The number in parentheses is the number of gene families in the MRCA as estimated with TreeFam software. Differences in genome size, GC content, gene number, and amount of repeat sequences between all invasive (red lettering) and non-invasive species (black lettering), as well as between the corresponding fly species and Lepidoptera each analyzed separately, are not significant by t-test.

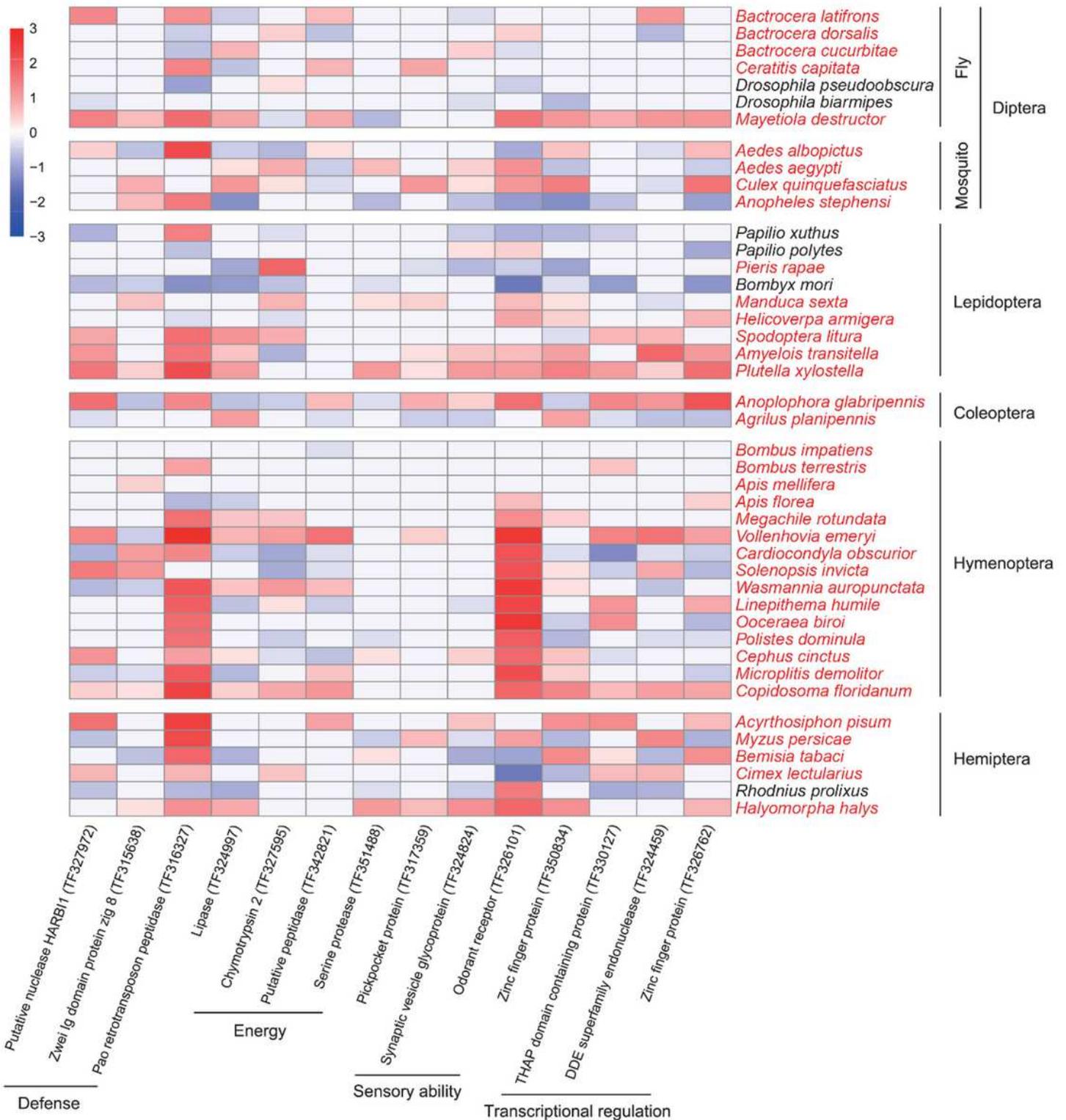


Figure 2

The comparison of expanded and contracted gene number in invasiveness-related gene families between invasive species (red lettering) and non-invasive species (black lettering). The expansion and contraction gene numbers were converted by \log_{10} . $y = \log_{10}(|x|)$ ($x \geq 1$) or $y = -\log_{10}(|x|)$ ($x \leq -1$), where x represents the expanded gene number ($x \geq 1$) or the contracted gene number ($x \leq -1$) and y was used in the heatmap.

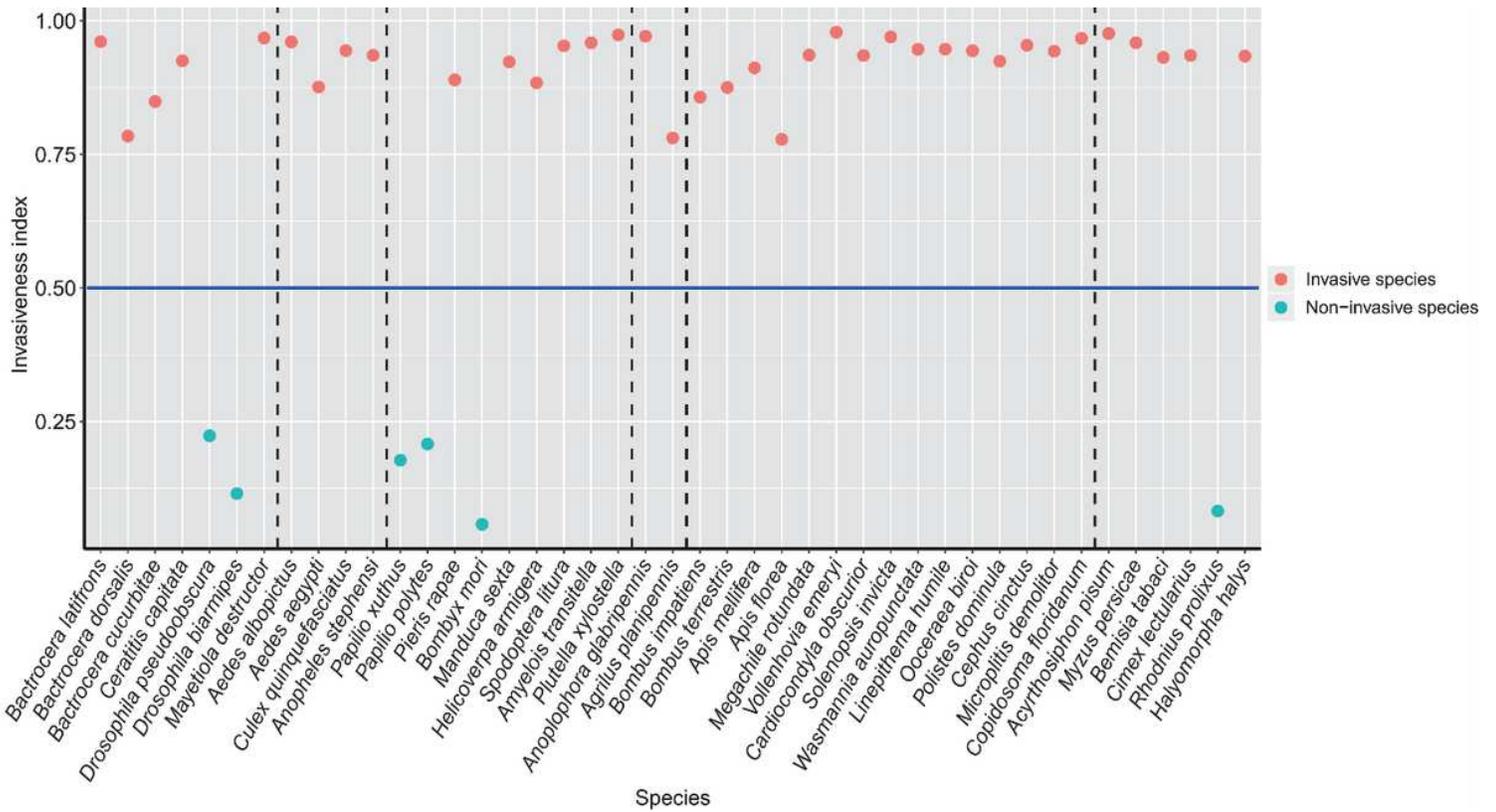
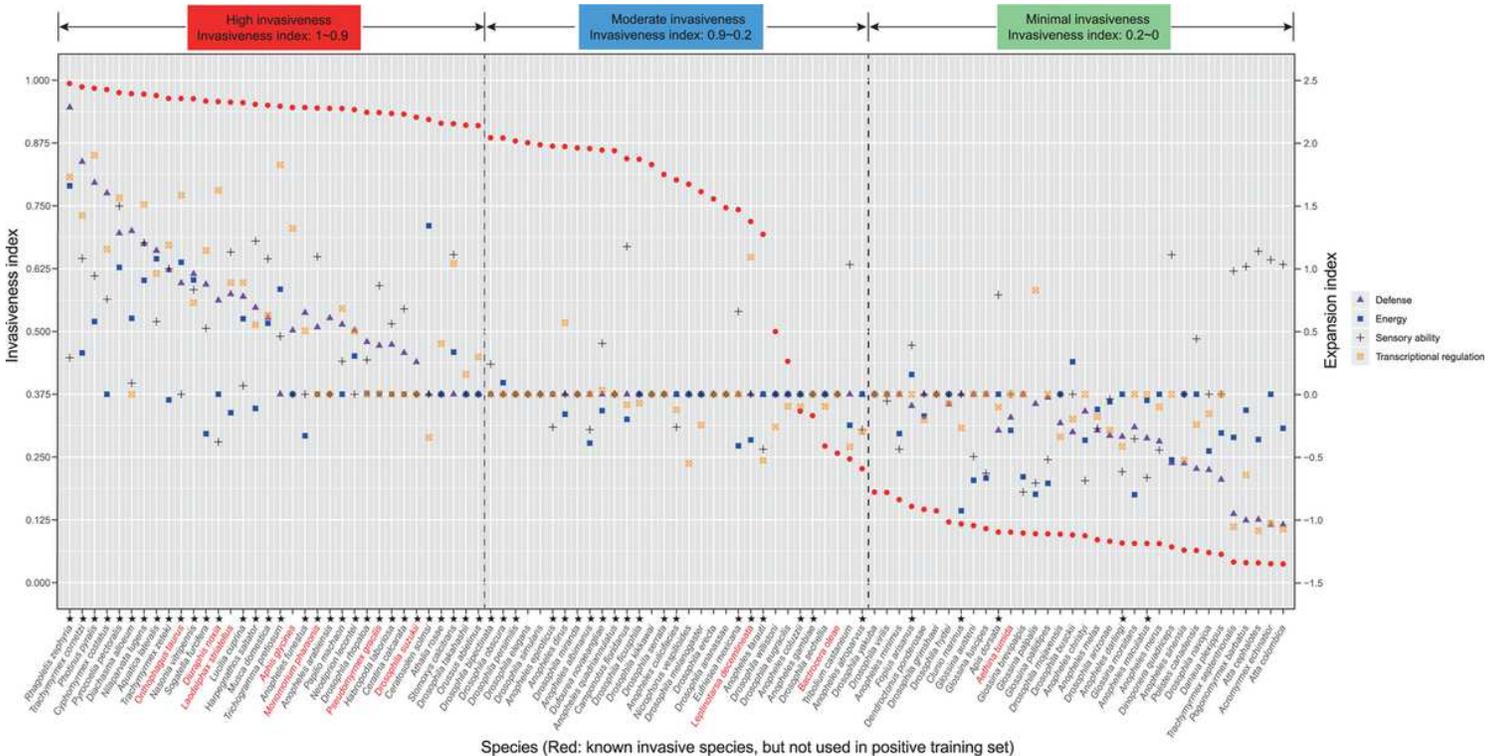


Figure 3

Invasiveness indexes of all forty-three species. Red dots represent invasive insects and green dots represent non-invasive insects. The dashed lines separate the species into different taxa: fly, mosquito, Lepidoptera, Coleoptera, Hymenoptera, and Hemiptera from left to right. The blue solid line represents the cutoff value of 0.5.



Species (Red: known invasive species, but not used in positive training set)

Figure 4

. Invasiveness indexes and gene family expansion indexes of the other 99 insect species. The symbol “ \boxtimes ” represents the invasiveness index. Three levels of invasiveness (high, moderate, and low) were classified by the invasiveness index cutoffs at 0.9 and 0.2. Fourteen identified invasiveness-related gene families are categorized into four function groups as defense, energy, chemosensory function, and transcriptional regulation. The symbol “ \boxtimes ” represents the expansion index of gene families in defense function group, the symbol “ \boxplus ” represents the expansion index of gene families in the energy function group, the symbol “+” represents the expansion index of gene families in the chemosensory function group, and the symbol “ \boxminus ” represents the expansion index of gene families in the transcriptional regulation function group. The species in red lettering were confirmed to be invasive but excluded in the 43-species sample set because of their relatively low-quality genome assemblies (a scaffold N50 < 400 Kb), while the ones in black were species with no evidence to confirm them as either invasive nor non-invasive (generally because they have not been confirmed to have been introduced anywhere). The symbol “ \boxtimes ” above a species’ scientific name means this species was predicted to be invasive by DIGS (Determine Invasiveness based on Genome Sequences).

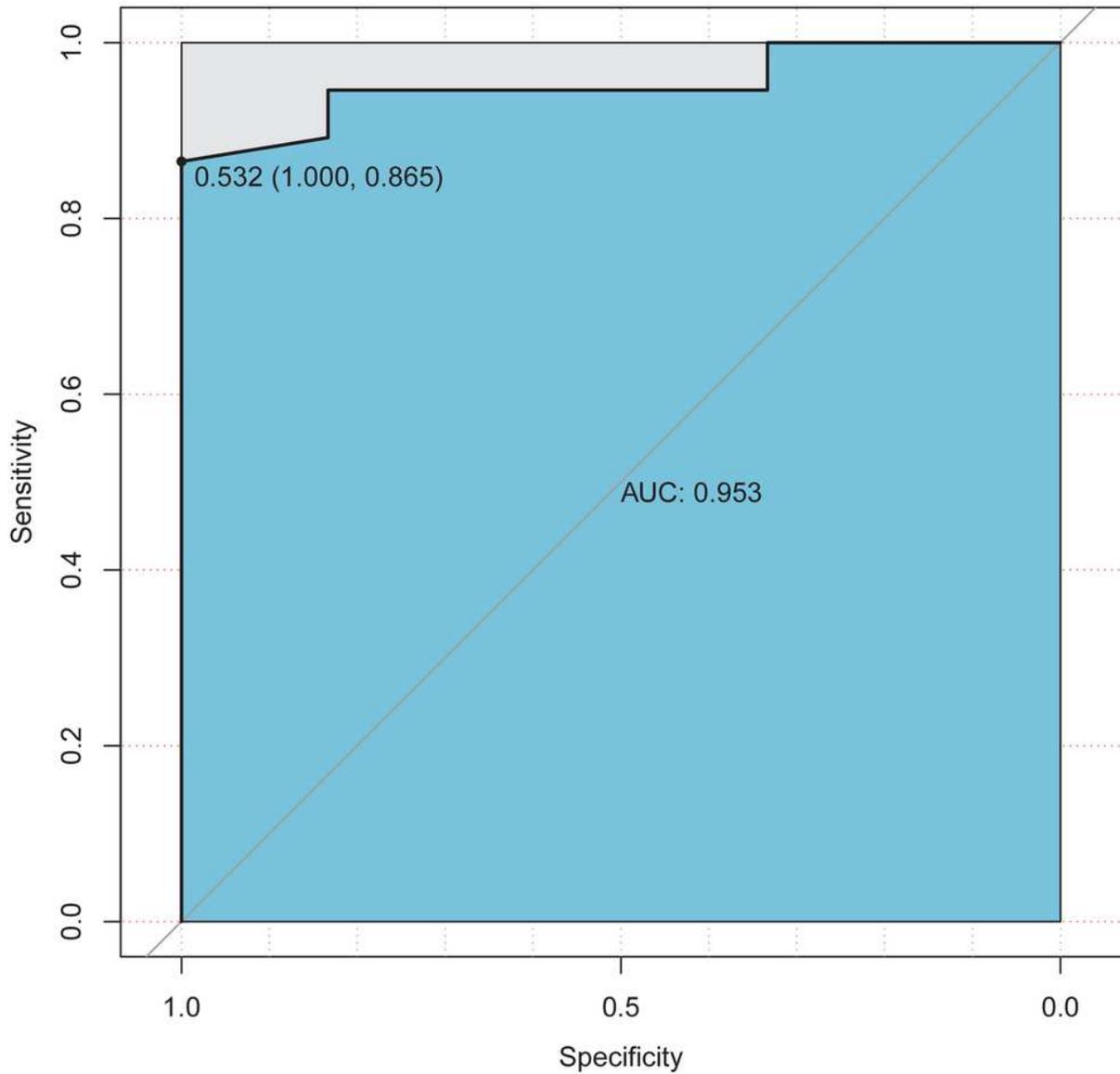


Figure 5

The ROC curve of the DIGS. ROC curve: the receiver operating characteristic curve, AUC: the area under the curve of ROC, sensitivity: true positive rate, specificity: true negative rate, DIGS: Determine Invasiveness based on Genome Sequences.

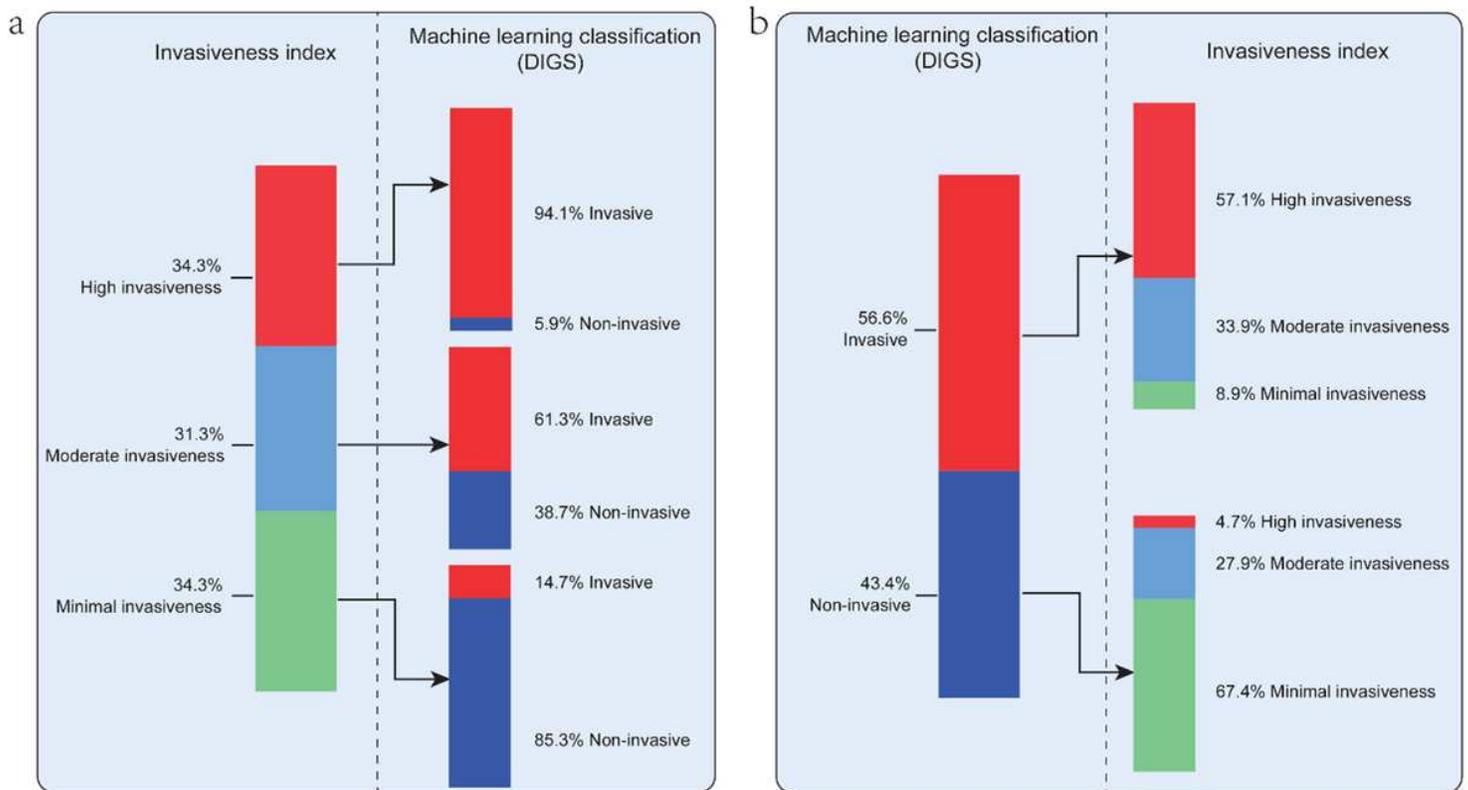


Figure 6

Invasiveness classification by DIGS and invasiveness level assessment by invasiveness index. a) The percentage of invasive and non-invasive species classified by DIGS in each invasiveness level assessed by the invasiveness index was calculated. b) The proportions of species with different levels of invasiveness in invasive and non-invasive categories classified by DIGS is shown.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [BMCGPredictingInsectInvasivenessSI20200707.docx](#)